

Question 1

Classification algorithm	weka classifier	Correct classified	Incorrectly classified
Decision Tree	j48	77.03%	22.96%
Neural Network	MLP (hidden layer = 3)	73.05%	26.94%
Neural Network	MLP (hidden layer = 5)	73.62%	26.37%
k-NN	k=3	80.83%	19.16%

The best classification algorithm was **k-Nearest Neighbours** with k= 3 having 80.83% accuracy as highlighted above.

Decision tree

Classifier

Choose **J48 -C 0.25 -M 2**

Test options

☐ Use training set

☐ Supplied test set

☒ Cross-validation Folds

☐ Percentage split %

(Nom) Diabetes category

Result list (right-click for options)

18:23:54 - trees.J48

Classifier output

```

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

Correctly Classified Instances      406           77.0398 %
Incorrectly Classified Instances    121           22.9602 %
Kappa statistic                    0.586
Mean absolute error                0.1572
Root mean squared error            0.3752
Relative absolute error            42.7033 %
Root relative squared error        87.501 %
Total Number of Instances         527

=== Detailed Accuracy By Class ===

```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.806	0.285	0.814	0.806	0.810	0.520	0.773	0.780	mid
	0.679	0.077	0.626	0.679	0.651	0.583	0.822	0.582	high
	0.740	0.069	0.765	0.740	0.752	0.678	0.860	0.706	low
Weighted Avg.	0.770	0.201	0.773	0.770	0.771	0.567	0.801	0.731	

```

=== Confusion Matrix ===
 a  b  c  <-- classified as
258 34 28 | a = mid
 27 57  0 | b = high
 32  0 91 | c = low

```

NN with 3 hidden layers

Choose **MultilayerPerceptron -L 0.3 -M 0.2 -N 500 -V 0 -S 0 -E 20 -H 3**

Test options

☐ Use training set

☐ Supplied test set

☒ Cross-validation Folds

☐ Percentage split %

(Nom) Diabetes category

Result list (right-click for options)

18:23:54 - trees.J48

18:48:51 - functions.MultilayerPerceptron

Classifier output

```

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

Correctly Classified Instances      385           73.055 %
Incorrectly Classified Instances    142           26.945 %
Kappa statistic                    0.4507
Mean absolute error                0.2498
Root mean squared error            0.3753
Relative absolute error            67.8686 %
Root relative squared error        87.5295 %
Total Number of Instances         527

=== Detailed Accuracy By Class ===

```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.894	0.522	0.726	0.894	0.801	0.418	0.726	0.767	mid
	0.179	0.020	0.625	0.179	0.278	0.278	0.772	0.441	high
	0.683	0.062	0.771	0.683	0.724	0.649	0.855	0.730	low
Weighted Avg.	0.731	0.334	0.720	0.731	0.700	0.450	0.764	0.706	

```

=== Confusion Matrix ===
 a  b  c  <-- classified as
286  9 25 | a = mid
 69 15  0 | b = high
 39  0 84 | c = low

```

NN with 5 hidden layers

Classifier

Choose **MultilayerPerceptron -L 0.3-M 0.2-N 500-V 0-S 0-E 20-H 5**

Test options

☐ Use training set
☐ Supplied test set
☒ Cross-validation Folds
☐ Percentage split %

(Nom) Diabetes category ▼

Result list (right-click for options)

- 18:23:54 - trees.J48
- 18:48:51 - functions.MultilayerPerceptron
- 18:52:11 - functions.MultilayerPerceptron

Classifier output

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

Correctly Classified Instances      388             73.6243 %
Incorrectly Classified Instances    139             26.3757 %
Kappa statistic                    0.4852
Mean absolute error                 0.2434
Root mean squared error             0.37
Relative absolute error             66.129 %
Root relative squared error        86.3064 %
Total Number of Instances         527

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall  F-Measure  MCC   ROC Area  PRC Area  Class
          0.856   0.449   0.747    0.856   0.798    0.432  0.729    0.765    mid
          0.286   0.043   0.558   0.286   0.378    0.325  0.792    0.454    high
          0.732   0.067   0.769   0.732   0.750    0.677  0.862    0.691    low
Weighted Avg.   0.736   0.295   0.722   0.736   0.720    0.472  0.770    0.698

=== Confusion Matrix ===

  a   b   c   <-- Classified as
274  19  27 | a = mid
 60  24   0 | b = high
 33   0  90 | c = low
```

k-NN with k=3

Choose **IBk -K 3-W 0-A "weka.core.neighboursearch.LinearNNSearch-A "weka.core.EuclideanDistance -R first-last"**

Test options

☐ Use training set
☐ Supplied test set
☒ Cross-validation Folds
☐ Percentage split %

(Nom) Diabetes category ▼

Result list (right-click for options)

- 18:23:54 - trees.J48
- 18:48:51 - functions.MultilayerPerceptron
- 18:52:11 - functions.MultilayerPerceptron
- 19:14:21 - lazy.IBk

Classifier output

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

Correctly Classified Instances      426             80.8349 %
Incorrectly Classified Instances    101             19.1651 %
Kappa statistic                    0.6475
Mean absolute error                 0.1533
Root mean squared error             0.3195
Relative absolute error             41.635 %
Root relative squared error        74.5194 %
Total Number of Instances         527

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall  F-Measure  MCC   ROC Area  PRC Area  Class
          0.859   0.271   0.831   0.859   0.845   0.595  0.844    0.849    mid
          0.750   0.050   0.741   0.750   0.746   0.697  0.903    0.722    high
          0.715   0.057   0.793   0.715   0.752   0.683  0.890    0.753    low
Weighted Avg.   0.808   0.165   0.808   0.808   0.807   0.632  0.865    0.806

=== Confusion Matrix ===

  a   b   c   <-- Classified as
275  22  23 | a = mid
 21  63   0 | b = high
 35   0  88 | c = low
```

Question 2

Classification algorithm	Correct classified	Incorrectly classified
Average Probabilities	81.59%	18.40%
Majority Voting	81.02%	18.97%
Minimum Probability	80.26%	19.73%
Maximum Probability	80.64%	19.35%

I used below 3 classifiers for the above 4 different combination rules

Classification algorithm	weka classifier
Decision Tree	j48
Neural Network	MLP (hidden layer = 3)
k-NN	k=3

The small differences in accuracy when using above 4 combination rules is because

- of higher diversity of the dataset
- less disagreement between the decisions within the voters leading to small differences in accuracies in all 4 combinations
- possible weighting and bias impact on the output of the neuron impacting activation functions
- depending on the problem, complex decisions may require long chains of computational stages based on voting pattern

Average Probabilities

Choose: Vote -S 1 -B "weka.classifiers.lazy.IBk -K 3 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A "weka.core.EuclideanDistance -R first-last" -B "weka.classifiers.functions.MultilayerPerceptron -L 0.3 -M 0.2 -N 500 -V 0 -S 0 -

Test options

☐ Use training set
☐ Supplied test set
☒ Cross-validation Folds: 10
☐ Percentage split %: 66

More options...

(Nom) Diabetes category

Start Stop

Result list (right-click for options)

- 18:23:54 - trees.J48
- 18:48:51 - functions.MultilayerPerceptron
- 18:52:11 - functions.MultilayerPerceptron
- 19:14:21 - lazy.IBk
- 19:39:42 - meta.Vote

Classifier output

```

=== Stratified cross-validation ===
=== Summary ===
Correctly Classified Instances      430      81.5939 %
Incorrectly Classified Instances    97      18.4061 %
Kappa statistic                    0.6559
Mean absolute error                 0.1868
Root mean squared error             0.3079
Relative absolute error             50.7356 %
Root relative squared error         71.808 %
Total Number of Instances          527

=== Detailed Accuracy By Class ===
               TP Rate  FP Rate  Precision  Recall   F-Measure  MCC      ROC Area  PRC Area  Class
               ----
               0.881    0.285    0.827    0.861    0.853    0.609    0.850    0.653    mid
               0.655    0.038    0.764    0.655    0.705    0.657    0.927    0.768    high
               0.756    0.052    0.816    0.756    0.785    0.723    0.916    0.836    low
Weighted Avg.   0.816    0.191    0.814    0.816    0.814    0.644    0.878    0.835

=== Confusion Matrix ===
      a   b   c   <-- classified as
282  17   21 |   a = mid
 29   55   0 |   b = high
 30   0   93 |   c = low

```

Majority Voting

Choose: Vote -S 1 -B "weka.classifiers.lazy.IBk -K 3 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A "weka.core.EuclideanDistance -R first-last" -B "weka.classifiers.functions.MultilayerPerceptron -L 0.3 -M 0.2 -N 500 -V 0 -S 0 -

Test options

☐ Use training set
☐ Supplied test set
☒ Cross-validation Folds: 10
☐ Percentage split %: 66

More options...

(Nom) Diabetes category

Start Stop

Result list (right-click for options)

- 18:23:54 - trees.J48
- 18:48:51 - functions.MultilayerPerceptron
- 18:52:11 - functions.MultilayerPerceptron
- 19:14:21 - lazy.IBk
- 19:39:42 - meta.Vote
- 19:46:49 - meta.Vote

Classifier output

```

=== Stratified cross-validation ===
=== Summary ===
Correctly Classified Instances      427      81.0247 %
Incorrectly Classified Instances    100      18.9753 %
Kappa statistic                    0.6426
Mean absolute error                 0.1265
Root mean squared error             0.3557
Relative absolute error             34.3659 %
Root relative squared error         82.957 %
Total Number of Instances          527

=== Detailed Accuracy By Class ===
               TP Rate  FP Rate  Precision  Recall   F-Measure  MCC      ROC Area  PRC Area  Class
               ----
               0.884    0.304    0.818    0.884    0.850    0.597    0.790    0.794    mid
               0.631    0.036    0.768    0.631    0.693    0.645    0.797    0.543    high
               0.740    0.052    0.813    0.740    0.774    0.711    0.844    0.662    low
Weighted Avg.   0.810    0.203    0.809    0.810    0.807    0.631    0.804    0.723

=== Confusion Matrix ===
      a   b   c   <-- classified as
283  16   21 |   a = mid
 31   53   0 |   b = high
 32   0   91 |   c = low

```

Minimum Probability

Classifier

Choose **Vote** -S 1 -B "weka.classifiers.lazy.IBK-K 3 -W 0 -A "weka.core.neighboursearch.LinearNNSearch-A "weka.core.EuclideanDistance-R first-last" -B "weka.classifiers.functions.MultilayerPerceptron-L 0.3-M 0.2-N 500-V 0-S

Test options

☐ Use training set
☐ Supplied test set Set...
☒ Cross-validation Folds **10**
☐ Percentage split % 66
More options...

(Nom) Diabetes category

Start Stop

Result list (right-click for options)

- 18:23:54 - trees.J48
- 18:48:51 - functions.MultilayerPerceptron
- 18:52:11 - functions.MultilayerPerceptron
- 19:14:21 - lazy.IBK
- 19:39:42 - meta.Vote
- 19:46:49 - meta.Vote
- 19:51:19 - meta.Vote**

Classifier output

```
=== stratified cross-validation ===
=== Summary ===
Correctly Classified Instances      423          80.2657 %
Incorrectly Classified Instances    104          19.7343 %
Kappa statistic                    0.6404
Mean absolute error                 0.1329
Root mean squared error             0.3509
Relative absolute error             36.0967 %
Root relative squared error         61.8451 %
Total Number of Instances          527

=== Detailed Accuracy By Class ===
          TP Rate  FP Rate  Precision  Recall   F-Measure  MCC      ROC Area  PRC Area  Class
          0.844    0.261    0.833     0.844    0.839     0.595    0.783    0.784    mid
          0.702    0.061    0.696     0.702    0.694     0.635    0.837    0.633    high
          0.764    0.057    0.803     0.764    0.783     0.720    0.858    0.707    low
Weighted Avg.   0.803    0.181    0.803     0.803    0.803     0.624    0.809    0.742

=== Confusion Matrix ===
      a  b  c  <-- classified as
270  27  23  |  a = mid
 25  59   0  |  b = high
 29   0  94  |  c = low
```

Maximum Probability

Classifier

Choose **Vote** -S 1 -B "weka.classifiers.lazy.IBK-K 3 -W 0 -A "weka.core.neighboursearch.LinearNNSearch-A "weka.core.EuclideanDistance-R first-last" -B "weka.classifiers.functions.MultilayerPerceptron-L 0.3-M 0.2-N 500-V 0-S

Test options

☐ Use training set
☐ Supplied test set Set...
☒ Cross-validation Folds **10**
☐ Percentage split % 66
More options...

(Nom) Diabetes category

Start Stop

Result list (right-click for options)

- 18:23:54 - trees.J48
- 18:48:51 - functions.MultilayerPerceptron
- 18:52:11 - functions.MultilayerPerceptron
- 19:14:21 - lazy.IBK
- 19:39:42 - meta.Vote
- 19:46:49 - meta.Vote
- 19:51:19 - meta.Vote
- 19:55:52 - meta.Vote**

Classifier output

```
=== stratified cross-validation ===
=== Summary ===
Correctly Classified Instances      425          80.6452 %
Incorrectly Classified Instances    102          19.3548 %
Kappa statistic                    0.6461
Mean absolute error                 0.2324
Root mean squared error             0.323
Relative absolute error             63.1466 %
Root relative squared error         75.3366 %
Total Number of Instances          527

=== Detailed Accuracy By Class ===
          TP Rate  FP Rate  Precision  Recall   F-Measure  MCC      ROC Area  PRC Area  Class
          0.850    0.261    0.834     0.850    0.842     0.592    0.843    0.851    mid
          0.702    0.052    0.720     0.702    0.711     0.657    0.928    0.751    high
          0.764    0.062    0.790     0.764    0.777     0.711    0.923    0.843    low
Weighted Avg.   0.806    0.181    0.806     0.806    0.806     0.630    0.875    0.833

=== Confusion Matrix ===
      a  b  c  <-- classified as
272  23  25  |  a = mid
 25  59   0  |  b = high
 29   0  94  |  c = low
```

Question 3

Weka choses feature "f_grains" as it has the highest weight assigned to it to minimize the residual errors.

Equation: $\text{predicted_carb} = 78.54 * f_grains + 6.4$

Choose **SimpleLinearRegression**

Test options

☐ Use training set
☐ Supplied test set Set...
☒ Cross-validation Folds **10**
☐ Percentage split % 66
More options...

(Num) carb

Start Stop

Result list (right-click for options)

- 21:14:21 - meta.RandomSubSpace
- 21:14:40 - meta.RandomSubSpace
- 21:14:54 - meta.RandomSubSpace
- 21:15:07 - meta.RandomSubSpace
- 21:34:13 - meta.RandomSubSpace
- 21:36:40 - meta.RandomSubSpace
- 21:37:05 - meta.RandomSubSpace
- 21:37:52 - meta.RandomSubSpace
- 21:38:22 - meta.RandomSubSpace
- 23:07:58 - functions.SimpleLinearRegression
- 23:08:20 - functions.SimpleLinearRegression
- 23:08:25 - functions.SimpleLinearRegression**

Classifier output

```
f_spirits
f_sweets
f_tea_coffee
f_water
f_wine

Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===
Linear regression on f_grains
78.54 * f_grains + 6.4

Predicting 0 if attribute value is missing.

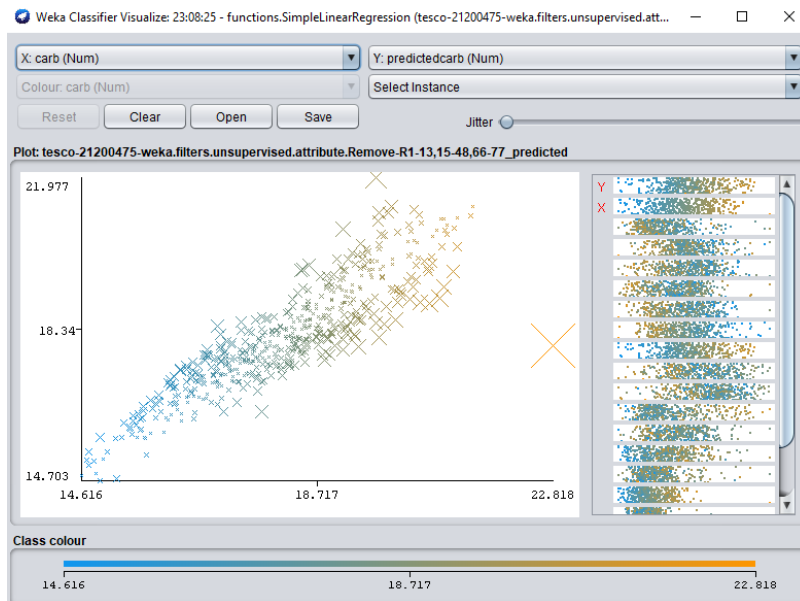
Time taken to build model: 0.03 seconds

=== Cross-validation ===
=== Summary ===
Correlation coefficient              0.8956
Mean absolute error                  0.4866
Root mean squared error              0.6342
Relative absolute error              41.9363 %
Root relative squared error          44.4492 %
Total Number of Instances            527
```

The model quality is good when we consider the correlation coefficient 89.56% with the predicted values of the feature 'carb'. However, there is 63% root mean squared error which is very high indicating that a combination of 'f_grains' feature alone would not produce the best Linear Regression model.

A possible solution to improve the model would be to choose few more features out of 17 attributes and assigning weights for each selected attribute. It would minimize the squared error on training data and improve prediction accuracy.

Analysing predicted vs actual values for the above Linear regression model shows that errors are high.



Question 4

I used the below "Bagging" method for all 3 classifiers (decision tree, Neural network and k-NN)

Classifier

Choose
Bagging -P 100 -S 1 -num-slots 1 -I 20 -W weka.classifiers.functions.MultilayerPerceptron -- -L 0.3 -M 0.2 -N 500 -V 0 -S 0 -E 20 -H 3

The below highlighted green cells have best accuracies in each of the iterations.

No of iterations	Accuracy in %ages		
	j48	NN (hidden layer = 3)	k-NN (k=3)
2	76.4706	78.7476	79.1271
4	79.6964	78.9374	80.6452
6	80.6452	79.1271	81.4042
8	79.8861	78.9374	81.5939
10	79.8861	79.1271	81.9734
12	80.2657	79.3169	81.7837
14	81.0247	79.5066	81.4042
16	81.0247	79.6964	81.0247
18	82.1632	79.6964	81.4042
20	81.9734	79.8861	81.5939

Among all 3 classifiers, bagging method significantly improved the accuracy from 76% to 82% for j48 decision tree.

We can see that level of improvement in accuracy often "level off" after an ensemble has been increased to a certain size. This is visible across all classifiers. For example: j48 decision tree classifier plateau's around interaction size = 14.

For MultiLevelPerceptron with “hidden layer = 3” - As you increase the number of iterations from 2 to 20 in MLP, time taken to build the model in Weka keeps increasing from 30 secs to 5 minutes. Additionally, time taken for validation also keeps increasing. In my case for MLP, accuracy improved marginally from 78 to 79% although number of iterations increased from 2 to 20.

For k-NN (k=3) - accuracy improved from 79 to 81% within first few iterations (2 to 6), after no new diversity is added, so ensemble accuracy will create a plateau.

Consistent improvement was seen in j48 decision tree from 76% to 82%, making it the best classifier with **iteration size = 18** having maximum accuracy 82% among the three classifiers.

I changed the **bagSizePercent** from 20%, 40% upto 100% in upward steps of 20% for iteration size = 18

The accuracy performance increases from 79% to 82% when we increase the **bagSizePercent** from 20% to 100% respectively.

j48 bagSizePercent	accuracy in %
20%	79.5066
40%	80.2657
60%	81.0247
80%	81.0247
100%	82.1632

Question 5

I used the “RandomSubSpace” method for all 3 classifiers (decision tree, Naïve Bayes and k-NN).

Note: Neural network was taking a long time to build model. Hence, I used Naïve Bayes for this question.

The below highlighted green cells have best accuracies in each of the iterations.

No of iterations	Accuracy in %ages		
	j48	Naïve Bayes	k-NN (k=3)
2	76.8501	72.1063	80.8349
4	79.3169	73.055	82.1632
6	81.7837	74.1935	83.112
8	83.112	74.0038	82.9222
10	82.1632	73.055	82.5427
12	83.3017	72.8653	82.9222
14	82.5427	72.6755	82.7324
16	83.112	72.8653	82.3529
18	83.112	72.8653	83.112
20	83.112	73.055	82.7324

Among all 3 classifiers, RandomSubSpacing method significantly improved the accuracy from 76% to 83% for j48 decision tree

For Naïve Bayes classifier, as you increase the number of iterations from 2 to 8, accuracy improved marginally from 72 to 74%. After that it was flat and lower sometimes.

For k-NN (k=3), accuracy improved from 80 to 83% within first few iterations (2 to 8), after no new diversity is added, so ensemble accuracy will create a plateau.

We can consider that **iteration size = 8** has maximum accuracy 83% among the 2 of the 3 classifiers.

I changed the **subSpaceSize** from 0.2 to 1.0 in upward steps of 0.2 for iteration size = 8

The accuracy performance decreases from 82% to 65% when we increase the **subSpaceSize** from 0.2 to 1.0 respectively.

j48 bagSizePercent	accuracy in %
0.2	82.5427
0.4	82.7324
0.6	82.5427
0.8	80.8349
1.0	65.4649

The screenshot shows the 'weka.gui.GenericObjectEditor' window for the 'weka.classifiers.meta.RandomSubSpace' classifier. The 'About' section describes the method as a decision tree based classifier that maintains highest accuracy on training data and improves on generalization accuracy as it grows in complexity. The configuration fields are as follows:

- batchSize: 100
- classifier: Choose J48 -C 0.25 -M 2
- debug: False
- doNotCheckCapabilities: False
- numDecimalPlaces: 2
- numExecutionSlots: 1
- numIterations: 8 (highlighted with a red circle)
- seed: 1
- subSpaceSize: 0.2 (highlighted with a red circle)

Question 6

Set of classifiers is expected to benefit from bagging techniques:

Bagging encourages diversity in the ensemble, works better for "unstable" classifiers - e.g. decision trees, neural networks.

Set of classifiers is expected to benefit from random subspacing techniques:

RandomSubspacing adds more instability into the classifier (diversity) since different features are used when calculating distances.

k-NN would benefit from it

For my dataset, determine the best ensemble strategy for each of these classifiers. Discuss if this is in line with what you expected.

RandomSubspacing method worked well for k-NN classifier since it improved accuracy from 80% to approx. 83%.

Bagging method significantly improved the accuracy from 76% to 82% for j48 decision tree.

Bagging methods worked ordinary for Neural Networks. Accuracy did not improve drastically (78% to 79%) compared to other classifiers.

I used RandomSubSpace for Naïve Bayes and found it helpful to improve the accuracy from 72% to 74%. However, not as much as other classifiers stated above.

The results are in line with my expectations based on the below ensemble theories.

More accurate classifiers contribute more to the ensemble strategies.

Bagging can often reduce variance part of error.

Boosting can often reduce variance and bias, since it focuses on misclassified examples.

Boosting may sometimes increase error, as it is susceptible to noise and may lead to overfitting.