# OpenML,
# metafeatures
# and
# analyzing the results of experiments with filtering anomalies

Katarína Švecová
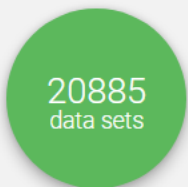
PV115

2020

# OpenML

# OpenML beta_2

## Machine learning, better, together

**20885**
data sets

**111443**
tasks

**14769**
flows

**10025717**
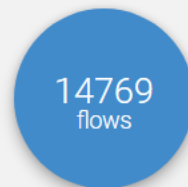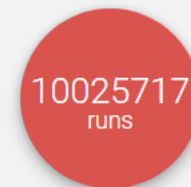runs

Find or add **data** to analyse

Download or create scientific **tasks**

Find or add data analysis **flows**

Upload and explore all **results** online.

ython API 0.10 is released

t started now **or** read the paper first :)

H

Bri

**Explore**

| | | |
|---|---|---|
| 🗄 | Data | 20885 |
| 🏆 | Task | 111443 |
| ⚙ | Flow | 14769 |
| ⭐ | Run | 10025717 |
| ⚗ | Study | 131 |
| 🚩 | Task type | 8 |
| 📊 | Measure | 226 |
| 👥 | People | 7559 |

📗 **Help**

📶 **Blog**

📣 **Contact**

❤ **Please cite us**

- Main goal – to make machine learning accesible
- open source project on GitHub
- Datasets, Tasks, Flows, Runs
- Study, task type, measure
- API (REST, Python, R, Java, C#) – enable downloading datasets, tasks and sharing results

# OpenML

Search

Press F11 to exit full screen

## Explore

| | | |
|---|---|---|
| 🗄 | **Data** | **2992** |
| 🏆 | Task | |
| ⚙ | Flow | |
| ⭐ | Run | |
| ⚗ | Study | |
| 🚩 | Task type | 8 |
| 👥 | People | |

📖 **Help**

📡 **Blog**

📣 **Contact**

❤ **Please cite us**

---

## 2992 results

🔽 FILTERS    SORT: MOST RUNS ▾    ☰ ID'S    ⊞ TABLE    ➕ ADD NEW

Only showing active (verified) datasets.

---

🗄 **credit-g (1)**    This dataset classifies people described by a set of attributes as good or bad credit risks. This da…
★ **505316 runs**    ❤ **16 likes**    ☁ **195 downloads**    📶 **211 reach**    ⚡ **12 impact**
1000 instances - 21 features - 2 classes - 0 missing values

---

🗄 **blood-transfusion-service-c...**    Data taken from the Blood Transfusion Service Center in Hsin-Chu City in Taiwan -- this is a classi…
★ **464739 runs**    ❤ **5 likes**    ☁ **67 downloads**    📶 **72 reach**    ⚡ **29 impact**
748 instances - 5 features - 2 classes - 0 missing values

---

🗄 **monks-problems-2 (1)**    Once upon a time, in July 1991, the monks of Corsendonk Priory were faced with a school held in …
★ **394292 runs**    ❤ **1 likes**    ☁ **21 downloads**    📶 **22 reach**    ⚡ **28 impact**
601 instances - 7 features - 2 classes - 0 missing values

---

🗄 **tic-tac-toe (1)**    This database encodes the complete set of possible board configurations at the end of tic-tac-toe…
★ **385593 runs**    ❤ **1 likes**    ☁ **65 downloads**    📶 **66 reach**    ⚡ **2 impact**
958 instances - 10 features - 2 classes - 0 missing values

---

🗄 **monks-problems-1 (1)**    Once upon a time, in July 1991, the monks of Corsendonk Priory were faced with a school held in …
★ **358449 runs**    ❤ **2 likes**    ☁ **18 downloads**    📶 **20 reach**    ⚡ **31 impact**
556 instances - 7 features - 2 classes - 0 missing values

---

🗄 **steel-plates-fault (1)**    A dataset of steel plates' faults, classified into 7 different types. The goal was to train machine lea…
★ **277313 runs**    ❤ **1 likes**    ☁ **38 downloads**    📶 **39 reach**    ⚡ **18 impact**
1941 instances - 34 features - 2 classes - 0 missing values

---

🗄 **kr-vs-kp (1)**    1. Title: Chess End-Game -- King+Rook versus King+Pawn on a7 (usually abbreviated KRKPA7). Th…
★ **270777 runs**    ❤ **0 likes**    ☁ **36 downloads**    📶 **36 reach**    ⚡ **5 impact**
3196 instances - 37 features - 2 classes - 0 missing values

ARFF   CSV   JSON   XML   RDF

# 🗄 tic-tac-toe

**active**   ⊞ **ARFF**   CC **Publicly available**   ⦸ Visibility: public   ☁ Uploaded 06-04-2014 by **Jan van Rijn**

♥ 1 likes   ☁ downloaded by 65 people , 74 total downloads   ⚠ 0 issues   👎 0 downvotes

🏷  **mythbusting_1**  **OpenML-CC18**  **OpenML100**  **study_1**  **study_123**  **study_135**  **study_14**  **study_144**  **study_15**  **study_20**  **study_29**  **study_30**  **study_37**  **study_41**  **study_52**  **study_7**  **study_70**  **study_89**  **study_98**  **study_99**  **uci**  **study_234**   ➕ Add tag

⟳ Loading wiki

Author: David W. Aha
Source: [UCI](https://archive.ics.uci.edu/ml/datasets/Tic-Tac-Toe+Endgame) - 1991
Please cite: [UCI](http://archive.ics.uci.edu/ml/citation_policy.html)

Tic-Tac-Toe Endgame database
This database encodes the complete set of possible board configurations at the end of tic-tac-toe games, where "x" is assumed to have played first. The target concept is "win for x" (i.e., true when "x" has one of 8 possible ways to create a "three-in-a-row").

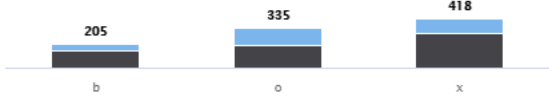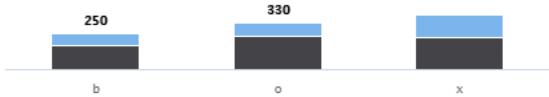▾ Show all

## 10 features

| Class **(target)** | nominal | 2 unique values 0 missing |
|---|---|---|

332 negative   626 positive

## 10 features

| Class **(target)** | nominal | 2 unique values<br>0 missing | |
| --- | --- | --- | --- |



| top-left-square | nominal | 3 unique values<br>0 missing | |
| --- | --- | --- | --- |



| top-middle-square | nominal | 3 unique values<br>0 missing | |
| --- | --- | --- | --- |



**❤ Show all 10 features**

## 107 properties

| | | | |
| --- | --- | --- | --- |
| 📊 **NumberOfInstances** | 958 | Number of instances (rows) of the dataset. | |
| 📊 **NumberOfFeatures** | 10 | Number of attributes (columns) of the dataset. | |
| 📊 **NumberOfClasses** | 2 | Number of distinct values of the target attribute (if it is nominal). | |
| 📊 **NumberOfMissingVal...** | 0 | Number of missing values in the dataset. | |
| 📊 **NumberOfInstancesW...** | 0 | Number of instances with at least one value missing. | |

**❤ Show all 107 properties**

## 23 tasks

| 🏆 | **Supervised Classification on tic-tac-toe**<br>**270253 runs** - estimation_procedure: 10-fold Crossvalidation - target_feature: Class |
| --- | --- |
| 🏆 | **Supervised Classification on tic-tac-toe** |

# Tasks

- Dataset with a specific task – clustering/classification and a method of evaluation

Search

Explore

Data

**Task**

Flow

Run

Study

Task type

Measure

People

# 🏆 Supervised Classification on tic-tac-toe

🏆 Task 145804    🏳 Supervised Classification    🗄 tic-tac-toe    ★ 270253 runs submitted

♥ 0 likes    ☁ downloaded by 8 people , 8 total downloads    ⚠ 0 issues

👁 Visibility: Public

🏷 study_107    ➕ Add tag

| EVALUATIONS | 👥 PEOPLE | ☰ RUNS | ➕ ADD RESULTS |

Metric:    PREDICTIVE ACCURACY ▾

## 270253 Runs

Help

Blog

Contact

♥ Please cite us

### Evaluations per flow (multiple parameter settings)

every point is a run, click for details

Predictive accuracy

| | 0.65 | 0.675 | 0.7 | 0.725 | 0.75 | 0.775 | 0.8 | 0.825 | 0.85 | 0.875 | 0.9 | 0.925 | 0.95 | 0.975 | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

mlr.classif.svm(6)

mlr.classif.xgboost(6)

mlr.classif.xgboost(9)

mlr.classif.ranger(13)

mlr.classif.ranger(9)

mlr.classif.ranger(16)

mlr.classif.ranger(15)

weka.kf.AttributeSelectio...liefF-Standardize-IBk5(1)

weka.kf.AdaBoostM1-IBk5(1)

mlr.classif.kknn(10)

weka.kf.Bagging-IBk5(1)

mlr.classif.glmnet(11)

weka.kf.AttributeSelection-Ranker-ReliefF-SMO(1)

weka.kf.Bagging-SMO(1)

# Flows

- Flow – a specific algorithm in a specific implementation

Search

Explore

Data

Task

Flow

Run

Study

Task type

Measure

People

Help

Blog

Contact

Please cite us

V. 1

# ⚙️ weka.kf.ReplaceMissingValues-J48

👁 Visibility: public   ☁ Uploaded 27-02-2018 by William Raynaut   ⛓ Weka_3.9.2   ⭐ 812 runs

❤ 0 likes   ☁ downloaded by 0 people   ⚠ 0 issues   👎 0 downvotes , 0 total downloads

🏷 mDataExp  Modelage  study_107   ➕ Add tag

A Weka KnowledgeFlow using ReplaceMissingValues-J48.kf

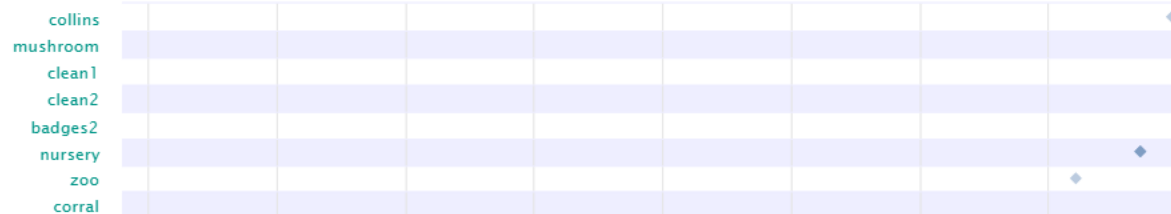⟳ Loading wiki

## Parameters

## 818 Runs

| SUPERVISED CLASSIFICATION ▼ | PREDICTIVE ACCURACY ▼ | Parameter: NONE ▼ | ☰ LIST ALL RUNS |

### Evaluations per dataset (multiple parameter settings)

every point is a run, click for details

collins
mushroom
clean1
clean2
badges2
nursery
zoo
corral

# Run

- A specific run with specific settings

Search

**Explore**

- Data
- Task
- Flow
- Run
- Study
- Task type
- Measure
- People

**Help**

**Blog**

**Contact**

**Please cite us**

JSON  XML  RDF

# Run 23504

🏆 Task 2076 (Supervised Classification)  🗄 kropt  ☁ Uploaded 25-06-2014 by Joaquin Vanschoren

♥ 1 likes  ☁ downloaded by 0 people  ⚠ 0 issues  👎 0 downvotes , 0 total downloads

🏷  ➕ Add tag

## Flow

| | |
|---|---|
| weka.RandomForest(1) | Leo Breiman (2001). Random Forests. Machine Learning. 45(1):5-32. |
| weka.RandomForest(1)_I | 11 |
| weka.RandomForest(1)_K | 0 |
| weka.RandomForest(1)_S | 1 |
| weka.RandomForest(1)_num-slots | 1 |

## Result files

⬇ **Description**                                                    xml
XML file describing the run, including user-defined evaluation measures.

⬇ **Predictions**                                                   arff
ARFF file with instance-level predictions generated by the model.

19 Evaluation measures

# OpenML

HELP    SIGN IN

**Explore**

- Data
- Task
- Flow
- **Run**
- Study
- Task type
- Measure
- People

**Help**

**Blog**

**Contact**

**Please cite us**

## Area under ROC curve

### 0.9444

#### Per class

| draw | zero | one | two | three | four | five | six | seven | eight | nine | ten | eleven | twelve |
|------|------|-----|-----|-------|------|------|-----|-------|-------|------|-----|--------|--------|
| 0.9897 | 0.9993 | 0.992 | 0.9987 | 0.9978 | 0.9944 | 0.9867 | 0.9804 | 0.9641 | 0.9684 | 0.9579 | 0.9335 | 0.9177 | 0.9167 |

#### Cross-validation details (10-fold Crossvalidation)



## F measure

### 0.6426

#### Per class

| draw | zero | one | two | three | four | five | six | seven | eight | nine | ten | eleven | twelve |
|------|------|-----|-----|-------|------|------|-----|-------|-------|------|-----|--------|--------|
| 0.8442 | 0.4286 | 0.6406 | 0.8187 | 0.5547 | 0.6431 | 0.6694 | 0.6705 | 0.5631 | 0.6436 | 0.6087 | 0.5549 | 0.5584 | 0.598 |

#### Cross-validation details (10-fold Crossvalidation)

```python
import openml
from sklearn import impute, tree, pipeline

# Define a scikit-learn classifier or pipeline
clf = pipeline.Pipeline(
    steps=[
        ('imputer', impute.SimpleImputer()),
        ('estimator', tree.DecisionTreeClassifier())
    ]
)
# Download the OpenML task for the german credit card dataset with 10-fold
# cross-validation.
task = openml.tasks.get_task(31)
# Run the scikit-learn model on the task.
run = openml.runs.run_model_on_task(clf, task)
# Publish the experiment on OpenML (optional, requires an API key.
# You can get your own API key by signing up to OpenML.org)
run.publish()
print(f'View the run online: {openml.config.server}/run/{run.run_id}')
```

# Metafeatures

# Types of metafeatures

- Basic data description (number of instances, dimensionality, number of classes, majority class percentage…)

- Statistical methods (mean, median, skewness, kurtosis…)

- Landmarking – results with (quite a few) selected methods (
    random tree **depth**,

    DecisionStumpKappa, NaiveBayes**Kappa -** the agreement between two raters, similar to accuracy, but considering the probability of a chance agreement,

    J48 **error rate**,

    **A**rea **U**nder the ROC **C**urve, which is made by plotting true positive rate and false positive rate)

| | | |
|---|---|---|
| 📊 **NumberOfInstances** | **1000** | Number of instances (rows) of the dataset. |
| 📊 **NumberOfFeatures** | 21 | Number of attributes (columns) of the dataset. |
| 📊 **NumberOfClasses** | 2 | Number of distinct values of the target attribute (if it is nominal). |
| 📊 **NumberOfMissingVal...** | 0 | Number of missing values in the dataset. |
| 📊 **NumberOfInstancesW...** | 0 | Number of instances with at least one value missing. |
| 📊 **NumberOfNumericFe...** | 7 | Number of numeric attributes. |
| 📊 **NumberOfSymbolicFe...** | 14 | Number of nominal attributes. |
| 📊 **PercentageOfBinaryF...** | 14.29 | Percentage of binary attributes. |

Number of different things

| | | | |
|---|---|---|---|
| 📊 | **Quartile2StdDevOfNu...** | 1.12 | Second quartile (Median) of standard deviation of attributes of the n... |
| 📊 | **RandomTreeDepth1A...** | 0.66 | Area Under the ROC Curve achieved by the landmarker weka.classifie... |
| 📊 | **Dimensionality** | 0.02 | Number of attributes divided by the number of instances. |
| 📊 | **MaxMutualInformation** | 0.09 | Maximum mutual information between the nominal attributes and th... |
| 📊 | **MinNominalAttDistinc...** | 2 | The minimal number of distinct values among attributes of the nomi... |
| 📊 | **PercentageOfInstanc...** | 0 | Percentage of instances having missing values. |
| 📊 | **Quartile3AttributeEntr...** | 1.87 | Third quartile of entropy among attributes. |

| | **kNN1NKappa** | Kappa coefficient achieved by the landmarker weka.classifiers.lazy.IBk |
| | | data quality |
| | **DecisionStumpAUC** | Area Under the ROC Curve achieved by the landmarker weka.classifiers.trees.D... |
| | | data quality |
| | **J48.00001.ErrRate** | Error rate achieved by the landmarker weka.classifiers.trees.J48 -C .00001 |
| | | data quality |
| | **J48.00001.Kappa** | Kappa coefficient achieved by the landmarker weka.classifiers.trees.J48 -C .00... |
| | | data quality |
| | **J48.0001.ErrRate** | Error rate achieved by the landmarker weka.classifiers.trees.J48 -C .0001 |
| | | data quality |

# On the predictive power of meta-features in OpenML

- 2017

- Study on 61 metafeatures and 720 datasets

- Feature selection gains better results

- Information on response, mutual information, noise to signal, shape of extremes, **dimensionality,** minimum variability of numeric attributes, information of categorical attributes and mutual information

Analyzing the results of experiments with filtering anomalies

Dataset,"NumberOfInstances","NumberOfFeatures","NumberOfClasses","NumberOfMissingValues","NumberOfInstancesWithMissingValues","NumberOfNumericFeatures","NumberOfSymbolicFeatures","RandomTreeDepth1Kappa","J48.00001.AUC","MaxSkewnessOfNumericAtts","MinStdDevOfNumericAtts","PercentageOfMissingValues","Quartile3KurtosisOfNumericAtts","AutoCorrelation","RandomTreeDepth2AUC","J48.00001.ErrRate","MaxStdDevOfNumericAtts","MinorityClassPercentage","PercentageOfNumericFeatures","Quartile3MeansOfNumericAtts","CfsSubsetEval_DecisionStumpAUC","RandomTreeDepth2ErrRate","J48.00001.Kappa","MeanAttributeEntropy","MinorityClassSize","PercentageOfSymbolicFeatures","Quartile3MutualInformation","CfsSubsetEval_DecisionStumpErrRate","RandomTreeDepth2Kappa","J48.0001.AUC","MeanKurtosisOfNumericAtts","NaiveBayesAUC","Quartile1AttributeEntropy","Quartile3SkewnessOfNumericAtts","CfsSubsetEval_DecisionStumpKappa","RandomTreeDepth3AUC","J48.0001.ErrRate","MeanMeansOfNumericAtts","NaiveBayesErrRate","Quartile1KurtosisOfNumericAtts","Quartile3StdDevOfNumericAtts","CfsSubsetEval_NaiveBayesAUC","CfsSubsetEval_NaiveBayesErrRate","RandomTreeDepth3ErrRate","J48.0001.Kappa","MeanMutualInformation","NaiveBayesKappa","Quartile1MeansOfNumericAtts","REPTreeDepth1AUC","CfsSubsetEval_NaiveBayesKappa","RandomTreeDepth3Kappa","J48.001.AUC","MeanNoiseToSignalRatio","NumberOfBinaryFeatures","Quartile1MutualInformation","REPTreeDepth1ErrRate","CfsSubsetEval_kNN1NAUC","StdvNominalAttDistinctValues","J48.001.ErrRate","MeanNominalAttDistinctValues","Quartile1SkewnessOfNumericAtts","REPTreeDepth1Kappa","REPTreeDepth2AUC","CfsSubsetEval_kNN1NErrRate","kNN1NAUC","J48.001.Kappa","MeanSkewnessOfNumericAtts","Quartile1StdDevOfNumericAtts","REPTreeDepth2ErrRate","CfsSubsetEval_kNN1NKappa","kNN1NErrRate","MajorityClassPercentage","MeanStdDevOfNumericAtts","Quartile2AttributeEntropy","REPTreeDepth2Kappa","ClassEntropy","kNN1NKappa","MajorityClassSize","MinAttributeEntropy","Quartile2KurtosisOfNumericAtts","REPTreeDepth3AUC","DecisionStumpAUC","MaxAttributeEntropy","MinKurtosisOfNumericAtts","Quartile2MeansOfNumericAtts","REPTreeDepth3ErrRate","DecisionStumpErrRate","MaxKurtosisOfNumericAtts","MinMeansOfNumericAtts","Quartile2MutualInformation","REPTreeDepth3Kappa","DecisionStumpKappa","MaxMeansOfNumericAtts","MinMutualInformation","Quartile2SkewnessOfNumericAtts","RandomTreeDepth1AUC","Dimensionality","MaxMutualInformation","MinNominalAttDistinctValues","PercentageOfBinaryFeatures","Quartile2StdDevOfNumericAtts","RandomTreeDepth1ErrRate","EquivalentNumberOfAtts","MaxNominalAttDistinctValues","MinSkewnessOfNumericAtts","PercentageOfInstancesWithMissingValues","Quartile3AttributeEntropy"
"banknote-authentication",1372,5,2,0,0,4,1,2,20,3.576396338280576,0.8148164881029215,0.0036443148688046646,null,-1.022243043654092,0,null,0.2244897959183735,null,2,2.101013136739067,0,1.077230875627456,0.9992706053975201,0.5620538416256953,0.9854814767006583,1.0885685435217796,44.46064139941695,80,1.7911716198979,0.8344961490469428,0.9330450496966567,0.02113702623906705,5.8690467435803795,610,20,null,0.19314868804664723,0.09912536443148688,0.9571796328821847,null,0.9396067294866831,null,0.7790794824458168,0.6220644978877691,0.8005251139599922,0.9854814767006583,0.14479469144119225,0.1588921282798834,-0.67298891708841984,5.47929258026428,0.9528828363667657,0.9764080719418269,0.021137026239067054,0.6405147434402338,0.6763618735987396,-0.7853085765306612,0.8344961490469428,0.12463556851311954,0.0663265306122449,0.9571796328821847,null,null,1,null,0.19314868804664723,0.7466787022320504,0.8674539679304253,0.9854814767006583,2,-0.8652081513981438,0.6220644978877691,0.9382578202314874,0,0.021137026239067054,-0.119291418861707,2.286504991155923,0.9101996471752506,0.061224489795918366,0.9986676640419947,0.9571796328821847,3.7807131392201496,null,0.11443148688046648,0.876125584909549,0.001457259475218659,55.53935860058085,null,0.030142115780319223,0.7692363729859445,0.9911281257467461,0.9970492846881102,762,-0.7515813814324717,0.9156811869533534,0.9400972419431177,0.8344961490469428,null,-1.1916565211370262,null,0.07434402332361516,0.19314868804664723,1.2704759156366023,null,-0.27174558770618523,0.8515038072681896,0.6220644978877691,1.9223531209912545

# Example of data

"dataset","clf","clf_family","clf_params","od_name","od_params","removed","accuracy","od
_time","clf_time","total_time","gain_clf","gain_clfBest","gain_random"
"JM1","IBk","lazy","[]","CODB","{jar_path: data/java/WEKA-
CODB.jar}",0.5,0.76463,318.14848,6.88443,325.03291,0.00193,-0.05255,0.00165
"JM1","IBk","lazy","[]","CODB","{jar_path: data/java/WEKA-
CODB.jar}",1,0.76601,318.14848,6.76814,324.91662,0.00331,-0.05117,0.00257
"JM1","IBk","lazy","[]","CODB","{jar_path: data/java/WEKA-
CODB.jar}",2,0.77042,318.14848,6.76459,324.91307,0.00772,-0.04676,0.0079

```prolog
?- search_data(all, [gain_clfBest > 0], R),
|    stats(R, dataset, count, [dataset], Result),
|     make_csv(Result, "c:/users/katka/desktop/lab/prolog/
10jul/grafy/dataset_best_over_0_count.csv").
```

dataset_best_over_0_count.csv - Notepad

File  Edit  Format  View  Help

```
dataset,count:dataset
"wilt2",32
"wdbc",11
"wall-robot-navigation",7
"vehicle",32
"texture",8
"spf3",47
"qsar-biodeg",57
"pima-diabetes",98
"phoneme",13
"phishing",16
"pc4",28
"pc3",47
"madelon",45
"letter",1
"ilpd",63
"first order theorem",4
```

```
?- search_data(all, [gain_clfBest > 0], R), stats(R, gain_clfBest,
count, [], Result).
Result = [[count:gain_clfBest], [1178]].


                    ?- search_data(R), stats(R, dataset, count, Result).
                    Result = [['count:dataset'], [21373]]




?- search_data(all, [gain_clfBest > 0], R), stats(R, gain_clfBest,
count, [od_name], Result).
Result = [[od_name, count:gain_clfBest], ['"TDWithPrunning"', 145],
 ['"TD"', 97], ['"Random"', 122], ['"NearestNeighbors"', 122], ['"L
OF"', 115], ['"KDN"', 88], ['"IsolationForest"'|...], [...|...]|...
].
```
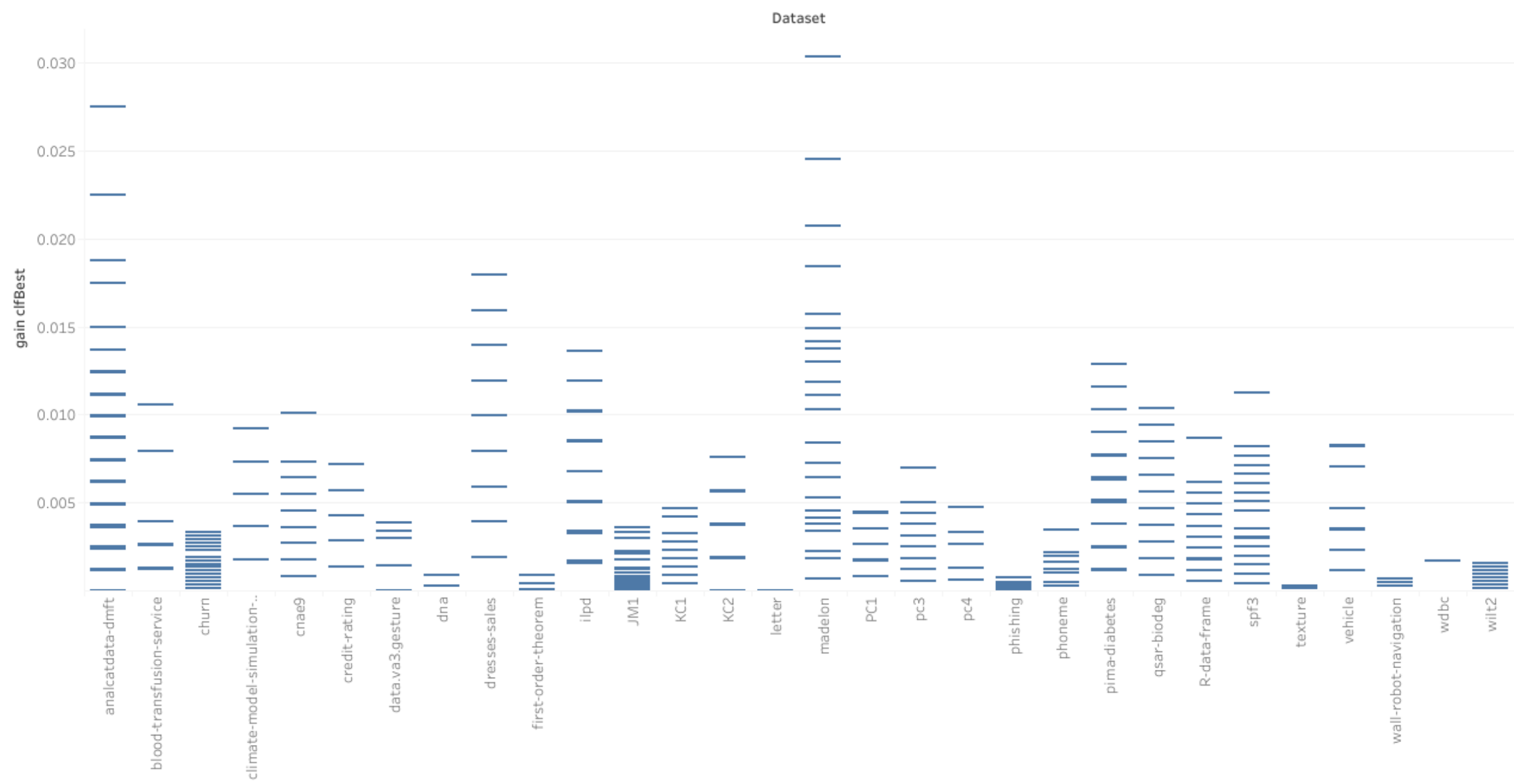
```
?- search_data(all, [gain_clfBest > 0], R1), stats(R1, gain_clfBest
, count, [od_name], R2), condi(R2, [_, BestCount], BestCount > 100,
 Result).
Result = [[od_name, 'count:gain_clfBest'], ['"TDWithPrunning"', 145
], ['"Random"', 122], ['"NearestNeighbors"', 122], ['"LOF"', 115],
['"DS"', 102], ['"ClassLikelihood"', 108]].
```
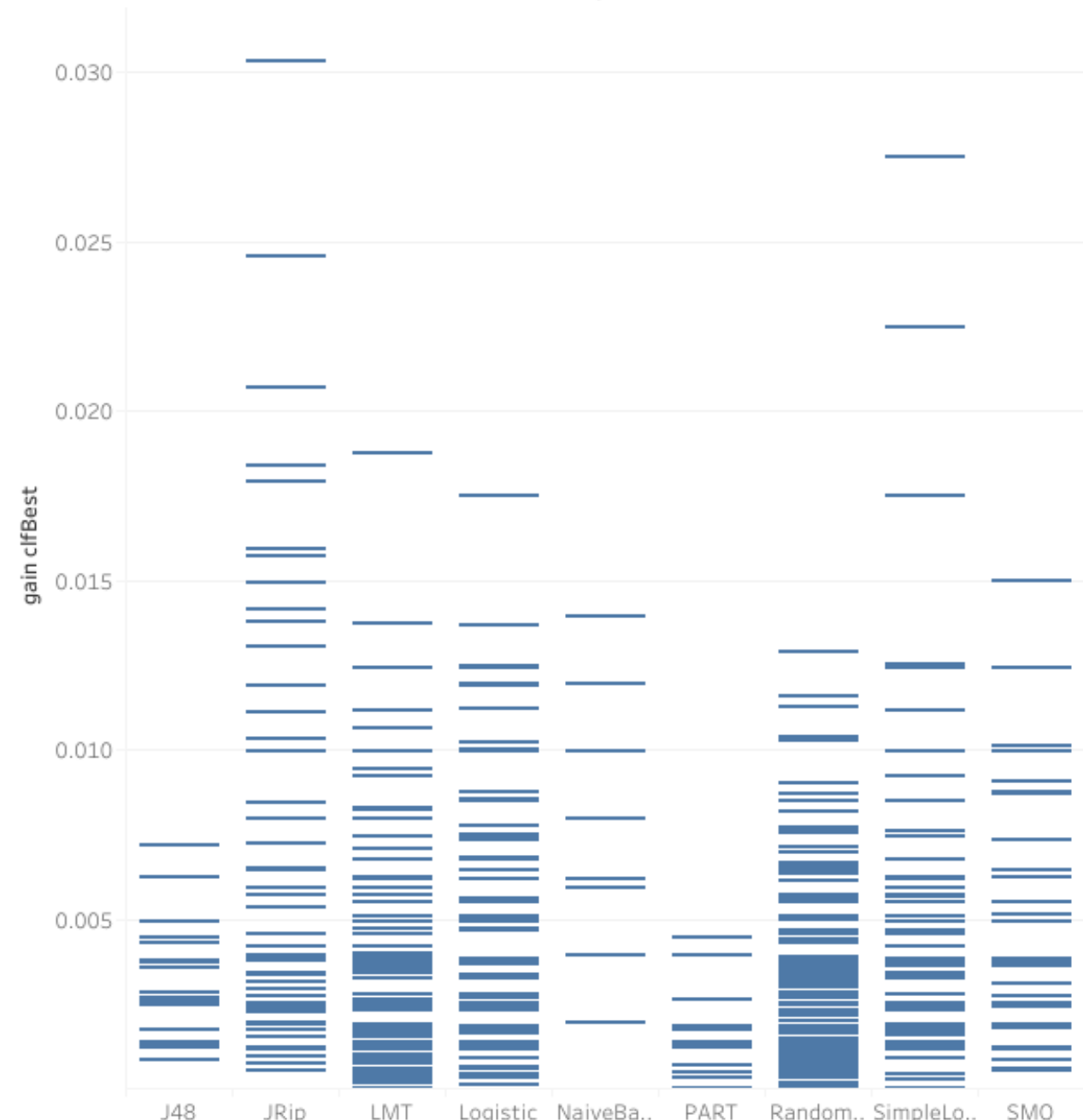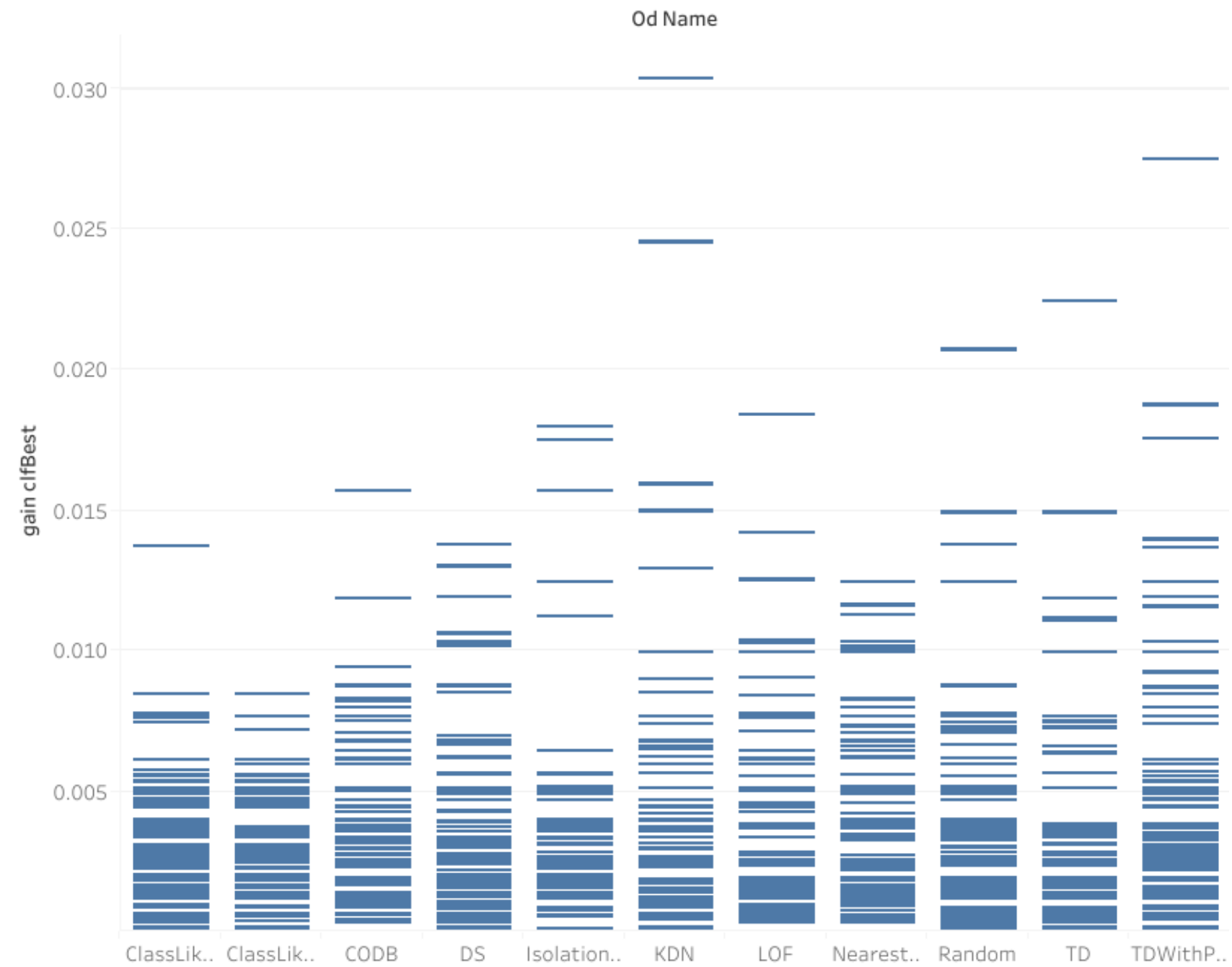
# Graphs

Positive
gain_clfBest

# Sheet 3

Sheet 3 (12)

# Sheet 3 (13)

# Sheet 3 (6)

# Sheet 3 (14)

# Sheet 3 (4)

# Sheet 3 (3)

# Sheet 3 (2)



Dataset

Od Name
- ClassLikelihood
- ClassLikelihoodDifference
- CODB
- DS
- IsolationForest
- KDN
- LOF
- NearestNeighbors
- Random
- TD
- TDWithPrunning

Sheet 3 (20)

Sheet 3 (17)

# Sheet 3 (17)

# Sheet 3 (16)



Clf

gain clfBest

Od Name
- ClassLikelihood
- ClassLikelihoodDifference
- CODB
- DS
- IsolationForest
- KDN
- LOF
- NearestNeighbors
- Random
- TD
- TDWithPrunning

J48  JRip  LMT  Logistic  NaiveBa..  PART  Random..  SimpleLo..  SMO
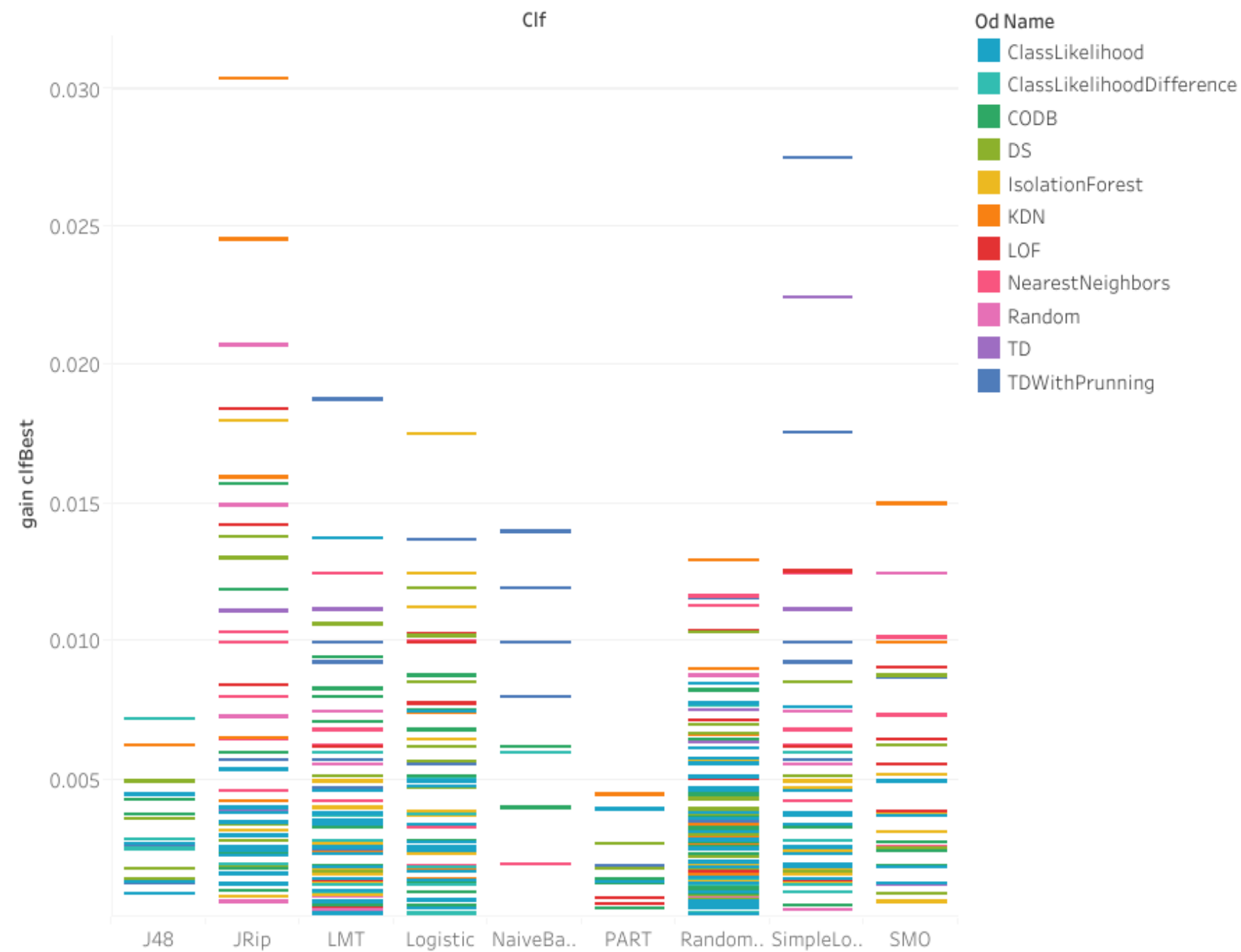
Sheet 3 (15)

Clf Family / Clf

| bayes | functions | rules | trees |

NaiveBa.. | Logistic | SimpleLo.. | SMO | JRip | PART | J48 | LMT | Random..

gain clfBest

Sheet 3 (11)

Clf Family

bayes | functions | rules | trees

gain clfBest

# Count of gain_clfBest

# Sheet 3 (27)

**Clf**

# Sheet 3 (26)



Od Name — Count of gain clfBest bar chart with categories: ClassLik.., ClassLik.., CODB, DS, Isolation.., KDN, LOF, Nearest.., Random, TD, TDWithP..

# Sheet 3 (25)



Removed

Count of gain clfBest

| | 0.5 | 1 | 2 | 3 | 4 | 5 |

# Sheet 3 (24)

# Sheet 3 (23)



Stacked bar chart. X-axis: Dataset. Y-axis: Count of gain clfBest. Legend titled "Od Name" with categories: ClassLikelihood, ClassLikelihoodDifference, CODB, DS, IsolationForest, KDN, LOF, NearestNeighbors, Random, TD, TDWithPrunning.

Datasets: analcatdata-dmft, blood-transfusion-service, churn, climate-model-simulation-.., cnae9, credit-rating, data.va3.gesture, dna, dresses-sales, first-order-theorem, ilpd, JM1, KC1, KC2, letter, madelon, PC1, pc3, pc4, phishing, phoneme, pima-diabetes, qsar-biodeg, R-data-frame, spf3, texture, vehicle, wall-robot-navigation, wdbc, wilt2

# Sheet 3 (22)

# Sheet 3 (29)

# Sheet 3 (30)



Bar chart titled "Od Name" with y-axis "Count of gain clfBest" ranging from 0 to 140. Legend "Removed" with values 0.5, 1, 2, 3, 4, 5. Categories along x-axis: ClassLik.., ClassLik.., CODB, DS, Isolation.., KDN, LOF, Nearest.., Random, TD, TDWithP..

# Sheet 3 (28)

# Other graphs

# Sheet 3 (3)



Clf Family
- bayes
- functions
- lazy
- rules
- trees

Sheet 3 (10)

Sheet 3 (18)

Sheet 3 (9)

Sheet 3 (8)

Sheet 3 (5)

Sheet 3 (5)

**Dataset**
- analcatdata-dmft
- banknote-authentication
- blood-transfusion-service
- churn
- climate-model-simulation-crashes
- cnae9
- credit-rating
- data.va3.gesture
- dna
- dresses-sales
- first-order-theorem
- ilpd
- JM1
- KC1
- KC2
- letter
- madelon
- mice-protein
- PC1
- pc3
- pc4
- phishing
- phoneme
- pima-diabetes
- qsar-biodeg
- R-data-frame
- spf3
- texture
- vehicle
- wall-robot-navigation
- wdbc
- wilt2

Sheet 3 (19)

# Sheet 3 (21)



Dataset

- https://docs.openml.org/

- https://www.openml.org/

- https://www.openml.org/search?type=data

- https://www.openml.org/d/50

- https://www.openml.org/search?type=task

- https://www.openml.org/t/145804

- https://www.openml.org/search?type=flow

- https://www.openml.org/f/7791

- https://www.openml.org/search?type=run

- https://www.openml.org/r/23504

- https://docs.openml.org/APIs/

- https://en.wikipedia.org/wiki/Cohen's_kappa

- https://en.wikipedia.org/wiki/Receiver_operating_characteristic#:~:text=A%20receiver%20operating%20characteristic%20curve,why%20it%20is%20so%20named.

- Besim BILALLI, Alberto ABELL´O, Tom´as ALUJA-BANET. On the predictive power of meta-features in OpenML. International Journal of Applied Mathematics and Computer Science. 2017, Vol 27, Iss 4, 2083-8492.