Table 1: Generalized Linear Models: Comparison of C2ST scores with a Random Forest (RF) and a Neural Network (NN). For the NN we follow the setup of Lueckmann et al., 2021. Evaluation across seven distinct scenarios on 50 synthetic and 17 real-world datasets. All results within two standard errors of the best average result in each scenario are marked in **bold**. Across all scenarios, both RF and NN classifiers yield quite consistent rankings of model performance with only insubstantial deviations in terms of the big picture. In particular, ICL is consistently among the top-performing approaches under both evaluation metrics. Out of the 14 total scenario–domain combinations (7 scenarios × 2 dataset types), the RF and NN metrics identify the same best-performing model in 12 cases.

| Scenario | Model | Synthetic Evaluation | | Real-World Evaluation | |
|---|---|---|---|---|---|
| | | C2ST RF ($\downarrow$) | C2ST NN ($\downarrow$) | C2ST RF ($\downarrow$) | C2ST NN ($\downarrow$) |
| Scenario 1 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.904 ($\pm$ 0.076) | 0.857 ($\pm$ 0.001) | 0.797 ($\pm$ 0.083) | 0.803 ($\pm$ 0.004) |
| | VI: MultivariateNormal | **0.750** ($\pm$ 0.128) | 0.780 ($\pm$ 0.002) | **0.607** ($\pm$ 0.070) | 0.713 ($\pm$ 0.004) |
| | VI: Structured Normal | **0.753** ($\pm$ 0.126) | 0.781 ($\pm$ 0.002) | **0.600** ($\pm$ 0.070) | 0.705 ($\pm$ 0.004) |
| | VI: IAF | **0.777** ($\pm$ 0.122) | 0.793 ($\pm$ 0.002) | 0.683 ($\pm$ 0.132) | 0.746 ($\pm$ 0.006) |
| | HMC | **0.745** ($\pm$ 0.130) | 0.777 ($\pm$ 0.002) | **0.595** ($\pm$ 0.075) | **0.702** ($\pm$ 0.004) |
| | **ICL (ours)** | **0.765** ($\pm$ 0.123) | **0.712** ($\pm$ 0.002) | **0.614** ($\pm$ 0.074) | **0.701** ($\pm$ 0.004) |
| Scenario 2 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.957 ($\pm$ 0.091) | 0.883 ($\pm$ 0.002) | 0.892 ($\pm$ 0.044) | 0.851 ($\pm$ 0.003) |
| | VI: MultivariateNormal | 0.910 ($\pm$ 0.131) | 0.860 ($\pm$ 0.002) | **0.820** ($\pm$ 0.031) | 0.815 ($\pm$ 0.003) |
| | VI: Structured Normal | 0.908 ($\pm$ 0.119) | 0.859 ($\pm$ 0.002) | **0.824** ($\pm$ 0.023) | 0.817 ($\pm$ 0.003) |
| | VI: IAF | 0.968 ($\pm$ 0.063) | 0.889 ($\pm$ 0.001) | 0.888 ($\pm$ 0.067) | 0.849 ($\pm$ 0.004) |
| | **ICL (ours)** | **0.839** ($\pm$ 0.072) | **0.824** ($\pm$ 0.001) | **0.768** ($\pm$ 0.033) | **0.789** ($\pm$ 0.003) |
| Scenario 3 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.866 ($\pm$ 0.101) | 0.838 ($\pm$ 0.002) | 0.797 ($\pm$ 0.083) | 0.803 ($\pm$ 0.004) |
| | VI: MultivariateNormal | 0.656 ($\pm$ 0.131) | 0.733 ($\pm$ 0.002) | **0.590** ($\pm$ 0.035) | 0.685 ($\pm$ 0.003) |
| | VI: Structured Normal | 0.653 ($\pm$ 0.125) | 0.731 ($\pm$ 0.002) | **0.582** ($\pm$ 0.028) | 0.681 ($\pm$ 0.003) |
| | VI: IAF | 0.751 ($\pm$ 0.148) | 0.780 ($\pm$ 0.002) | 0.673 ($\pm$ 0.141) | 0.741 ($\pm$ 0.006) |
| | **ICL (ours)** | **0.611** ($\pm$ 0.070) | **0.710** ($\pm$ 0.001) | **0.576** ($\pm$ 0.027) | **0.693** ($\pm$ 0.003) |
| Scenario 4 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.968 ($\pm$ 0.036) | 0.889 ($\pm$ 0.001) | 0.916 ($\pm$ 0.040) | 0.863 ($\pm$ 0.003) |
| | VI: MultivariateNormal | 0.855 ($\pm$ 0.123) | 0.832 ($\pm$ 0.002) | 0.771 ($\pm$ 0.017) | 0.790 ($\pm$ 0.002) |
| | VI: Structured Normal | 0.847 ($\pm$ 0.116) | 0.828 ($\pm$ 0.002) | 0.769 ($\pm$ 0.012) | 0.789 ($\pm$ 0.002) |
| | VI: IAF | 0.942 ($\pm$ 0.077) | 0.876 ($\pm$ 0.001) | 0.833 ($\pm$ 0.069) | 0.821 ($\pm$ 0.004) |
| | **ICL (ours)** | **0.753** ($\pm$ 0.049) | **0.781** ($\pm$ 0.001) | **0.762** ($\pm$ 0.015) | **0.786** ($\pm$ 0.002) |
| Scenario 5 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.866 ($\pm$ 0.085) | 0.838 ($\pm$ 0.002) | 0.810 ($\pm$ 0.036) | 0.810 ($\pm$ 0.003) |
| | VI: MultivariateNormal | 0.765 ($\pm$ 0.100) | 0.787 ($\pm$ 0.002) | 0.711 ($\pm$ 0.038) | 0.760 ($\pm$ 0.003) |
| | VI: Structured Normal | 0.758 ($\pm$ 0.098) | 0.784 ($\pm$ 0.002) | 0.705 ($\pm$ 0.032) | 0.757 ($\pm$ 0.003) |
| | VI: IAF | 0.814 ($\pm$ 0.105) | 0.812 ($\pm$ 0.002) | 0.777 ($\pm$ 0.106) | 0.793 ($\pm$ 0.005) |
| | **ICL (ours)** | **0.621** ($\pm$ 0.063) | **0.715** ($\pm$ 0.001) | **0.610** ($\pm$ 0.045) | **0.710** ($\pm$ 0.003) |
| Scenario 6 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.724 ($\pm$ 0.060) | 0.767 ($\pm$ 0.001) | 0.703 ($\pm$ 0.039) | 0.756 ($\pm$ 0.003) |
| | VI: MultivariateNormal | **0.534** ($\pm$ 0.018) | **0.672** ($\pm$ 0.001) | **0.538** ($\pm$ 0.019) | 0.674 ($\pm$ 0.002) |
| | VI: Structured Normal | **0.536** ($\pm$ 0.016) | **0.673** ($\pm$ 0.001) | **0.536** ($\pm$ 0.019) | 0.673 ($\pm$ 0.002) |
| | VI: IAF | 0.542 ($\pm$ 0.026) | 0.676 ($\pm$ 0.001) | **0.535** ($\pm$ 0.015) | 0.672 ($\pm$ 0.002) |
| | **ICL (ours)** | **0.532** ($\pm$ 0.019) | **0.671** ($\pm$ 0.001) | 0.556 ($\pm$ 0.017) | **0.653** ($\pm$ 0.002) |
| Scenario 7 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.998 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.938 ($\pm$ 0.074) | 0.874 ($\pm$ 0.001) | 0.936 ($\pm$ 0.024) | 0.873 ($\pm$ 0.003) |
| | VI: MultivariateNormal | 0.814 ($\pm$ 0.181) | 0.812 ($\pm$ 0.002) | **0.741** ($\pm$ 0.020) | 0.775 ($\pm$ 0.003) |
| | VI: Structured Normal | 0.824 ($\pm$ 0.177) | 0.817 ($\pm$ 0.002) | **0.734** ($\pm$ 0.025) | 0.772 ($\pm$ 0.003) |
| | VI: IAF | 0.939 ($\pm$ 0.091) | 0.874 ($\pm$ 0.002) | 0.864 ($\pm$ 0.093) | 0.837 ($\pm$ 0.005) |
| | **ICL (ours)** | **0.700** ($\pm$ 0.116) | **0.721** ($\pm$ 0.002) | 0.773 ($\pm$ 0.048) | **0.751** ($\pm$ 0.003) |

Table 2: Factor Analysis: Comparison of C2ST scores using a Random Forest (RF) and a Neural Network (NN) classifier across six different scenarios on 50 synthetic and 17 real-world datasets. For the NN we follow the setup of Lueckmann et al., 2021. All results within two standard errors of the best average result in each scenario are marked in **bold**. Across both classifiers and dataset types, the ICL method consistently achieves top-tier performance, often outperforming standard variational approximations and the Laplace baseline. The RF and NN classifiers yield largely consistent rankings, identifying ICL as the best-performing model in the majority of cases.

| Scenario | Model | Synthetic Evaluation | | Real-World Evaluation | |
|---|---|---|---|---|---|
| | | C2ST RF ($\downarrow$) | C2ST NN ($\downarrow$) | C2ST RF ($\downarrow$) | C2ST NN ($\downarrow$) |
| Scenario 1 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 1.000 ($\pm$ 0.001) | 0.997 ($\pm$ 0.000) | 0.979 ($\pm$ 0.008) | 0.950 ($\pm$ 0.001) |
| | VI: MultivariateNormal | 0.998 ($\pm$ 0.003) | 0.960 ($\pm$ 0.000) | 0.966 ($\pm$ 0.010) | 0.944 ($\pm$ 0.001) |
| | VI: Structured Normal | 0.997 ($\pm$ 0.004) | 0.959 ($\pm$ 0.000) | 0.979 ($\pm$ 0.010) | 0.950 ($\pm$ 0.001) |
| | VI: IAF | 0.953 ($\pm$ 0.104) | 0.937 ($\pm$ 0.001) | 0.849 ($\pm$ 0.075) | 0.885 ($\pm$ 0.003) |
| | **ICL (ours)** | **0.552** ($\pm$ 0.028) | **0.737** ($\pm$ 0.000) | **0.606** ($\pm$ 0.038) | **0.764** ($\pm$ 0.001) |
| Scenario 2 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.998 ($\pm$ 0.002) | 0.960 ($\pm$ 0.000) | 0.975 ($\pm$ 0.010) | 0.948 ($\pm$ 0.001) |
| | VI: MultivariateNormal | 0.989 ($\pm$ 0.009) | 0.955 ($\pm$ 0.000) | 0.951 ($\pm$ 0.025) | 0.936 ($\pm$ 0.001) |
| | VI: Structured Normal | 0.984 ($\pm$ 0.031) | 0.953 ($\pm$ 0.000) | 0.958 ($\pm$ 0.025) | 0.940 ($\pm$ 0.001) |
| | VI: IAF | 0.966 ($\pm$ 0.066) | 0.944 ($\pm$ 0.001) | 0.799 ($\pm$ 0.058) | 0.860 ($\pm$ 0.002) |
| | **ICL (ours)** | **0.542** ($\pm$ 0.006) | **0.732** ($\pm$ 0.000) | **0.622** ($\pm$ 0.032) | **0.772** ($\pm$ 0.001) |
| Scenario 3 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.999 ($\pm$ 0.002) | 0.960 ($\pm$ 0.000) | 0.951 ($\pm$ 0.007) | 0.936 ($\pm$ 0.001) |
| | VI: MultivariateNormal | 0.994 ($\pm$ 0.007) | 0.958 ($\pm$ 0.000) | 0.945 ($\pm$ 0.007) | 0.933 ($\pm$ 0.001) |
| | VI: Structured Normal | 0.997 ($\pm$ 0.003) | 0.959 ($\pm$ 0.000) | 0.942 ($\pm$ 0.009) | 0.932 ($\pm$ 0.001) |
| | VI: IAF | 0.990 ($\pm$ 0.011) | 0.987 ($\pm$ 0.000) | 0.928 ($\pm$ 0.015) | 0.925 ($\pm$ 0.001) |
| | **ICL (ours)** | **0.537** ($\pm$ 0.023) | **0.729** ($\pm$ 0.000) | **0.609** ($\pm$ 0.019) | **0.765** ($\pm$ 0.001) |
| Scenario 4 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) | 0.977 ($\pm$ 0.003) | 0.949 ($\pm$ 0.000) |
| | VI: MultivariateNormal | 0.999 ($\pm$ 0.001) | 0.960 ($\pm$ 0.000) | 0.973 ($\pm$ 0.008) | 0.947 ($\pm$ 0.001) |
| | VI: Structured Normal | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) | 0.973 ($\pm$ 0.007) | 0.947 ($\pm$ 0.001) |
| | VI: IAF | 0.999 ($\pm$ 0.001) | 0.960 ($\pm$ 0.000) | **0.961** ($\pm$ 0.018) | 0.941 ($\pm$ 0.001) |
| | **ICL (ours)** | **0.684** ($\pm$ 0.060) | **0.803** ($\pm$ 0.001) | 0.988 ($\pm$ 0.003) | **0.955** ($\pm$ 0.000) |
| Scenario 5 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.999 ($\pm$ 0.002) | 0.960 ($\pm$ 0.000) | 0.944 ($\pm$ 0.010) | 0.933 ($\pm$ 0.001) |
| | VI: MultivariateNormal | 0.995 ($\pm$ 0.007) | 0.958 ($\pm$ 0.000) | 0.930 ($\pm$ 0.017) | 0.926 ($\pm$ 0.001) |
| | VI: Structured Normal | 0.998 ($\pm$ 0.005) | 0.960 ($\pm$ 0.000) | 0.934 ($\pm$ 0.011) | 0.928 ($\pm$ 0.001) |
| | VI: IAF | 0.992 ($\pm$ 0.012) | 0.957 ($\pm$ 0.000) | 0.910 ($\pm$ 0.011) | 0.916 ($\pm$ 0.001) |
| | **ICL (ours)** | **0.535** ($\pm$ 0.016) | **0.728** ($\pm$ 0.000) | **0.886** ($\pm$ 0.017) | **0.904** ($\pm$ 0.001) |
| Scenario 6 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.998 ($\pm$ 0.002) | 0.960 ($\pm$ 0.000) | 0.949 ($\pm$ 0.008) | 0.935 ($\pm$ 0.001) |
| | VI: MultivariateNormal | 0.991 ($\pm$ 0.013) | 0.956 ($\pm$ 0.000) | 0.938 ($\pm$ 0.009) | 0.930 ($\pm$ 0.001) |
| | VI: Structured Normal | 0.997 ($\pm$ 0.005) | 0.959 ($\pm$ 0.000) | 0.944 ($\pm$ 0.006) | 0.933 ($\pm$ 0.001) |
| | VI: IAF | 0.989 ($\pm$ 0.029) | 0.955 ($\pm$ 0.000) | 0.865 ($\pm$ 0.027) | 0.893 ($\pm$ 0.001) |
| | **ICL (ours)** | **0.543** ($\pm$ 0.021) | **0.732** ($\pm$ 0.000) | **0.666** ($\pm$ 0.020) | **0.794** ($\pm$ 0.001) |

Table 3: Gaussian Mixture Models: Comparison of C2ST scores using a Random Forest (RF) and a Neural Network (NN) classifier across six distinct scenarios on 50 synthetic and 17 real-world datasets. All results within two standard errors of the best average result in each scenario are marked in **bold**. For the NN we follow the setup of Lueckmann et al., 2021. Both RF and NN classifiers yield consistent rankings, with ICL emerging as the top method in scenarios with more pronounced model mismatch.

| Scenario | Model | Synthetic Evaluation | | Real-World Evaluation | |
|---|---|---|---|---|---|
| | | C2ST RF ($\downarrow$) | C2ST NN ($\downarrow$) | C2ST RF ($\downarrow$) | C2ST NN ($\downarrow$) |
| Scenario 1 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.988 ($\pm$ 0.013) | 1.012 ($\pm$ 0.000) | 0.995 ($\pm$ 0.006) | 0.996 ($\pm$ 0.001) |
| | VI: MultivariateNormal | 0.988 ($\pm$ 0.013) | 1.012 ($\pm$ 0.000) | 0.994 ($\pm$ 0.007) | 0.993 ($\pm$ 0.001) |
| | VI: Structured Normal | 0.987 ($\pm$ 0.015) | 0.982 ($\pm$ 0.000) | 0.993 ($\pm$ 0.009) | 0.992 ($\pm$ 0.001) |
| | VI: IAF | 0.989 ($\pm$ 0.013) | 0.983 ($\pm$ 0.000) | 0.995 ($\pm$ 0.010) | 0.996 ($\pm$ 0.001) |
| | **ICL (ours)** | **0.760** ($\pm$ 0.092) | **0.825** ($\pm$ 0.001) | **0.847** ($\pm$ 0.082) | **0.869** ($\pm$ 0.003) |
| Scenario 2 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | VI: DiagonalNormal | 0.989 ($\pm$ 0.024) | 0.983 ($\pm$ 0.000) | 0.998 ($\pm$ 0.003) | 0.997 ($\pm$ 0.001) |
| | VI: MultivariateNormal | 0.991 ($\pm$ 0.021) | 0.991 ($\pm$ 0.000) | 0.999 ($\pm$ 0.002) | 1.002 ($\pm$ 0.001) |
| | VI: Structured Normal | 0.992 ($\pm$ 0.017) | 0.988 ($\pm$ 0.000) | 0.999 ($\pm$ 0.002) | 1.002 ($\pm$ 0.001) |
| | VI: IAF | 0.992 ($\pm$ 0.021) | 0.988 ($\pm$ 0.000) | 0.998 ($\pm$ 0.004) | 0.997 ($\pm$ 0.001) |
| | **ICL (ours)** | **0.812** ($\pm$ 0.061) | **0.851** ($\pm$ 0.001) | **0.937** ($\pm$ 0.041) | **0.915** ($\pm$ 0.002) |
| Scenario 3 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | VI: DiagonalNormal | **0.996** ($\pm$ 0.011) | 1.004 ($\pm$ 0.000) | **0.992** ($\pm$ 0.018) | 0.988 ($\pm$ 0.001) |
| | VI: MultivariateNormal | 0.997 ($\pm$ 0.009) | 1.007 ($\pm$ 0.000) | **0.993** ($\pm$ 0.016) | 0.992 ($\pm$ 0.001) |
| | VI: Structured Normal | **0.995** ($\pm$ 0.017) | 0.996 ($\pm$ 0.000) | **0.993** ($\pm$ 0.016) | 0.992 ($\pm$ 0.001) |
| | VI: IAF | **0.994** ($\pm$ 0.018) | 0.993 ($\pm$ 0.000) | **0.993** ($\pm$ 0.017) | 0.992 ($\pm$ 0.001) |
| | **ICL (ours)** | 1.000 ($\pm$ 0.000) | **0.997** ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | **0.997** ($\pm$ 0.000) |
| Scenario 4 | Laplace Approximation | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | VI: DiagonalNormal | **1.000** ($\pm$ 0.002) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | VI: MultivariateNormal | **1.000** ($\pm$ 0.002) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | VI: Structured Normal | **1.000** ($\pm$ 0.001) | 0.997 ($\pm$ 0.000) | **0.996** ($\pm$ 0.016) | 1.004 ($\pm$ 0.001) |
| | VI: IAF | **1.000** ($\pm$ 0.002) | 0.997 ($\pm$ 0.000) | 1.000 ($\pm$ 0.000) | 0.997 ($\pm$ 0.000) |
| | **ICL (ours)** | 1.000 ($\pm$ 0.000) | **0.997** ($\pm$ 0.000) | **1.000** ($\pm$ 0.000) | **0.997** ($\pm$ 0.000) |