# Supplementary materials for paper „An ensemble classifier based on kNN with an interval threshold strategy".
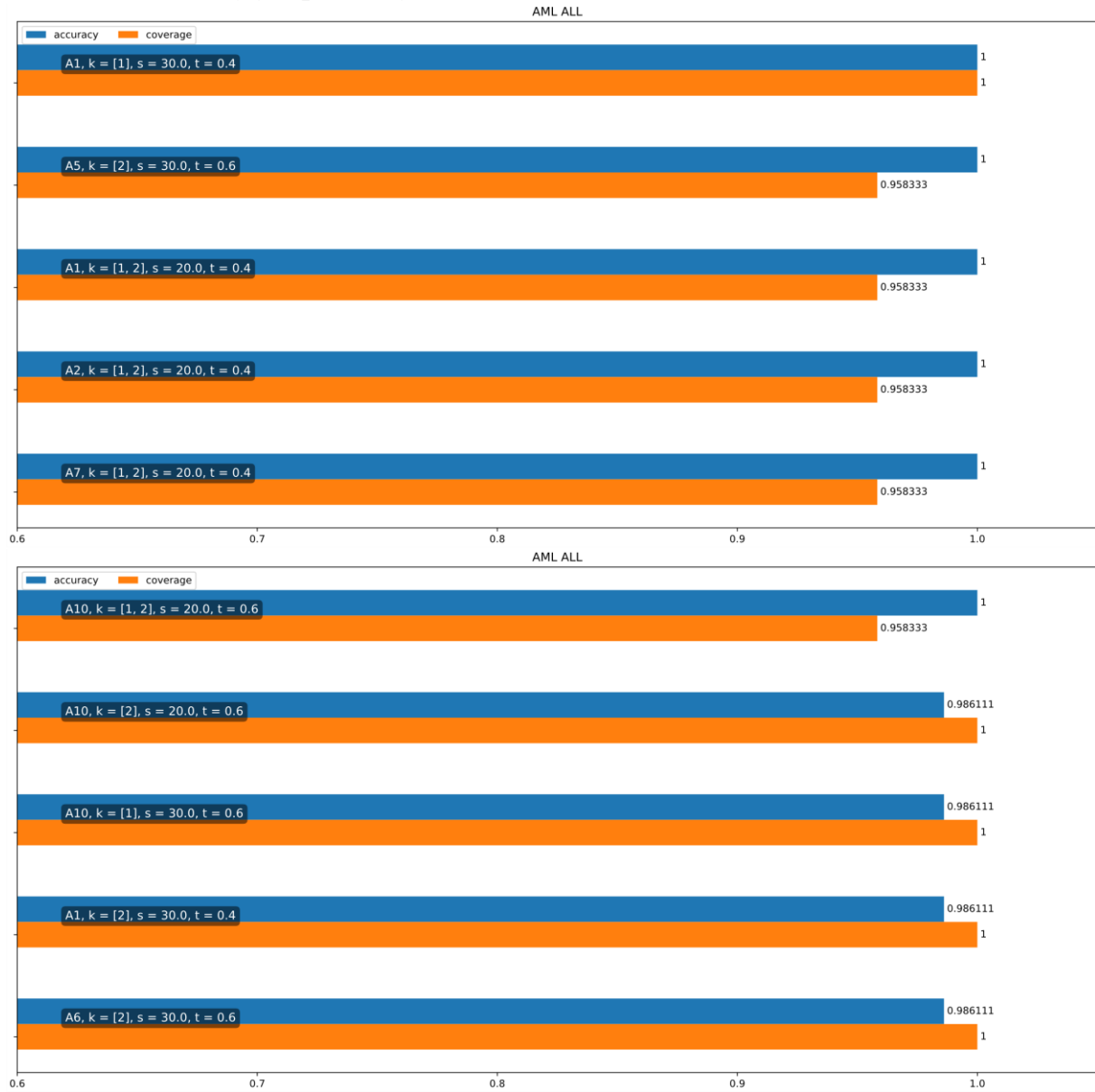
U. Bentkowska, M. Mrukowicz, W. Gałka, K. Lech

# 1. Figures

The results presented in Tables 2,4,6,8,10 are plotted here, respectively. **The data in the figures are sorted by accuracy. Average values from folds in LOO are presented for both quality measures (accuracy and coverage) for the established hyperparameters: aggregation, k, s and t.**

## 1.1. AML ALL (Lymphoma) dataset



AML ALL

- accuracy  coverage

A1, k = [1], s = 30.0, t = 0.4 — 1 / 1
A5, k = [2], s = 30.0, t = 0.6 — 1 / 0.958333
A1, k = [1, 2], s = 20.0, t = 0.4 — 1 / 0.958333
A2, k = [1, 2], s = 20.0, t = 0.4 — 1 / 0.958333
A7, k = [1, 2], s = 20.0, t = 0.4 — 1 / 0.958333



AML ALL

- accuracy  coverage

A10, k = [1, 2], s = 20.0, t = 0.6 — 1 / 0.958333
A10, k = [2], s = 20.0, t = 0.6 — 0.986111 / 1
A10, k = [1], s = 30.0, t = 0.6 — 0.986111 / 1
A1, k = [2], s = 30.0, t = 0.4 — 0.986111 / 1
A6, k = [2], s = 30.0, t = 0.6 — 0.986111 / 1

AML ALL

| | accuracy | coverage | |
|---|---|---|---|
| A10, k = [2], s = 30.0, t = 0.5 | | | 0.986111 |
| | | | 1 |
| A10, k = [1], s = 50.0, t = 0.6 | | | 0.986111 |
| | | | 1 |
| A10, k = [2], s = 50.0, t = 0.5 | | | 0.986111 |
| | | | 1 |
| A1, k = [2], s = 20.0, t = 0.4 | | | 0.985915 |
| | | | 0.986111 |
| A5, k = [2], s = 20.0, t = 0.6 | | | 0.985915 |
| | | | 0.986111 |

AML ALL

| | accuracy | coverage | |
|---|---|---|---|
| A1, k = [1], s = 50.0, t = 0.4 | | | 0.985915 |
| | | | 0.986111 |
| A1, k = [1], s = 5.0, t = 0.4 | | | 0.985714 |
| | | | 0.972222 |
| A1, k = [1], s = 30.0, t = 0.5 | | | 0.985714 |
| | | | 0.972222 |
| A5, k = [2], s = 30.0, t = 0.5 | | | 0.985714 |
| | | | 0.972222 |
| A10, k = [1, 2], s = 5.0, t = 0.6 | | | 0.985714 |
| | | | 0.972222 |

## 1.2. Colon



COLON

| | accuracy | coverage |
|---|---|---|
| A5, k = [2], s = 2.0, t = 0.5 | 0.912281 | 0.919355 |
| A6, k = [2], s = 2.0, t = 0.5 | 0.912281 | 0.919355 |
| A10, k = [2], s = 2.0, t = 0.5 | 0.912281 | 0.919355 |
| A1, k = [1, 2], s = 2.0, t = 0.4 | 0.909091 | 0.887097 |
| A5, k = [1, 2], s = 2.0, t = 0.4 | 0.909091 | 0.887097 |



COLON

| | accuracy | coverage |
|---|---|---|
| A5, k = [1, 2], s = 2.0, t = 0.5 | 0.909091 | 0.887097 |
| A6, k = [1, 2], s = 2.0, t = 0.5 | 0.909091 | 0.887097 |
| A10, k = [1, 2], s = 2.0, t = 0.4 | 0.909091 | 0.887097 |
| A10, k = [1, 2], s = 2.0, t = 0.5 | 0.909091 | 0.887097 |
| A3, k = [2], s = 2.0, t = 0.6 | 0.896552 | 0.935484 |

**COLON**

- accuracy
- coverage

A5, k = [1, 2], s = 2.0, t = 0.6 — 0.892857 / 0.903226
A6, k = [1, 2], s = 2.0, t = 0.6 — 0.892857 / 0.903226
A9, k = [1, 2], s = 2.0, t = 0.6 — 0.892857 / 0.903226
A10, k = [1, 2], s = 2.0, t = 0.6 — 0.892857 / 0.903226
A1, k = [1, 2], s = 5.0, t = 0.4 — 0.890909 / 0.887097

**COLON**

- accuracy
- coverage

A4, k = [2], s = 2.0, t = 0.6 — 0.896552 / 0.935484
A8, k = [2], s = 2.0, t = 0.6 — 0.896552 / 0.935484
A5, k = [2], s = 5.0, t = 0.6 — 0.896552 / 0.935484
A9, k = [2], s = 5.0, t = 0.6 — 0.892857 / 0.903226
A2, k = [1, 2], s = 2.0, t = 0.4 — 0.892857 / 0.903226

# 1.3. DLBCL



DLBCL

- A1, k = [1, 2], s = 20.0, t = 0.6 — accuracy 0.926829, coverage 0.87234
- A2, k = [1, 2], s = 20.0, t = 0.6 — accuracy 0.926829, coverage 0.87234
- A7, k = [1, 2], s = 20.0, t = 0.6 — accuracy 0.926829, coverage 0.87234
- A1, k = [1, 2], s = 10.0, t = 0.6 — accuracy 0.925, coverage 0.851064
- A2, k = [1, 2], s = 10.0, t = 0.6 — accuracy 0.925, coverage 0.851064

DLBCL

- A7, k = [1, 2], s = 10.0, t = 0.6 — accuracy 0.925, coverage 0.851064
- A1, k = [1, 2], s = 30.0, t = 0.5 — accuracy 0.925, coverage 0.851064
- A2, k = [1, 2], s = 30.0, t = 0.5 — accuracy 0.925, coverage 0.851064
- A7, k = [1, 2], s = 30.0, t = 0.5 — accuracy 0.925, coverage 0.851064
- A1, k = [2], s = 15.0, t = 0.6 — accuracy 0.909091, coverage 0.93617

DLBCL

accuracy coverage

A1, k = [1], s = 30.0, t = 0.5 — 0.906977 / 0.914894

A7, k = [1, 2], s = 2.0, t = 0.5 — 0.902439 / 0.87234

A1, k = [1, 2], s = 5.0, t = 0.5 — 0.902439 / 0.87234

A2, k = [1, 2], s = 5.0, t = 0.5 — 0.902439 / 0.87234

A7, k = [1, 2], s = 5.0, t = 0.5 — 0.902439 / 0.87234

DLBCL

accuracy coverage

A1, k = [1, 2], s = 15.0, t = 0.5 — 0.902439 / 0.87234

A1, k = [1, 2], s = 15.0, t = 0.6 — 0.902439 / 0.87234

A2, k = [1, 2], s = 15.0, t = 0.5 — 0.902439 / 0.87234

A2, k = [1, 2], s = 15.0, t = 0.6 — 0.902439 / 0.87234

A7, k = [1, 2], s = 15.0, t = 0.5 — 0.902439 / 0.87234

## 1.4.    Prostate



prostate

| | |
|---|---|
| A1, k = [1, 2], s = 20.0, t = 0.6 | 0.957265 |
| | 0.860294 |
| A2, k = [1, 2], s = 20.0, t = 0.6 | 0.957265 |
| | 0.860294 |
| A7, k = [1, 2], s = 20.0, t = 0.6 | 0.957265 |
| | 0.860294 |
| A1, k = [1, 2], s = 10.0, t = 0.6 | 0.956897 |
| | 0.852941 |
| A2, k = [1, 2], s = 10.0, t = 0.6 | 0.956897 |
| | 0.852941 |



prostate

| | |
|---|---|
| A7, k = [1, 2], s = 10.0, t = 0.6 | 0.956897 |
| | 0.852941 |
| A1, k = [1, 2], s = 15.0, t = 0.6 | 0.950413 |
| | 0.889706 |
| A2, k = [1, 2], s = 15.0, t = 0.6 | 0.950413 |
| | 0.889706 |
| A7, k = [1, 2], s = 15.0, t = 0.6 | 0.950413 |
| | 0.889706 |
| A1, k = [1, 2], s = 30.0, t = 0.6 | 0.95 |
| | 0.882353 |

prostate

| | accuracy | coverage |
|---|---|---|
| A2, k = [1, 2], s = 30.0, t = 0.6 | 0.95 | 0.882353 |
| A7, k = [1, 2], s = 30.0, t = 0.6 | 0.95 | 0.882353 |
| A1, k = [1, 2], s = 5.0, t = 0.5 | 0.939655 | 0.852941 |
| A2, k = [1, 2], s = 5.0, t = 0.5 | 0.939655 | 0.852941 |
| A7, k = [1, 2], s = 5.0, t = 0.5 | 0.939655 | 0.852941 |

prostate

| | accuracy | coverage |
|---|---|---|
| A1, k = [2], s = 15.0, t = 0.6 | 0.933333 | 0.992647 |
| A1, k = [1, 2], s = 15.0, t = 0.5 | 0.931624 | 0.860294 |
| A2, k = [1, 2], s = 15.0, t = 0.5 | 0.931624 | 0.860294 |
| A7, k = [1, 2], s = 15.0, t = 0.5 | 0.931624 | 0.860294 |
| A1, k = [1], s = 20.0, t = 0.6 | 0.931298 | 0.963235 |

## 1.5. Ovarian

ovarian

| series | accuracy | coverage |
|---|---|---|
| A1, k = [2], s = 50.0, t = 0.4 | 1 | 1 |
| A10, k = [2], s = 50.0, t = 0.6 | 1 | 1 |
| A1, k = [2], s = 15.0, t = 0.4 | 1 | 0.996047 |
| A1, k = [2], s = 15.0, t = 0.5 | 1 | 0.996047 |
| A1, k = [1, 2], s = 20.0, t = 0.4 | 1 | 0.996047 |

ovarian

| series | accuracy | coverage |
|---|---|---|
| A2, k = [1, 2], s = 20.0, t = 0.4 | 1 | 0.996047 |
| A7, k = [1, 2], s = 20.0, t = 0.4 | 1 | 0.996047 |
| A1, k = [1], s = 10.0, t = 0.5 | 1 | 0.992095 |
| A10, k = [1], s = 50.0, t = 0.6 | 1 | 0.992095 |
| A5, k = [2], s = 50.0, t = 0.6 | 1 | 0.992095 |

# 2. Statistical tests

## 2.1. Methodology

Based on the filtered results in Tables 3, 5, 7, 9, 11 there are obtained series of accuracy for fixed aggregation and interval-threshold but different s and t. However, only series with sufficient number of samples (above or equal to 10) and with mean accuracy comparable to the SOTA were considered.

To compare this series with SOTA models first, the Shapiro-Wilk test was performed to determine if one sample t-test is available (with additional check to not contain extreme outliers). If the series was not normally distributed, the Miao, Gel, and Gastwirth (MGG) symmetry test was performed and if the distribution was symmetric, the exact one sample Wilcoxon signed rank test was performed (due to the ties, sometimes present in data series). For some series of data neither t-test nor exact one sample Wilcoxon signed rank test was available. All tests were performed using a 0.05 significance level in R programming language.

*Table 1. Shapiro-Wilk test. Bold font was used in the cases where the distribution is normal*

| series | W | p-value |
|---|---|---|
| AML ALL, A1 k=[1] | 0,891843 | 0,024389 |

| series | statistic | p-value |
|---|---|---|
| AML ALL, A1 k=[2] | 0,862589 | 0,007054 |
| AML ALL, A1 k=[1, 2] | 0,739264 | 8,86E-05 |
| AML ALL, A2 k=[1, 2] | 0,739264 | 8,86E-05 |
| AML ALL, A5 k=[1, 2] | 0,854774 | 0,012692 |
| AML ALL, A7 k=[1, 2] | 0,744107 | 0,000103 |
| AML ALL, A10 k=[1, 2] | 0,905686 | 0,045165 |
| **COLON, A1 k=[1]** | **0,912078** | **0,125704** |
| **COLON, A1 k=[2]** | **0,937927** | **0,324388** |
| **COLON, A1 k=[1, 2]** | **0,968286** | **0,852925** |
| **COLON, A10 k=[2]** | **0,932202** | **0,294286** |
| **COLON, A10 k=[1, 2]** | **0,917158** | **0,333862** |
| **COLON, A2 k=[1, 2]** | **0,948711** | **0,540873** |
| **COLON, A5 k=[2]** | **0,931437** | **0,319595** |
| **COLON, A6 k=[2]** | **0,953122** | **0,7055** |
| **COLON, A7 k=[1, 2]** | **0,949114** | **0,547042** |
| OVARIAN A1 k=[1] | 0,775826 | 0,000287 |
| OVARIAN A1 k=[2] | 0,738545 | 8,67E-05 |
| OVARIAN A1 k=[1, 2] | 0,660697 | 9,26E-06 |
| OVARIAN A2 k=[1, 2] | 0,660356 | 9,18E-06 |
| OVARIAN A5 k=[1] | 0,794904 | 0,003165 |
| OVARIAN A5 k=[2] | 0,72897 | 8,99E-05 |
| OVARIAN A5 k=[1, 2] | 0,622027 | 0,000102 |
| OVARIAN A6 k=[1] | 0,820443 | 0,01748 |
| OVARIAN A6 k=[2] | 0,843133 | 0,010814 |
| OVARIAN A6 k=[1, 2] | 0,703955 | 0,000984 |
| OVARIAN A7 k=[1, 2]v | 0,660697 | 9,26E-06 |
| OVARIAN A10 k=[1] | 0,794389 | 0,000713 |
| OVARIAN A10 k=[2] | 0,754469 | 0,000143 |
| OVARIAN A10 k=[1, 2] | 0,77816 | 0,000415 |
| prostate A1 k=[1, 2] | 0,782761 | 0,001189 |
| **prostate A1 k=[1]** | **0,945853** | **0,335122** |
| prostate A1 k=[2] | 0,896272 | 0,029648 |
| prostate A2 k=[1, 2] | 0,804963 | 0,002375 |
| **prostate A5 k=[2]** | **0,921433** | **0,298003** |
| prostate A7 k=[1, 2] | 0,786981 | 0,001352 |
| **prostate A10 k=[2]** | **0,925** | **0,29275** |
| **DLBCL A1 k=[1, 2]** | **0,906593** | **0,120066** |
| **DLBCL A2 k=[1, 2]** | **0,91151** | **0,142824** |
| **DLBCL A7 k=[1, 2]** | **0,904138** | **0,110106** |

*Table 2. MGG Symmetry Test. Bold font was used in the cases where the data are symmetric*

| series | statistic | p-value |
|---|---|---|
| **AML ALL, A1 k=[1]** | **-1,23555** | **0,216626** |
| AML ALL, A1 k=[2] | -2,57292 | 0,010084 |
| **AML ALL, A1 k=[1, 2]** | **-0,91373** | **0,360856** |
| **AML ALL, A2 k=[1, 2]** | **-0,91373** | **0,360856** |

| | | |
|---|---|---|
| **AML ALL, A5 k=[1, 2]** | **-1,92073** | **0,054766** |
| **AML ALL, A7 k=[1, 2]** | **-0,88522** | **0,376036** |
| AML ALL, A10 k=[1, 2] | -2,69147 | 0,007114 |
| OVARIAN, A1 k=[1] | -3,56724 | 0,000361 |
| OVARIAN, A1 k=[2] | -3,23553 | 0,001214 |
| OVARIAN, A1 k=[1, 2] | -4,8396 | 1,3E-06 |
| OVARIAN, A2 k=[1, 2] | -4,8396 | 1,3E-06 |
| **OVARIAN, A5 k=[1]** | **-1,19308** | **0,232838** |
| OVARIAN, A5 k=[2] | -3,51416 | 0,000441 |
| **OVARIAN, A5 k=[1, 2]** | **-0,85229** | **0,394055** |
| OVARIAN, A6 k=[1] | -2,36617 | 0,017973 |
| OVARIAN, A6 k=[2] | -2,39562 | 0,016592 |
| OVARIAN, A6 k=[1, 2] | -2,55678 | 0,010565 |
| OVARIAN, A7 k=[1, 2]v | -4,8396 | 1,3E-06 |
| **OVARIAN, A10 k=[1]** | **-1,91238** | **0,055827** |
| OVARIAN, A10 k=[2] | -3,87287 | 0,000108 |
| OVARIAN, A10 k=[1, 2] | -2,96867 | 0,002991 |
| prostate, A1 k=[1, 2] | -2,69562 | 0,007026 |
| prostate, A1 k=[2] | -2,40168 | 0,01632 |
| prostate, A2 k=[1, 2] | -2,74545 | 0,006043 |
| prostate, A7 k=[1, 2] | -2,87105 | 0,004091 |

*Table 3 Exact Wilcoxon signed rank test. Bold font was used in the cases of significant difference*

| series | V | p-value |
|---|---|---|
| **AML ALL A1 k=[1] knn k=1** | **204** | **0,001157** |
| **AML ALL A1 k=[1] knn k=2** | **231** | **9,54E-07** |
| **AML ALL A1 k=[1] knn k=3** | **224** | **1,72E-05** |
| **AML ALL A1 k=[1] knn k=5** | **231** | **9,54E-07** |
| **AML ALL A1 k=[1] random forest** | **50** | **0,021001** |
| **AML ALL A1 k=[1, 2] knn k=1** | **213** | **0,00023** |
| **AML ALL A1 k=[1, 2] knn k=2** | **230** | **1,91E-06** |
| **AML ALL A1 k=[1, 2] knn k=3** | **230** | **1,91E-06** |
| **AML ALL A1 k=[1, 2] knn k=5** | **230** | **1,91E-06** |
| AML ALL A1 k=[1, 2] random forest | 101 | 0,626988 |
| **AML ALL A2 k=[1, 2] knn k=1** | **213** | **0,00023** |
| **AML ALL A2 k=[1, 2] knn k=2** | **230** | **1,91E-06** |
| **AML ALL A2 k=[1, 2] knn k=3** | **230** | **1,91E-06** |
| **AML ALL A2 k=[1, 2] knn k=5** | **230** | **1,91E-06** |
| AML ALL A2 k=[1, 2] random forest | 101 | 0,626988 |
| **AML ALL A5 k=[1, 2] knn k=1** | **153** | **1,53E-05** |
| **AML ALL A5 k=[1, 2] knn k=2** | **153** | **1,53E-05** |
| **AML ALL A5 k=[1, 2] knn k=3** | **153** | **1,53E-05** |
| **AML ALL A5 k=[1, 2] knn k=5** | **153** | **1,53E-05** |
| AML ALL A5 k=[1, 2] random forest | 90 | 0,538208 |
| **AML ALL A7 k=[1, 2] knn k=1** | **213** | **0,00023** |
| **AML ALL A7 k=[1, 2] knn k=2** | **231** | **9,54E-07** |

| | | |
|---|---|---|
| **AML ALL A7 k=[1, 2] knn k=3** | **230** | **1,91E-06** |
| **AML ALL A7 k=[1, 2] knn k=5** | **231** | **9,54E-07** |
| AML ALL A7 k=[1, 2] random forest | 101 | 0,626988 |
| OVARIAN A5 k=[1] knn k=1 | 91 | 0,080994 |
| OVARIAN A5 k=[1] knn k=2 | 91 | 0,080994 |
| OVARIAN A5 k=[1] knn k=3 | 87 | 0,131592 |
| OVARIAN A5 k=[1] knn k=5 | 91 | 0,080994 |
| **OVARIAN A5 k=[1] random forest** | **5** | **0,00061** |
| **OVARIAN A5 k=[1, 2] knn k=1** | **90** | **0,000488** |
| **OVARIAN A5 k=[1, 2] knn k=2** | **90** | **0,000488** |
| **OVARIAN A5 k=[1, 2] knn k=3** | **88** | **0,001221** |
| **OVARIAN A5 k=[1, 2] knn k=5** | **90** | **0,000488** |
| **OVARIAN A5 k=[1, 2] random forest** | **14** | **0,023682** |
| OVARIAN A10 k=[1] knn k=1 | 124 | 0,491716 |
| OVARIAN A10 k=[1] knn k=2 | 113 | 0,776831 |
| OVARIAN A10 k=[1] knn k=3 | 100 | 0,862623 |
| OVARIAN A10 k=[1] knn k=5 | 113 | 0,776831 |
| **OVARIAN A10 k=[1] random forest** | **8** | **4,39E-05** |

*Table 4 One sample t-test. Bold font was used in the cases of significant difference*

| series | statistic | p |
|---|---|---|
| COLON A1 k=[1] knn k=1 | 1,734633 | 0,103 |
| COLON A1 k=[1] knn k=2 | 1,734633 | 0,103 |
| COLON A1 k=[1] knn k=3 | 1,734633 | 0,103 |
| **COLON A1 k=[1] knn k=5** | **7,314081** | **2,55E-06** |
| **COLON A1 k=[1] random forest** | **7,314081** | **2,55E-06** |
| COLON A1 k=[2] knn k=1 | 1,459003 | 0,165 |
| COLON A1 k=[2] knn k=2 | 1,459003 | 0,165 |
| COLON A1 k=[2] knn k=3 | 1,459003 | 0,165 |
| **COLON A1 k=[2] knn k=5** | **6,053918** | **2,21E-05** |
| **COLON A1 k=[2] random forest** | **6,053918** | **2,21E-05** |
| **COLON A1 k=[1, 2] knn k=1** | **5,253814** | **0,000156** |
| **COLON A1 k=[1, 2] knn k=2** | **5,253814** | **0,000156** |
| **COLON A1 k=[1, 2] knn k=3** | **5,253814** | **0,000156** |
| **COLON A1 k=[1, 2] knn k=5** | **10,22664** | **1,39E-07** |
| **COLON A1 k=[1, 2] random forest** | **10,22664** | **1,39E-07** |
| COLON A10 k=[2] knn k=1 | 0,786319 | 0,445 |
| COLON A10 k=[2] knn k=2 | 0,786319 | 0,445 |
| COLON A10 k=[2] knn k=3 | 0,786319 | 0,445 |
| **COLON A10 k=[2] knn k=5** | **4,855304** | **0,000255** |
| **COLON A10 k=[2] random forest** | **4,855304** | **0,000255** |
| **COLON A10 k=[1, 2] knn k=1** | **2,77546** | **0,0216** |
| **COLON A10 k=[1, 2] knn k=2** | **2,77546** | **0,0216** |
| **COLON A10 k=[1, 2] knn k=3** | **2,77546** | **0,0216** |
| **COLON A10 k=[1, 2] knn k=5** | **5,849461** | **0,000244** |
| **COLON A10 k=[1, 2] random forest** | **5,849461** | **0,000244** |
| **COLON A2 k=[1, 2] knn k=1** | **5,58961** | **8,78E-05** |
| **COLON A2 k=[1, 2] knn k=2** | **5,58961** | **8,78E-05** |
| **COLON A2 k=[1, 2] knn k=3** | **5,58961** | **8,78E-05** |
| **COLON A2 k=[1, 2] knn k=5** | **11,37143** | **3,97E-08** |
| **COLON A2 k=[1, 2] random forest** | **11,37143** | **3,97E-08** |
| **COLON A5 k=[2] knn k=1** | **2,404602** | **0,0318** |
| **COLON A5 k=[2] knn k=2** | **2,404602** | **0,0318** |
| **COLON A5 k=[2] knn k=3** | **2,404602** | **0,0318** |
| **COLON A5 k=[2] knn k=5** | **6,778901** | **0,000013** |
| **COLON A5 k=[2] random forest** | **6,778901** | **0,000013** |
| COLON A6 k=[2] knn k=1 | 1,081534 | 0,308 |
| COLON A6 k=[2] knn k=2 | 1,081534 | 0,308 |
| COLON A6 k=[2] knn k=3 | 1,081534 | 0,308 |
| **COLON A6 k=[2] knn k=5** | **4,226342** | **0,00222** |
| **COLON A6 k=[2] random forest** | **4,226342** | **0,00222** |
| **COLON A7 k=[1, 2] knn k=1** | **5,212132** | **0,000168** |
| **COLON A7 k=[1, 2] knn k=2** | **5,212132** | **0,000168** |
| **COLON A7 k=[1, 2] knn k=3** | **5,212132** | **0,000168** |
| **COLON A7 k=[1, 2] knn k=5** | **10,64265** | **8,69E-08** |

| | | |
|---|---|---|
| **COLON A7 k=[1, 2] random forest** | **10,64265** | **8,69E-08** |
| **prostate A1 k=[1] knn k=1** | **10,1644** | **6,94E-09** |
| **prostate A1 k=[1] knn k=2** | **10,1644** | **6,94E-09** |
| **prostate A1 k=[1] knn k=3** | **10,1644** | **6,94E-09** |
| **prostate A1 k=[1] knn k=5** | **11,89539** | **5,81E-10** |
| **prostate A1 k=[1] random forest** | **-8,13463** | **1,93E-07** |
| **prostate A5 k=[2] knn k=1** | **6,550198** | **4,13E-05** |
| **prostate A5 k=[2] knn k=2** | **6,550198** | **4,13E-05** |
| **prostate A5 k=[2] knn k=3** | **6,550198** | **4,13E-05** |
| **prostate A5 k=[2] knn k=5** | **8,024943** | **6,34E-06** |
| **prostate A5 k=[2] random forest** | **-9,03996** | **2,01E-06** |
| **prostate A10 k=[2] knn k=1** | **6,174766** | **4,76E-05** |
| **prostate A10 k=[2] knn k=2** | **6,174766** | **4,76E-05** |
| **prostate A10 k=[2] knn k=3** | **6,174766** | **4,76E-05** |
| **prostate A10 k=[2] knn k=5** | **7,667817** | **5,79E-06** |
| **prostate A10 k=[2] random forest** | **-9,60891** | **5,5E-07** |
| **DLBCL A1 k=[1, 2] knn k=1** | **6,043671** | **3,02E-05** |
| DLBCL A1 k=[1, 2] knn k=2 | 1,495912 | 0,157 |
| **DLBCL A1 k=[1, 2] knn k=3** | **10,38471** | **5,84E-08** |
| **DLBCL A1 k=[1, 2] knn k=5** | **10,38471** | **5,84E-08** |
| **DLBCL A1 k=[1, 2] random forest** | **-7,18617** | **4,66E-06** |
| **DLBCL A2 k=[1, 2] knn k=1** | **4,04551** | **0,0012** |
| DLBCL A2 k=[1, 2] knn k=2 | 0,507355 | 0,62 |
| **DLBCL A2 k=[1, 2] knn k=3** | **7,422839** | **3,24E-06** |
| **DLBCL A2 k=[1, 2] knn k=5** | **7,422839** | **3,24E-06** |
| **DLBCL A2 k=[1, 2] random forest** | **-6,2473** | **2,14E-05** |
| **DLBCL A7 k=[1, 2] knn k=1** | **6,358447** | **1,77E-05** |
| DLBCL A7 k=[1, 2] knn k=2 | 1,66072 | 0,119 |
| **DLBCL A7 k=[1, 2] knn k=3** | **10,84264** | **3,4E-08** |
| **DLBCL A7 k=[1, 2] knn k=5** | **10,84264** | **3,4E-08** |
| **DLBCL A7 k=[1, 2] random forest** | **-7,30767** | **3,86E-06** |

*Table 5 The summary of statistical differences. Only the series when t-test or exact one sample Wilcoxon signed rank test was available are presented.*

| | kNN, k=1 | KNN, k=2 | kNN, k=3 | kNN, k=5 | Random Forest |
|---|---|---|---|---|---|
| AML ALL, A1, k=[1] | + | + | + | + | - |
| AML ALL, A1, k=[1, 2] | + | + | + | + | x |
| AML ALL, A2, k=[1, 2] | + | + | + | + | x |
| AML ALL, A5, k=[1, 2] | + | + | + | + | x |
| AML ALL, A7, k=[1, 2] | + | + | + | + | x |
| COLON, A1 k=[1] | x | x | x | + | + |
| COLON, A1 k=[2] | x | x | x | + | + |
| COLON, A1, k=[1, 2] | + | + | + | + | + |
| COLON, A2, k=[1, 2] | + | + | + | + | + |
| COLON, A5, k=[2] | + | + | + | + | + |
| COLON, A6, k=[2] | x | x | x | + | + |

| | | | | | |
|---|---|---|---|---|---|
| COLON, A7, k=[1, 2] | **+** | **+** | **+** | **+** | **+** |
| COLON, A10, k=[2] | x | x | x | **+** | **+** |
| COLON, A10, k=[1, 2] | **+** | **+** | **+** | **+** | **+** |
| OVARIAN, A5, k=[1] | x | x | x | x | - |
| OVARIAN, A5, k=[1, 2] | **+** | **+** | **+** | **+** | - |
| OVARIAN, A10, k=[1] | x | x | x | x | - |
| DLBCL, A1, k=[1, 2] | **+** | x | **+** | **+** | - |
| DLBCL, A2, k=[1, 2] | **+** | x | **+** | **+** | - |
| DLBCL, A7, k=[1, 2] | **+** | x | **+** | **+** | - |
| PROSTATE, A1, k=[1] | **+** | **+** | **+** | **+** | - |
| PROSTATE, A5, k=[2] | **+** | **+** | **+** | **+** | - |
| PROSTATE, A10, k=[2] | **+** | **+** | **+** | **+** | - |

Legend:

*-:* statistically significant difference in favour of SOTA model

***+:* statistically significant difference in favour of the proposed model**

*x:* no statistically significant difference