# Dissimilarités entre jeux de données

-

# Listings

15 février 2017

## 1   Méta-attributs généraux des jeux de données

TABLE 1 – Méta-attributs simples des jeux de données

| Méta-attribut | Description | $\delta^{\emptyset}(x)$ |
|---|---|---|
| DefaultAccuracy | The predictive accuracy obtained by predicting the majority class. | 0 |
| Dimensionality | Number of attributes divided by the number of instances. | 0 |
| MajorityClassPercentage | Percentage of rows with the class with the most assigned index. | $|x - 1|$ |
| MajorityClassSize | The number of instances that have the majority class. | 0 |
| MinorityClassPercentage | Percentage of rows with the class with the least assigned index. | $|x - 1|$ |
| MinorityClassSize | The number of instances that have the minority class. | 0 |
| NumberOfBinaryFeatures | Count of binary attributes. | 0 |
| NumberOfClasses | The number of classes in the class attribute. | 0 |
| NumberOfFeatures | Number of attributes (columns) of the dataset. | 0 |
| NumberOfInstances | Number of instances (rows) of the dataset. | 0 |
| NumberOfInstancesWithMissingValues | Number of instances with at least one value missing. | 0 |
| NumberOfMissingValues | Number of missing values in the dataset. | 0 |
| NumberOfNumericFeatures | Count of categorical attributes. | 0 |
| NumberOfSymbolicFeatures | Count of nominal attributes. | 0 |
| PercentageOfBinaryFeatures | Percentage of binary attributes. | 0 |
| PercentageOfInstancesWithMissingValues | Percentage of instances with missing values. | $|x - 1|$ |
| PercentageOfMissingValues | Percentage of missing values. | $|x - 1|$ |
| PercentageOfNumericFeatures | Percentage of numerical attributes. | 0 |
| PercentageOfSymbolicFeatures | Percentage of nominal attributes. | 0 |

TABLE 2 – Méta-attributs généraux décrivant les attributs numériques

| Méta-attribut | Description | $\delta^{\emptyset}(x)$ |
|---|---|---|
| MeanMeansOfNumericAtts | Mean of means among numeric attributes. | $0$ |
| MeanStdDevOfNumericAtts | Mean standard deviation of numeric attributes. | $|x|$ |
| MeanKurtosisOfNumericAtts | Mean kurtosis among numeric attributes. | $|x+1,2|$ |
| MeanSkewnessOfNumericAtts | Mean skewness among numeric attributes. | $|x|$ |
| MinMeansOfNumericAtts | Min of means among numeric attributes. | $0$ |
| MinStdDevOfNumericAtts | Min standard deviation of numeric attributes. | $|x|$ |
| MinKurtosisOfNumericAtts | Min kurtosis among numeric attributes. | $|x+1,2|$ |
| MinSkewnessOfNumericAtts | Min skewness among numeric attributes. | $|x|$ |
| MaxMeansOfNumericAtts | Max of means among numeric attributes. | $0$ |
| MaxStdDevOfNumericAtts | Max standard deviation of numeric attributes. | $|x|$ |
| MaxKurtosisOfNumericAtts | Max kurtosis among numeric attributes. | $|x+1,2|$ |
| MaxSkewnessOfNumericAtts | Max skewness among numeric attributes. | $|x|$ |
| Quartile1MeansOfNumericAtts | First quartile of means among numeric attributes. | $0$ |
| Quartile1StdDevOfNumericAtts | First quartile of standard deviation of numeric attributes. | $|x|$ |
| Quartile1KurtosisOfNumericAtts | First quartile of kurtosis among numeric attributes. | $|x+1,2|$ |
| Quartile1SkewnessOfNumericAtts | First quartile of skewness among numeric attributes. | $|x|$ |
| Quartile2MeansOfNumericAtts | Second quartile of means among numeric attributes. | $0$ |
| Quartile2StdDevOfNumericAtts | Second quartile of standard deviation of numeric attributes. | $|x|$ |
| Quartile2KurtosisOfNumericAtts | Second quartile of kurtosis among numeric attributes. | $|x+1,2|$ |
| Quartile2SkewnessOfNumericAtts | Second quartile of skewness among numeric attributes. | $|x|$ |
| Quartile3MeansOfNumericAtts | Third quartile of means among numeric attributes. | $0$ |
| Quartile3StdDevOfNumericAtts | Third quartile of standard deviation of numeric attributes. | $|x|$ |
| Quartile3KurtosisOfNumericAtts | Third quartile of kurtosis among numeric attributes. | $|x+1,2|$ |
| Quartile3SkewnessOfNumericAtts | Third quartile of skewness among numeric attributes. | $|x|$ |

TABLE 3 – Méta-attributs généraux décrivant les attributs nominaux

| Méta-attribut | Description | $\delta^{\emptyset}(x)$ |
|---|---|---|
| ClassEntropy | Entropy of the target attribute. | 0 |
| EquivalentNumberOfAtts | An estimate of the amount of useful attributes. | 0 |
| NoiseToSignalRatio | An estimate of the amount of non-useful information in the attributes regarding the class. | 0 |
| MeanAttributeEntropy | Mean of entropy among attributes. | $|x|$ |
| MeanMutualInformation | Mean of mutual information between the nominal attributes and the target attribute. | $|x|$ |
| MinAttributeEntropy | Min of entropy among attributes. | $|x|$ |
| MinMutualInformation | Min of mutual information between the nominal attributes and the target attribute. | $|x|$ |
| MaxAttributeEntropy | Max of entropy among attributes. | $|x|$ |
| MaxMutualInformation | Max of mutual information between the nominal attributes and the target attribute. | $|x|$ |
| Quartile1AttributeEntropy | First quartile of entropy among attributes. | $|x|$ |
| Quartile1MutualInformation | First quartile of mutual information between the nominal attributes and the target attribute. | $|x|$ |
| Quartile2AttributeEntropy | Second quartile of entropy among attributes. | $|x|$ |
| Quartile2MutualInformation | Second quartile of mutual information between the nominal attributes and the target attribute. | $|x|$ |
| Quartile3AttributeEntropy | Third quartile of entropy among attributes. | $|x|$ |
| Quartile3MutualInformation | Third quartile of mutual information between the nominal attributes and the target attribute. | $|x|$ |
| MaxNominalAttDistinctValues | The maximum number of distinct values among attributes of the nominal type. | $|x|$ |
| MinNominalAttDistinctValues | The minimal number of distinct values among attributes of the nominal type. | $|x|$ |
| MeanNominalAttDistinctValues | Mean of number of distinct values among the attributes of the nominal type. | $|x|$ |
| StdvNominalAttDistinctValues | Standard deviation of the number of distinct values among nominal attributes. | $|x|$ |

TABLE 4 – "Aire sous la courbe" des landmarkers

| Méta-attribut | Description | $\delta^{\emptyset}(x)$ |
|---|---|---|
| DecisionStumpAUC | Area Under ROC achieved by the landmarker weka.classifiers.trees.DecisionStump | $|x - 0,5|$ |
| J48AUC | Area Under ROC achieved by the landmarker weka.classifiers.trees.J48 | $|x - 0,5|$ |
| JRipAUC | Area Under ROC achieved by the landmarker weka.classifiers.rules.Jrip | $|x - 0,5|$ |
| kNN_3NAUC | Area Under ROC achieved by the landmarker weka.classifiers.lazy.IBk -K 3 | $|x - 0,5|$ |
| NaiveBayesAUC | Area Under ROC achieved by the landmarker weka.classifiers.bayes.NaiveBayes | $|x - 0,5|$ |
| NBTreeAUC | Area Under ROC achieved by the landmarker weka.classifiers.trees.NBTree | $|x - 0,5|$ |
| RandomTreeDepth3AUC | Area Under ROC achieved by the landmarker weka.classifiers.trees.RandomTree -depth 3 | $|x - 0,5|$ |
| REPTreeDepth3AUC | Area Under ROC achieved by the landmarker weka.classifiers.trees.REPTree -L 3 | $|x - 0,5|$ |
| SimpleLogisticAUC | Area Under ROC achieved by the landmarker weka.classifiers.functions.SimpleLogistic | $|x - 0,5|$ |

TABLE 5 – Taux d'erreur des landmarkers

| Méta-attribut | Description | $\delta^{\emptyset}(x)$ |
|---|---|---|
| DecisionStumpErrRate | Error rate achieved by the landmarker weka.classifiers.trees.DecisionStump | $\lvert x - 1 \rvert$ |
| J48ErrRate | Error rate achieved by the landmarker weka.classifiers.trees.J48 | $\lvert x - 1 \rvert$ |
| JRipErrRate | Error rate achieved by the landmarker weka.classifiers.rules.Jrip | $\lvert x - 1 \rvert$ |
| kNN_3NErrRate | Error rate achieved by the landmarker weka.classifiers.lazy.IBk -K 3 | $\lvert x - 1 \rvert$ |
| NaiveBayesErrRate | Error rate achieved by the landmarker weka.classifiers.bayes.NaiveBayes | $\lvert x - 1 \rvert$ |
| NBTreeErrRate | Error rate achieved by the landmarker weka.classifiers.trees.NBTree | $\lvert x - 1 \rvert$ |
| RandomTreeDepth3ErrRate | Error rate achieved by the landmarker weka.classifiers.trees.RandomTree -depth 3 | $\lvert x - 1 \rvert$ |
| REPTreeDepth3ErrRate | Error rate achieved by the landmarker weka.classifiers.trees.REPTree -L 3 | $\lvert x - 1 \rvert$ |
| SimpleLogisticErrRate | Error rate achieved by the landmarker weka.classifiers.functions.SimpleLogistic | $\lvert x - 1 \rvert$ |

TABLE 6 – Kappa de Cohen des landmarkers

| Méta-attribut | Description | $\delta^{\emptyset}(x)$ |
|---|---|---|
| DecisionStumpKappa | Kappa coefficient achieved by the landmarker weka.classifiers.trees.DecisionStump | $\lvert x \rvert$ |
| J48Kappa | Kappa coefficient achieved by the landmarker weka.classifiers.trees.J48 | $\lvert x \rvert$ |
| JRipKappa | Kappa coefficient achieved by the landmarker weka.classifiers.rules.Jrip | $\lvert x \rvert$ |
| kNN_3NKappa | Kappa coefficient achieved by the landmarker weka.classifiers.lazy.IBk -K 3 | $\lvert x \rvert$ |
| NaiveBayesKappa | Kappa coefficient achieved by the landmarker weka.classifiers.bayes.NaiveBayes | $\lvert x \rvert$ |
| NBTreeKappa | Kappa coefficient achieved by the landmarker weka.classifiers.trees.NBTree | $\lvert x \rvert$ |
| RandomTreeDepth3Kappa | Kappa coefficient achieved by the landmarker weka.classifiers.trees.RandomTree -depth 3 | $\lvert x \rvert$ |
| REPTreeDepth3Kappa | Kappa coefficient achieved by the landmarker weka.classifiers.trees.REPTree -L 3 | $\lvert x \rvert$ |
| SimpleLogisticKappa | Kappa coefficient achieved by the landmarker weka.classifiers.functions.SimpleLogistic | $\lvert x \rvert$ |

# 2 Méta-attributs des attributs

TABLE 7 – Méta-attributs simples des attributs

| Méta-attribut | Description | $\delta^{\emptyset}(x)$ |
|---|---|---|
| ValuesCount | Number of values. | $|x-1|$ |
| NonMissingValuesCount | Number of non missing values. | $|x|$ |
| MissingValuesCount | Number of missing values. | $0$ |
| Distinct | Number of distinct values. | $|x-1|$ |
| AverageClassCount | Average count of occurrences among different classes. | $0$ |
| Entropy | Entropy of the values. | $|x-1|$ |
| MostFequentClassCount | Count of the most probable class. | $0$ |
| LeastFequentClassCount | Count of the least probable class. | $0$ |
| ModeClassCount | Mode of the number of distinct values. | $0$ |
| MedianClassCount | Median of the number of distinct values. | $0$ |
| PearsonCorrellationCoefficient | Pearson Correlation Coeeficient of the values with the target attribute. | $|x|$ |
| SpearmanCorrelationCoefficient | Spearman Correlation Coeeficient of the values with the target attribute. | $|x|$ |
| CovarianceWithTarget | Covariance of the values with the target attribute. | $|x|$ |

TABLE 8 – Méta-attributs spécifiques aux attributs nominaux

| Méta-attribut | Description | $\delta^{\emptyset}(x)$ |
|---|---|---|
| UniformDiscrete | Result of Pearson's chi-squared test for discrete uniform distribution. | $|x-1|$ |
| ChiSquareUniformDistribution | Statistic value for the Pearson's chi-squared test. | $|x|$ |
| RationOfDistinguishing CategoriesByKolmogorov SmirnoffSlashChiSquare | Ratio of attribute values that after sub-setting the dataset to that attribute value leads to different distribution of the target as indicated by the Kolmogorov-Smirnoff test. | $0$ |
| RationOfDistinguishing CategoriesByUtest | Ratio of attribute values that after sub-setting the dataset to that attribute value leads to different distribution of the target as indicated by the Mann-Whitney U-test. | $0$ |

TABLE 9 – Méta-attributs spécifiques aux attributs numériques

| Méta-attribut | Description | $\delta^{\emptyset}(x)$ |
|---|---|---|
| IsUniform | Whether statistical test did or not reject that the attribute values corresponds to a uniform distribution. | $\|x-1\|$ |
| IntegersOnly | Whether attribute values are only integers. | $0$ |
| Min | Minimal value of the attribute values. | $0$ |
| Max | Maximal value of the attribute values. | $0$ |
| Kurtosis | Kurtosis of the values. | $\|x+1,2\|$ |
| Mean | Mean of the values. | $0$ |
| Skewness | Skewness of the values. | $\|x\|$ |
| StandardDeviation | Standard deviation of the values. | $\|x\|$ |
| Variance | Variance of the values. | $\|x\|$ |
| Mode | Mode of the values. | $0$ |
| Median | Median of the values. | $0$ |
| ValueRange | Difference between maximum and minimum of the values. | $\|x\|$ |
| LowerOuterFence | Lower outer fence of the values. | $0$ |
| HigherOuterFence | Higher outer fence of the values. | $0$ |
| LowerQuartile | Lower quartile of the values. | $0$ |
| HigherQuartile | Higher quartile of the values. | $0$ |
| HigherConfidence | Higher confidence interval of the values. | $0$ |
| LowerConfidence | Lower confidence interval of the values. | $0$ |
| PositiveCount | Number of positive values. | $\|x\|$ |
| NegativeCount | Number of negative values. | $\|x\|$ |

TABLE 10 – Méta-attributs normalisés selon le nombre d'instances

| Méta-attribut | Description | $\delta^{\emptyset}(x)$ |
|---|---|---|
| MissingValues | 1 if count of missing values is greater than 0, 0 otherwise. | $\|x-1\|$ |
| AveragePercentageOfClass | Percentage of the occurrences among classes. | $\|x-1\|$ |
| PercentageOfMissing | Percentage of missing values in the attribute. | $\|x-1\|$ |
| PercentageOfNonMissing | Percentage of non missing values in the attribute. | $\|x\|$ |
| PercentageOfMostFrequentClass | Percentage of the most frequent class. | $\|x-1\|$ |
| PercentageOfLeastFrequentClass | Percentage of the least frequent class. | $\|x-1\|$ |
| ModeClassPercentage | Percentage of mode of class count calculated as ModeFrequentClassCount / ValuesCount. | $\|x-1\|$ |
| MedianClassPercentage | Percentage of median of class count calculated as MedianFrequentClassCount / ValuesCount. | $\|x-1\|$ |

TABLE 11 – Méta-attributs normalisés spécifiques aux attributs numériques

| Méta-attribut | Description | $\delta^{\emptyset}(x)$ |
|---|---|---|
| PositivePercentage | Percentage of positive values. | $\|x\|$ |
| NegativePercentage | Percentage of negative values. | $\|x\|$ |
| HasPositiveValues | 1 if attribute has at least one positive value, 0 otherwise. | $\|x\|$ |
| HasNegativeValues | 1 if attribute has at least one negative value, 0 otherwise. | $\|x\|$ |

# 3   Algorithmes de la *baseline*

Table 12 – Algorithmes d'apprentissage traditionnels de la *baseline*

| Implémentation Weka | Description |
| --- | --- |
| **GaussianProcesses** | Gaussian Processes for regression. See [**?**]. |
| **LinearRegression** | Linear regression for prediction. Uses the Akaike criterion for model selection, and is able to deal with weighted instances. See [**?**]. |
| **RBFRegressor** | Radial basis function networks, trained in a fully supervised manner by minimizing squared error with the BFGS method. See [**?**]. |
| **SMOreg** | Sequential minimal optimization algorithm for training a support vector regression model. See [**?**]. |
| **RandomForest** | Ensemble of decision trees outputting the mean prediction of the individual trees. See [**?**]. |
| **KStar** | Instance-based classifier where the class of a test instance is based upon the class of those training instances similar to it, as determined by an entropy-based distance function. See [**?**]. |
| **M5Rules** | Generates a decision list for regression problems using separate-and-conquer. Builds a model tree using M5 In each iteration and makes the best leaf into a rule. See [**?**]. |

# 4   Meta-Dataset

Table 13 – Chaines de traitement weka employés comme classifieurs dans les expériences comparatives

| | |
| --- | --- |
| A1DE | LMT |
| AdaBoostM1_DecisionStump | Logistic |
| Bagging_REPTree | LogitBoost_DecisionStump |
| BayesNet_K2 | LWL_DecisionStump |
| BFTree | MultiBoostAB_DecisionStump |
| ConjunctiveRule | NaiveBayes |
| Dagging_DecisionStump | NBTree |
| DecisionStump | OLM |
| DecisionTable_BestFirst | OneR |
| END_ND_J48 | PART |
| FT | RacedIncrementalLogitBoost_DecisionStump |
| FURIA | RandomForest |
| Grading_ZeroR | RandomSubSpace_REPTree |
| HoeffdingTree | RandomTree |
| HyperPipes | RBFClassifier |
| IB1 | REPTree |
| IBk | Ridor |
| J48 | RotationForest_PrincipalComponents_J48 |
| J48graft | SimpleCart |
| JRip | SimpleLogistic |
| KStar | SMO_PolyKernel |
| LADTree | SMO_RBFKernel |
| LibLINEAR | VFI |
| LibSVM | ZeroR |

TABLE 14 – Jeux de données OpenML utilisés dans les expériences comparatives

| | | | |
|---|---|---|---|
| anneal | rmftsa_sleepdata | diggle_table_a2 | socmob |
| anneal | sleuth_ex2016 | delta_elevators | fri_c1_250_10 |
| kr-vs-kp | sleuth_ex2015 | chatfield_4 | fri_c3_500_10 |
| labor | visualizing_livestock | house_16H | fri_c3_500_50 |
| audiology | diggle_table_a2 | cal_housing | lowbwt |
| autos | fruitfly | houses | fri_c0_500_10 |
| lymph | fri_c3_1000_25 | fri_c1_500_10 | echoMonths |
| balance-scale | fri_c3_100_50 | boston_corrected | kidney |
| breast-cancer | rmftsa_ladata | sensory | visualizing_ethanol |
| mfeat-fourier | veteran | disclosure_x_noise | arsenic-male-bladder |
| breast-w | abalone | fri_c1_100_5 | quake |
| mfeat-karhunen | pwLinear | fri_c2_250_10 | arsenic-female-bladder |
| mfeat-morphological | pol | autoMpg | arsenic-female-lung |
| car | fri_c4_1000_25 | fri_c3_250_25 | arsenic-male-lung |
| mfeat-zernike | analcatdata_vineyard | bank32nh | tae |
| cmc | bank8FM | analcatdata_vehicle | braziltourism |
| mushroom | fri_c2_100_5 | sleuth_case1102 | segment |
| nursery | 2dplanes | fri_c1_1000_50 | nursery |
| colic | analcatdata_supreme | fri_c4_500_25 | postoperative-patient-data |
| optdigits | visualizing_slope | kdd_el_nino-small | analcatdata_broadwaymult |
| credit-a | fri_c1_250_5 | autoHorse | mfeat-morphological |
| page-blocks | baskball | stock | heart-h |
| credit-g | fri_c0_250_50 | analcatdata_runshoes | pasture |
| pendigits | machine_cpu | house_8L | cars |
| postoperative-patient-data | ailerons | breastTumor | analcatdata_birthday |
| dermatology | cpu_small | fri_c0_1000_10 | iris |
| segment | visualizing_environmental | elevators | analcatdata_authorship |
| diabetes | space_ga | wind | mfeat-fourier |
| ecoli | sleep | schlvote | squash-stored |
| sonar | fri_c3_1000_10 | fri_c0_1000_25 | wine |
| glass | rmftsa_sleepdata | fri_c0_1000_50 | hayes-roth |
| soybean | fri_c1_1000_5 | analcatdata_gsssexsurvey | autos |
| haberman | fri_c3_250_5 | housing | kdd_JapaneseVowels |
| spambase | auto_price | fishcatch | letter |
| tae | fri_c1_250_25 | fri_c4_500_10 | waveform-5000 |
| heart-c | servo | bolts | optdigits |
| tic-tac-toe | analcatdata_wildcat | hungarian | kdd_internet_usage |
| heart-h | fri_c3_500_5 | vinnie | heart-c |
| heart-statlog | pm10 | auto93 | cmc |
| vehicle | fri_c4_1000_10 | sleuth_ex2016 | squash-unstored |
| hepatitis | puma32H | fri_c4_250_10 | analcatdata_marketing |
| vote | wisconsin | sleuth_ex2015 | fl2000 |
| ionosphere | fri_c0_100_5 | fri_c2_1000_50 | anneal |
| waveform-5000 | sleuth_ex1605 | visualizing_livestock | eucalyptus |
| iris | autoPrice | fri_c4_100_25 | car |
| BNG(cmc,nominal,55296) | meta | fri_c2_500_10 | kdd_ipums_la-97-small |
| electricity | cpu_act | fri_c1_500_5 | vehicle |
| primary-tumor | fri_c2_100_10 | pollen | mfeat-zernike |
| solar-flare | fri_c0_250_10 | boston | prnn_fglass |
| solar-flare | analcatdata_apnea3 | fri_c3_250_50 | balance-scale |
| adult | analcatdata_apnea2 | rabe_131 | audiology |
| yeast | fri_c1_500_50 | analcatdata_chlamydia | hypothyroid |
| satimage | analcatdata_apnea1 | fri_c1_100_50 | kdd_ipums_la-98-small |
| abalone | fri_c3_100_25 | fri_c2_250_50 | primary-tumor |
| braziltourism | fri_c1_250_50 | fri_c4_100_10 | glass |
| eucalyptus | strikes | fri_c2_500_25 | lymph |
| BNG(breast-w) | quake | mu284 | white-clover |
| meta_all | fri_c0_250_25 | mv | dermatology |
| meta_batchincremental | disclosure_x_bias | pollution | ecoli |
| meta_ensembles | fri_c2_100_25 | fri_c0_500_5 | analcatdata_challenger |
| meta_instanceincremental | fri_c0_250_5 | transplant | analcatdata_dmft |
| lung-cancer | sleuth_ex1714 | no2 | confidence |
| wine | bodyfat | mbagrade | kdd_ipums_la-99-small |
| hypothyroid | fri_c1_500_25 | fri_c0_500_50 | pendigits |
| shuttle-landing-control | rabe_265 | fri_c0_100_25 | mfeat-karhunen |
| Australian | rabe_266 | cloud | page-blocks |
| sick | fri_c3_100_10 | sleuth_case1202 | soybean |
| monks-problems-1 | newton_hema | visualizing_hamster | analcatdata_germangss |
| monks-problems-2 | wind_correlations | rabe_148 | grub-damage |
| monks-problems-3 | triazines | chscase_geyser1 | ada_prior |
| SPECT | fri_c1_100_10 | fri_c3_500_25 | ada_agnostic |
| SPECTF | elusage | hip | eye_movements |
| grub-damage | diabetes_numeric | analcatdata_negotiation | kc1-top5 |
| pasture | fri_c2_500_5 | chscase_census6 | mozilla4 |
| squash-stored | fri_c3_250_10 | fried | jEdit_4.2_4.3 |
| squash-unstored | fri_c2_250_25 | sleuth_case2002 | pc4 |
| white-clover | disclosure_x_tampered | fri_c2_1000_25 | pc3 |
| aids | cpu | fri_c0_1000_50 | jm1 |
| JapaneseVowels | fri_c4_1000_50 | chscase_census5 | mc2 |
| ipums_la-97-small | cholesterol | chscase_census4 | cm1_req |
| analcatdata_boxing2 | fri_c0_1000_5 | chscase_census3 | mc1 |
| prnn_crabs | pyrim | chscase_census2 | ar1 |
| analcatdata_boxing1 | pbcseq | fri_c1_1000_10 | ar3 |
| analcatdata_lawsuit | delta_ailerons | fri_c2_250_5 | ar4 |
| irish | hutsof99_logis | fri_c2_1000_5 | ar5 |
| analcatdata_broadwaymult | fri_c4_500_50 | fri_c2_1000_10 | kc2 |
| analcatdata_authorship | fri_c3_1000_50 | plasma_retinol | ar6 |
| analcatdata_asbestos | kin8nm | fri_c3_100_5 | kc3 |
| analcatdata_creditscore | fri_c0_100_10 | fri_c1_1000_25 | kc1-binary |
| prnn_synth | mushroom | fri_c4_250_50 | kc1 |
| analcatdata_cyyoung8092 | pbc | analcatdata_seropositive | pc1 |
| schizo | rmftsa_ctoarrivals | fri_c2_100_50 | pc2 |
| confidence | fri_c1_100_25 | visualizing_soil | mw1 |
| analcatdata_dmft | fri_c3_1000_5 | visualizing_galaxy | jEdit_4.0_4.2 |
| profb | chscase_vine2 | fri_c0_500_25 | desharnais |
| lupus | chscase_vine1 | disclosure_z | datatrieve |
| analcatdata_germangss | puma8NH | fri_c4_100_50 | teachingAssistant |
| prnn_viruses | diggle_table_a1 | fri_c4_250_25 | pc1_req |
| biomed | | | |