Classification results

The best classification results in this dataset were obtained using Naïve Bayes algorithm. The best overall accuracy was 77.81% and it was acquired in the data set, where cluster taxonomies were used; the taxonomies were built using DTWARP distance among curves, Ward linkage among clusters and Information Gain as the metric to select the best cut in the taxonomies. This combination also showed a good sensitivity (64.05%) and specificity (85.04%). The combination that identified cancer-specific breaths most accurately included Euclidean distance, Ward linkage and Symmetrical Uncertainty-based feature selection approach, identifying 66.54% of cancer-specific breaths correctly. The feature that allowed detecting the healthy participants' breaths most accurately was the average signal, however, several combinations of taxonomy-based classification approaches showed specificities above 85% and the differences in specificity here were not statistically significant. The best AU-ROC was also achieved using a taxonomy-based dataset using a combination of DTWARP distance, complete linkage and Information Gain metric. This means that with similar specificity the proposed approach improves overall accuracy, sensitivity and/or ROC.

**Table S1.** Classification results (and 95%CI) using Naïve Bayes classifiers.

| Feature | Overall Accuracy | Sensitivity | Specificity | AU-ROC |
|---|---|---|---|---|
| Minimum | 72.18% | 46.9% | 85.51% | 0.774 |
| | (71.49%…72.87%) | (45.39%…48.41%) | (84.76%…86.26%) | (0.765…0.782) |
| Average | 74.21% | 51.85% | 86.02% | 0.790 |
| | (73.5%…74.91%) | (50.35%…53.34%) | (85.27%…86.76%) | (0.782…0.798) |
| Maximum | 73.7% | 53.44% | 84.38% | 0.783 |
| | (72.96%…74.44%) | (51.94%…54.94%) | (83.6%…85.16%) | (0.775…0.791) |
| Average of the last 10 time points | 73.74% | 53.00% | 84.67% | 0.788 |
| | (73.02%…74.45%) | (51.51%…54.49%) | (83.9%…85.44%) | (0.78…0.797) |
| Area under the curve | 73.75% | 50.77% | 85.88% | 0.785 |
| | (73.04%…74.47%) | (49.28%…52.26%) | (85.13%…86.64%) | (0.776…0.793) |
| Cluster (DTWARP distance, complete linkage, InfoGain) | 77.33% | 63.32% | 84.70% | 0.830 |
| | (76.62%...78.04%) | (61.93%...64.72%) | (83.93%...85.48%) | (0.823…0.838) |
| Cluster (DTWARP distance, complete linkage, ReliefF) | 73.79% | 59.26% | 81.42% | 0.809 |
| | (73.06%...74.51%) | (57.88%...60.63%) | (80.55%...82.28%) | (0.801…0.817) |
| Cluster (DTWARP distance, complete linkage, Symm.Unc.) | 77.01% | 64.10% | 83.81% | 0.823 |
| | (76.33%…77.70%) | (62.72%...65.47%) | (83.01%...84.60%) | (0.815…0.831) |
| Cluster (DTWARP distance, Ward linkage, InfoGain) | 77.81% | 64.05% | 85.04% | 0.808 |
| | (77.15%…78.48%) | (62.66%…65.44%) | (84.29%…85.78%) | (0.8…0.816) |
| Cluster (DTWARP distance, Ward linkage, ReliefF) | 75.07% | 60.92% | 82.5% | 0.796 |
| | (74.39%…75.75%) | (59.53%…62.3%) | (81.71%…83.29%) | (0.788…0.804) |
| Cluster (DTWARP distance, Ward linkage, Symm.Unc.) | 77.5% | 62.2% | 85.53% | 0.809 |
| | (76.83%…78.18%) | (60.79%…63.6%) | (84.79%…86.28%) | (0.801…0.817) |
| Cluster (Euclidean distance, complete linkage, InfoGain) | 75.40% | 61.58% | 82.67% | 0.817 |
| | (74.70%...76.09%) | (60.17%...62.98%) | (81.87%...83.46%) | (0.809…0.825) |
| Cluster (Euclidean distance, complete linkage, ReliefF) | 76.03% | 58.79% | 85.10% | 0.792 |
| | (75.36%...76.70%) | (57.41%...60.16%) | (84.35%...85.85%) | (0.784…0.801) |
| Cluster (Euclidean distance, complete linkage, Symm.Unc.) | 75.59% | 58.71% | 84.47% | 0.814 |
| | (74.89%…76.28%) | (57.29%...60.13%) | (83.69%...85.24%) | (0.806…0.822) |
| Cluster (Euclidean distance, Ward linkage, InfoGain) | 76.39% | 65.02% | 82.36% | 0.821 |
| | (75.69%…77.1%) | (63.64%…66.4%) | (81.55%…83.17%) | (0.813…0.829) |
| Cluster (Euclidean distance, Ward linkage, ReliefF) | 74.47% | 64.61% | 79.66% | 0.793 |
| | (73.76%…75.18%) | (63.26%…65.96%) | (78.8%…80.52%) | (0.785…0.802) |
| Cluster (Euclidean distance, Ward linkage, Symm.Unc.) | 77.1% | 66.54% | 82.64% | 0.817 |
| | (76.41%…77.79%) | (65.21%…67.87%) | (81.83%…83.45%) | (0.81…0.825) |

The second-best classification results in this dataset were obtained using Random Forest algorithm. The best overall accuracy was 75.87%, obtained using the data set, where cluster taxonomies were used; the taxonomies were built using DTWARP distance among curves, Ward linkage among clusters and Symmetrical uncertainty as the metric to select the best cut in the taxonomies. All of the feautes and approaches showed poor sensitivity (less than 50%). However, specificity was significantly improved using taxonomies, reaching up to 92.39% specificity for the same combination that showed the best overall accuracy. The two best AU-ROCs (0.800 and 0.799) were acquired using taxonomies: Euclidean distance, complete linkage, ReliefF-based feature selection, and DTWARP distance, complete linkage, ReliefF-based feature selection respectively.

**Table S2.** Classification results using Random Forests.

| Feature | Overall Accuracy | Sensitivity | Specificity | AU-ROC |
|---|---|---|---|---|
| Minimum | 69.89% | 45.3% | 82.93% | 0.763 |
| | (69.19%…70.59%) | (43.87%…46.73%) | (82.12%…83.74%) | (0.755…0.771) |
| Average | 70.58% | 47% | 83.13% | 0.781 |
| | (69.86%…71.29%) | (45.56%…48.45%) | (82.3%…83.95%) | (0.773…0.789) |
| Maximum | 70.33% | 43.23% | 84.72% | 0.771 |
| | (69.66%…71.01%) | (41.82%…44.64%) | (83.95%…85.49%) | (0.762…0.779) |
| Average of the last 10 time points | 70.97% | 48.79% | 82.78% | 0.785 |
| | (70.25%…71.69%) | (47.32%…50.25%) | (81.96%…83.6%) | (0.777…0.793) |
| Area under the curve | 70.51% | 46.54% | 83.27% | 0.781 |
| | (69.79%…71.23%) | (45.1%…47.99%) | (82.44%…84.1%) | (0.773…0.789) |
| Cluster (DTWARP distance, complete linkage, InfoGain) | 74.90% | 45.77% | 90.37% | 0.783 |
| | (74.26%...7.54%) | (44.38%...47.17%) | (89.72%...91.02%) | (0.775…0.791) |
| Cluster (DTWARP distance, complete linkage, ReliefF) | 75.64% | 45.95% | 91.39% | 0.799 |
| | (75.01%...76.26%) | (44.56%...47.34%) | (90.80%...91.99%) | (0.792…0.807) |
| Cluster (DTWARP distance, complete linkage, Symm.Unc.) | 74.47% | 45.98% | 89.58% | 0.786 |
| | (73.83%...75.11%) | (44.60%...47.36%) | (88.91%...90.24%) | (0.778…0.794) |
| Cluster (DTWARP distance, Ward linkage, InfoGain) | 75.37% | 44.38% | 91.79% | 0.787 |
| | (74.76%…75.97%) | (42.99%…45.77%) | (91.21%…92.38%) | (0.778…0.795) |
| Cluster (DTWARP distance, Ward linkage, ReliefF) | 75.01% | 46.52% | 90.12% | 0.761 |
| | (74.39%…75.63%) | (45.13%…47.91%) | (89.5%…90.74%) | (0.752…0.769) |
| Cluster (DTWARP distance, Ward linkage, Symm.Unc.) | 75.87% | 44.69% | 92.39% | 0.777 |
| | (75.27%...76.47%) | (43.31%…46.07%) | (91.84%…92.94%) | (0.769…0.786) |
| Cluster (Euclidean distance, complete linkage, InfoGain) | 72.29% | 40.57% | 89.18% | 0.792 |
| | (71.65%...72.94%) | (39.20%...41.95%) | (88.52%...89.84%) | (0.784…0.800) |
| Cluster (Euclidean distance, complete linkage, ReliefF) | 74.51% | 46.04% | 89.66% | 0.800 |
| | (73.90%...75.13%) | (44.68%...47.41%) | (89.01%...90.30%) | (0.792…0.808) |
| Cluster (Euclidean distance, complete linkage, Symm.Unc.) | 72.04% | 39.59% | 89.27% | 0.791 |
| | (71.41%...72.66%) | (38.21%...40.98%) | (88.61%...89.93%) | (0.783…0.799) |
| Cluster (Euclidean distance, Ward linkage, InfoGain) | 72.88% | 42.57% | 89.02% | 0.75 |
| | (72.25%…73.52%) | (41.18%…43.96%) | (88.36%…89.69%) | (0.741…0.759) |
| Cluster (Euclidean distance, Ward linkage, ReliefF) | 72.78% | 46.13% | 86.93% | 0.747 |
| | (72.12%…73.45%) | (44.75%…47.52%) | (86.22%…87.64%) | (0.738…0.756) |
| Cluster (Euclidean distance, Ward linkage, Symm.Unc.) | 72.49% | 44.57% | 87.33% | 0.745 |
| | (71.83%…73.14%) | (43.17%…45.98%) | (86.63%...88.03%) | (0.737…0.754) |

The third-best classification results in this dataset were obtained using SVM algorithm. The best overall accuracy was 75.10%, obtained using the data set, where cluster taxonomies were not applied. The data set used the average signal value of the last 10 time-points of the curve. The second-best result (74.87%) was acquired in the data set that used DTWARP distance, complete linkage and Information Gain metric for feature selection. The best sensitivity (61.05%) was acquired in the data set, where the combination of Euclidean distance, Ward linkage and Information Gain metric for feature selection. Specificity was lower in the data sets where taxonomies were used, and the best specificity (91.14%) was obtained using the minimum value of the sensor response curve, however sensitivity in this case was unacceptably low (41.31%). The best AU-ROC (0.716) was acquired using taxonomies: DTWARP distance, complete linkage and Information Gain for feature selection.

**Table S3.** Classification results using SVMs.

| Feature | Overall Accuracy | Sensitivity | Specificity | AU-ROC |
|---|---|---|---|---|
| Minimum | 73.84% | 41.31% | 91.14% | 0.662 |
| | (73.23%…74.45%) | (39.99%…42.64%) | (90.51%…91.78%) | (0.655…0.669) |
| Maximum | 73.45% | 45.72% | 88.2% | 0.67 |
| | (72.79%…74.11%) | (44.31%…47.14%) | (87.53%…88.88%) | (0.662…0.677) |
| Average | 74.26% | 43.16% | 90.74% | 0.653 |
| | (73.64%…74.87%) | (41.77%…44.55%) | (90.12%…91.37%) | (0.645…0.66) |
| Average of the last 10 time points | 75.1% | 48.33% | 89.27% | 0.688 |
| | (74.47%…75.74%) | (46.94%…49.71%) | (88.62%…89.92%) | (0.68…0.696) |
| Area under the curve | 72.75% | 40.86% | 89.68% | 0.67 |
| | (72.13%…73.37%) | (39.48%…42.25%) | (89.02%…90.33%) | (0.662…0.677) |
| Cluster (DTWARP distance, complete linkage, InfoGain) | 74.87% | 60.73% | 82.48% | 0.716 |
| | (74.16%...75.59%) | (59.31%...62.15%) | (81.66%...83.30%) | (0.708…0.724) |
| Cluster (DTWARP distance, complete linkage, ReliefF) | 73.10% | 59.66% | 80.34% | 0.700 |
| | (72.34%...73.86%) | (58.15%...61.16%) | (79.45%...81.22%) | (0.691…0.709) |
| Cluster (DTWARP distance, complete linkage, Symm.Unc.) | 73.81% | 58.13% | 82.18% | 0.702 |
| | (73.07%...74.55%) | (56.69%...59.56%) | (81.32%...83.05%) | (0.693…0.710) |
| Cluster (DTWARP distance, Ward linkage, InfoGain) | 72.22% | 58.25% | 79.71% | 0.69 |
| | (71.45%…72.99%) | (56.76%…59.74%) | (78.83%…80.59%) | (0.681…0.698) |
| Cluster (DTWARP distance, Ward linkage, ReliefF) | 70.28% | 54.33% | 78.76% | 0.665 |
| | (69.54%…71.02%) | (52.9%…55.76%) | (77.85%…79.67%) | (0.657…0.674) |
| Cluster (DTWARP distance, Ward linkage, Symm.Unc.) | 72.67% | 58.46% | 80.28% | 0.694 |
| | (71.91%…73.42%) | (56.99%…59.92%) | (79.43%…81.13%) | (0.685…0.702) |
| Cluster (Euclidean distance, complete linkage, InfoGain) | 71.45% | 55.94% | 79.76% | 0.679 |
| | (70.72%...72.17%) | (54.48%...57.39%) | (78.89%...80.62%) | (0.670…0.687) |
| Cluster (Euclidean distance, complete linkage, ReliefF) | 70.53% | 57.50% | 77.53% | 0.675 |
| | (69.78%...71.28%) | (56.00%...58.99%) | (76.63%...78.43%) | (0.667…0.684) |
| Cluster (Euclidean distance, complete linkage, Symm.Unc.) | 73.67% | 58.21% | 81.92% | 0.701 |
| | (72.97%...74.37%) | (56.80%...59.62%) | (81.08%...82.77%) | (0.693…0.709) |
| Cluster (Euclidean distance, Ward linkage, InfoGain) | 73.86% | 61.05% | 80.72% | 0.709 |
| | (73.07%…74.65%) | (59.52%…62.58%) | (79.84%…81.6%) | (0.7…0.718) |
| Cluster (Euclidean distance, Ward linkage, ReliefF) | 71.24% | 55.58% | 79.58% | 0.676 |
| | (70.49%…71.99%) | (54.12%…57.03%) | (78.69%…80.46%) | (0.667…0.684) |
| Cluster (Euclidean distance, Ward linkage, Symm.Unc.) | 74.6% | 59.97% | 82.41% | 0.712 |
| | (73.88%…75.32%) | (58.52%…61.41%) | (81.59%…83.23%) | (0.704…0.72) |