

OpenML Use Cases

Jan N. van Rijn and Joaquin Vanschoren

March 23th, 2019

Table of Contents

BENCHMARKING



**THE
BUSTING
MYTHS**



BENCHMARKING



OpenML Benchmark Suites

OpenML Benchmark Suites

Common limitations of experimental evaluations:

- Often done on a small selection of datasets
 - Not sure about generalization to more datasets
 - 'Cherry picking'
- Publication bias
 - Every paper only reports good results
 - More useful to know WHEN (on what type of data) a new algorithm performs well
- Hard to compare conclusions across papers
- Abused datasets (example: liver-disorders)

OpenML Benchmark Suites

OpenML-100: A curated benchmark suite [Bischl et al., 2017].

Inclusion criteria:

- Basic data properties (500–100000 observations, < 5000 features, ≥ 2 classes)
- The ratio of the minority class and the majority class > 0.05
- Scientific publication introducing the dataset and learning task
- Pure classification tasks only (no data streams, multi-class tasks)
- No artificial data, binarized versions of regression datasets or sub-samples of bigger datasets
- Disadvantage: Very hard to get a good consensus about the inclusion criteria
- Next: The OpenML-CC18 (approx. 71 datasets)

OpenML Benchmark Suites

Benchmark suite:

- A collection of tasks
 - Terminology: A Study is a collection of runs
 - Collection of flows, datasets, ...
- Immutable (once closed)
- Standardized train-test splits are provided to ensure that results can be objectively compared
- results can be shared in a reproducible way through the APIs
- <https://docs.openml.org/benchmark/>

[B. Bischl, G. Casalicchio, M. Feurer, F. Hutter, M. Lang, R. G. Mantovani, J. N. van Rijn, and J. Vanschoren. [Openml benchmarking suites and the openml100](#).

arXiv preprint arXiv:1708.03731, 2017]

Requirements

- A set of tasks
 - Upload dataset and create tasks (task creation soon available)
 - Use existing tasks
- Latest version of Python API (Development version)
 - `openml.dataset.list_datasets(...)`
 - `openml.study.create_benchmark_suite(...)`
- OpenML account and user id (needs to be configured)
- Also available for Java (see docs) and R (coming soon)

Create Benchmark Suite

```
1 import openml
2
3 tasks = openml.tasks.list_tasks(
4     number_instances='100..500',
5     number_features='4..20',
6     size=20)
7 # task is a Dict[int, OpenMLTask]
8
9 study = openml.study.create_benchmark_suite(
10     alias=None,
11     name="Benchmark Example",
12     description="illustrates creating benchmark suites",
13     task_ids=tasks.keys()
14 )
15 study_id = study.publish()
16 print('Uploaded study with id=%d' % study_id)
```

Find Benchmark Suite

```
1 import openml
2
3 studies = openml.study.list_studies(
4     main_entity_type='task',
5     creator=1,
6     status='all'
7 )
```

Lists all studies / benchmark suites that comply to a set of filters

- Legal filters: `main_entity_type`, `uploader`, `status`, ...
- `studies` is now a `Dict[int, OpenMLStudy]`
- Note: Recently created studies are 'in preparation'

Common operations

Attach additional tasks

```
1 tasks = openml.tasks.list_tasks(data_name='letter', size=1)
2 # tasks_new is Dict[int, OpenMLTask]
3 openml.study.attach_to_study(study_id, tasks.keys())
```

Common operations

Attach additional tasks

```
1 tasks = openml.tasks.list_tasks(data_name='letter', size=1)
2 # tasks_new is Dict[int, OpenMLTask]
3 openml.study.attach_to_study(study_id, tasks.keys())
```

Detach tasks

```
1 # given a variable study_id (int)
2 task_id = [2, 3, 4]
3 openml.study.detach_from_study(study_id, task_id)
```

Common operations

Attach additional tasks

```
1 tasks = openml.tasks.list_tasks(data_name='letter', size=1)
2 # tasks_new is Dict[int, OpenMLTask]
3 openml.study.attach_to_study(study_id, tasks.keys())
```

Detach tasks

```
1 # given a variable study_id (int)
2 task_id = [2, 3, 4]
3 openml.study.detach_from_study(study_id, task_id)
```

Activate benchmark suite (no mutations possible afterwards)

```
1 # given a variable study_id (int)
2 openml.study.status_update(study_id, 'active')
```

Benchmark Suite

OpenML
HELP SIGN IN

Explore

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

Help
 Blog
 Contact
 Please cite us

Collaborative, reproducible benchmarking and analysis

Created 21-02-2019 by Jan van Rijsen · Visibility: public

SEARCH THESE TASKS IN MORE DETAIL

DESCRIPTION
 101 DATA SETS
 101 TASKS
 374 FLOWS
 1290 RUNS

Supervised Classification on SpeedDating

22360 runs
 1 likes
 0 downloads
 0 reach
 4084 impact

uploader_id: 1 · estimation_procedure: 10-fold Crossvalidation · target_feature: match · reuse: 4084

Supervised Classification on connect-4

7914 runs
 0 likes
 0 downloads
 0 reach
 3009 impact

uploader_id: 1 · estimation_procedure: 10-fold Crossvalidation · target_feature: class · reuse: 3009

Supervised Classification on Australian

207732 runs
 0 likes
 2 downloads
 2 reach
 200384 impact

uploader_id: 1 · estimation_procedure: 10-fold Crossvalidation · target_feature: Y · reuse: 200384

Supervised Classification on texture

14509 runs
 0 likes
 0 downloads
 0 reach
 8975 impact

uploader_id: 1 · estimation_procedure: 10-fold Crossvalidation · target_feature: Class · reuse: 8975

Supervised Classification on LED-display-domain-7digit

11959 runs
 0 likes
 1 downloads
 1 reach
 4790 impact

uploader_id: 1 · estimation_procedure: 10-fold Crossvalidation · target_feature: Class · reuse: 4790

Supervised Classification on dresses-sales

16153 runs
 0 likes
 0 downloads
 0 reach
 8950 impact

uploader_id: 1 · estimation_procedure: 10-fold Crossvalidation · target_feature: Class · reuse: 8950

Supervised Classification on Amazon_employee_access

21099 runs
 0 likes
 0 downloads
 0 reach
 15795 impact

uploader_id: 1 · estimation_procedure: 10-fold Crossvalidation · target_feature: target · reuse: 15795

Supervised Classification on MiceProtein

9738 runs
 0 likes
 0 downloads
 0 reach
 4061 impact

uploader_id: 1 · estimation_procedure: 10-fold Crossvalidation · target_feature: class · reuse: 4061 · reach_of_reuse: 1

Benchmark Suite

```
1 import openml # development branch
2 import sklearn # version 0.20.0 and up
3
4 # given a variable study_id (int, str)
5 benchmark_suite = openml.study.get_study(study_id, 'tasks')
6 # build a sklearn classifier
7 clf = sklearn.pipeline.make_pipeline(
8     sklearn.preprocessing.Imputer(),
9     sklearn.tree.DecisionTreeClassifier()
10 )
11 # iterate over all tasks
12 for task_id in benchmark_suite.tasks:
13     task = openml.tasks.get_task(task_id) # download the OpenML task
14     X, y = task.get_X_and_y() # get the data (not used in this example)
15     # run classifier on splits (requires API key)
16     run = openml.runs.run_model_on_task(clf, task)
17     # print accuracy score
18     score = run.get_metric_score(sklearn.metrics.accuracy_score)
19     print('Data set: %s; Accuracy: %0.2f' % (task.get_dataset().name,
20                                             score.mean()))
21     run.publish() # publish the experiment on OpenML (optional)
22     print('URL for run: %s/run/%d' % (openml.config.server, run.run_id))
```

**BUSTING
MYTHS**

Myth Busting for Data Mining

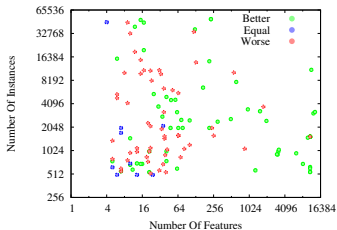
- Papers are generally build upon claims that are not well grounded, e.g.,
 - “We performed data transformation X because it is common practise.”
 - “We set hyperparameter Y to value Z because the authors recommended these values.”
- We can empirically analyze the validity of these claims on the meta-data from OpenML
- In this case: “The Importance of Feature Selection”

Effect of Feature Selection

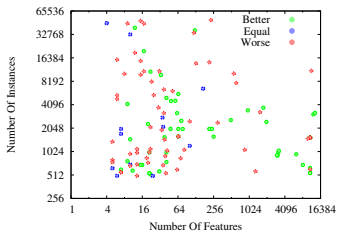
Experimental Setting by Post et al. [2016]:

- 400 binary classification datasets from OpenML
- 12 algorithms from Weka
- Correlation-based Feature Subset Selection
- We added runs that not existed on OpenML
- Recorded Area Under the ROC curve
- Limitation: Hyperparameter Optimization

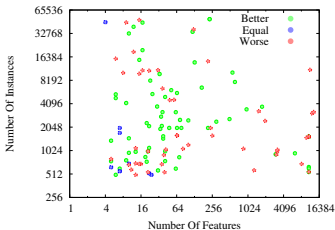
Effect of Feature Selection



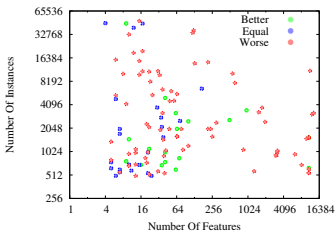
k-NN



J48



Naive Bayes



SMO

Effect of Feature Selection

Initial conclusions:

- Feature selection is often beneficial for the classifiers for which we expect it to be: k-NN and Naive Bayes
 - Surprisingly also for Decision Trees
- Feature selection is beneficial in 41% of the cases, only statistically significant in 10%
- Whether or not to use feature selection can be learned (see paper)
- Low amount of datasets on which feature selection significantly effects performance potentially indicates data bias
- Current Work: Linear vs. Non Linear classifiers

[M. J. Post, P. van der Putten, and J. N. van Rijn. [Does Feature Selection Improve Classification? A Large Scale Experiment in OpenML](#).

In *Advances in Intelligent Data Analysis XV*, pages 158–170. Springer, 2016]

Linear vs. Non-linear

Reoccurring topic in literature, see [Strang et al., 2018] for several examples

	Linear	Non-linear
ease of tuning		
interpretability		
fit risk		
performance		

Linear vs. Non-linear

Reoccurring topic in literature, see [Strang et al., 2018] for several examples

	Linear	Non-linear
ease of tuning	+	-
interpretability		
fit risk		
performance		

Linear vs. Non-linear

Reoccurring topic in literature, see [Strang et al., 2018] for several examples

	Linear	Non-linear
ease of tuning	+	-
interpretability	+	-
fit risk		
performance		

Linear vs. Non-linear

Reoccurring topic in literature, see [Strang et al., 2018] for several examples

	Linear	Non-linear
ease of tuning	+	-
interpretability	+	-
fit risk	underfit	overfit
performance		

Linear vs. Non-linear

Reoccurring topic in literature, see [Strang et al., 2018] for several examples

	Linear	Non-linear
ease of tuning	+	-
interpretability	+	-
fit risk	underfit	overfit
performance	-	+

Linear vs. Non-linear

Reoccurring topic in literature, see [Strang et al., 2018] for several examples

	Linear	Non-linear
ease of tuning	+	-
interpretability	+	-
fit risk	underfit	overfit
performance	-	+
Tree	Decision Stump	Decision Tree

Linear vs. Non-linear

Reoccurring topic in literature, see [Strang et al., 2018] for several examples

	Linear	Non-linear
ease of tuning	+	-
interpretability	+	-
fit risk	underfit	overfit
performance	-	+
Tree	Decision Stump	Decision Tree
SVM	Linear Kernel	Gaussian Kernel

Linear vs. Non-linear

Reoccurring topic in literature, see [Strang et al., 2018] for several examples

	Linear	Non-linear
ease of tuning	+	-
interpretability	+	-
fit risk	underfit	overfit
performance	-	+
Tree	Decision Stump	Decision Tree
SVM	Linear Kernel	Gaussian Kernel
Neural Network	Perceptron	MLP

Requirements

- A Benchmark suite
 - This case: the OpenML-100
- Latest version of Python API (Development version)
- Run results of the linear and non-linear classifier
 - We will generate them
- Listing function `openml.evaluation.list_evaluations(...)`
- Plotting library (matplotlib)

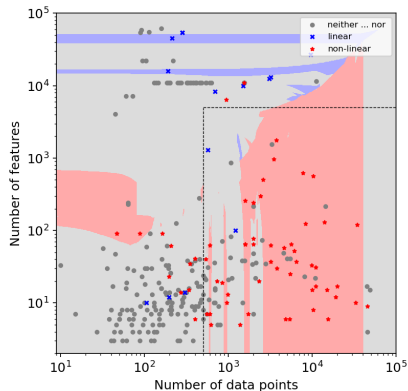
Generate Results

```
1 import openml # development branch
2 import sklearn # version 0.20.0 and up
3
4 # given a variable study_id (int, str)
5 benchmark_suite = openml.study.get_study('OpenML100', 'tasks')
6 # build a sklearn classifier
7 clfs = [
8     sklearn.pipeline.make_pipeline( # non-linear
9         sklearn.preprocessing.Imputer(),
10        sklearn.svm.SVC()
11    ),
12    sklearn.pipeline.make_pipeline( # linear
13        sklearn.preprocessing.Imputer(),
14        sklearn.svm.LinearSVC()
15    ),
16 ]
17 run_ids = list()
18 for task_id in benchmark_suite.tasks:
19     task = openml.tasks.get_task(task_id)
20     for clf in clfs:
21         run = openml.runs.run_model_on_task(clf, task)
22         run.publish()
23         print('URL for run: %s/run/%d' % (openml.config.server,
24                                           run.run_id))
25         run.push_tag('linear-vs-nonlinear')
```

Plot Results

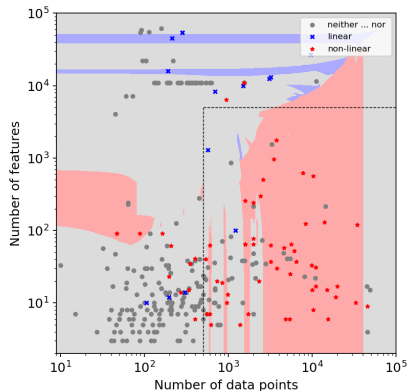
```
1 import openml # development branch
2
3 # given a variable study_id (int, str)
4 suite = openml.study.get_study('OpenML100', 'tasks')
5 tasks = openml.tasks.list_tasks(task_ids=suite.tasks)
6 # given a variable uploader_id (int)
7 evals = openml.evaluation.evaluation_list(
8     'predictive_accuracy',
9     uploader=uploader_id,
10    tag='linear-vs-nonlinear'
11 )
12
13 # organization
14 results = collections.defaultdict(dict)
15 for evaluation in evals.values():
16     results[evaluation.setup_id][evaluation.task_id] = evaluation.value
17 # plot
18 for setup_id in results.keys():
19     # find the tasks on which this evaluation is best
20     res_x = []
21     res_y = []
22     for task_id in results[setup_id].keys():
23         if results[setup_id][task_id] == [results[sid][task_id] for sid in results.keys()]:
24             res_x.append(tasks[task_id]['NumberOfInstances'])
25             res_y.append(tasks[task_id]['NumberOfFeatures'])
26     plt.scatter(res_x, res_y)
27 # set labels, titles and (log-)scales
```

Linear vs. Non-Linear

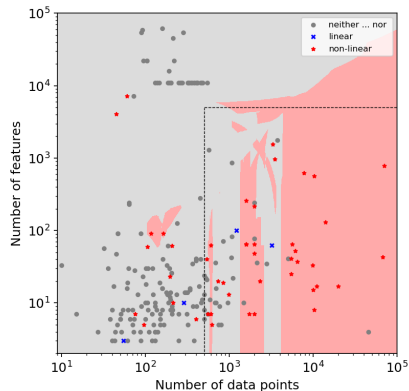


SVM

Linear vs. Non-Linear



SVM



Neural Networks

Linear vs. Non-Linear

- Non-linear models are almost exclusively better than linear models (as expected)
- Statistical equivalent for low data regimes
- Limitation: Conservative statistical test
- benchmark suites such as OpenML-100 have limitations

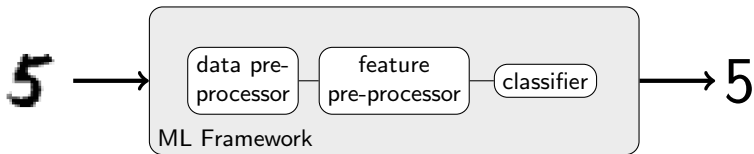
[B. Strang, P. van der Putten, J. N. van Rijn, and F. Hutter. [Don't rule out simple models prematurely: A large scale benchmark comparing linear and non-linear classifiers in openml](#).

In *International Symposium on Intelligent Data Analysis*, pages 303–315. Springer, 2018]

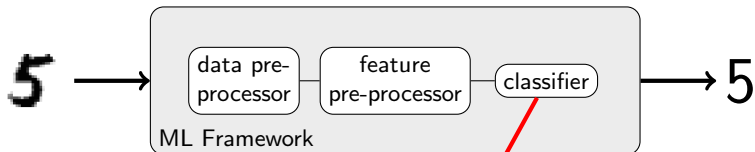


Warm Starting Bayesian Optimization

The Machine Learning Pipeline

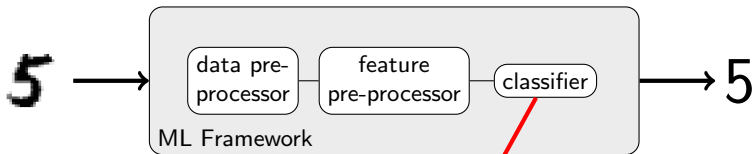


The Machine Learning Pipeline



classifier	# λ
Adaboost	4
Bernoulli Naive Bayes	2
Decision Tree	4
Extra Trees	5
Gradient Boosting	6
k-NN	3
LDA	4
...	

The Machine Learning Pipeline



classifier	# λ
Adaboost	4
Bernoulli Naive Bayes	2
Decision Tree	4
Extra Trees	5
Gradient Boosting	6
k-NN	3
LDA	4
...	



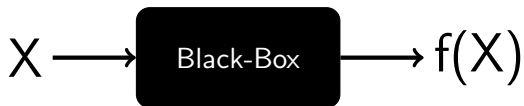
WWW.PHDCOMICS.COM

The Problem Definition

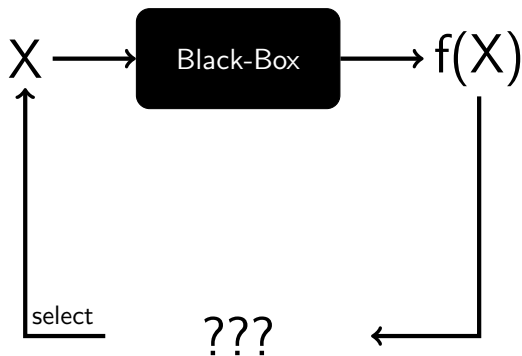
$$A^*, \lambda_* \in \arg \min_{A \in \mathbf{A}, \lambda \in \Lambda} L(A_\lambda, D_{train}, D_{valid})$$

Given a dataset D , find an algorithm A^* and its hyperparameters λ_* that minimizes a given loss function L

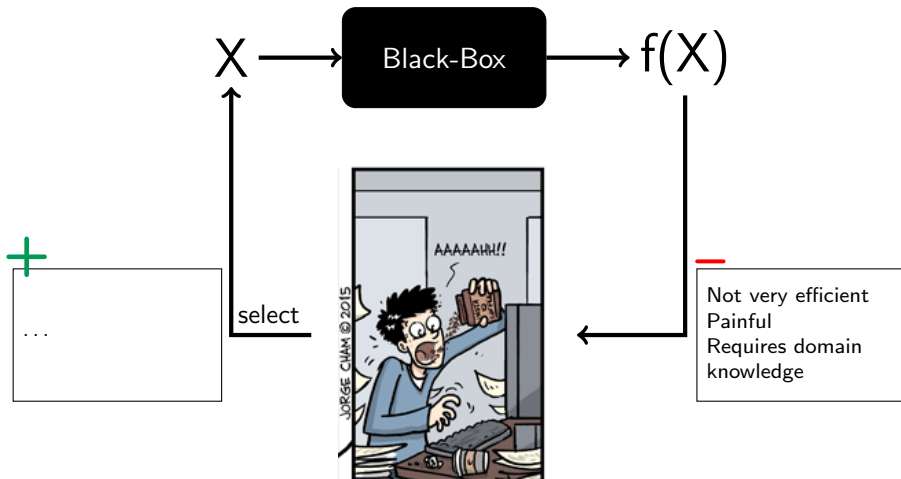
Black-Box Optimization



The Loop



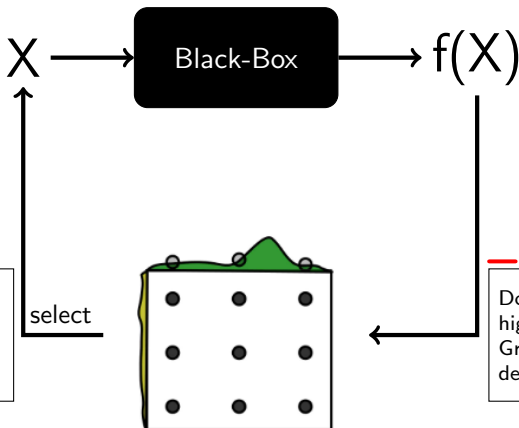
The Human in the Loop



Grid Search



WWW.PHDCOMICS.COM



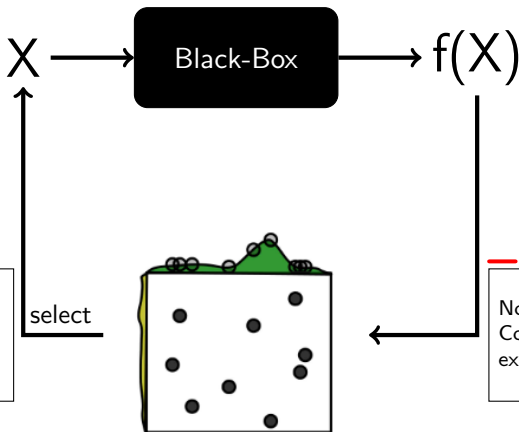
Very simple approach
Can be used to study the problem

Does not scale to high dimensions
Grid needs to be defined

Random Search



WWW.PHDCOMICS.COM



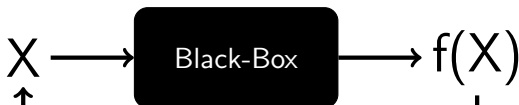
Even more simple
Easily parallelizable
Eventually converges to optimum

Not data efficient
Computationally expensive

Bayesian Optimization

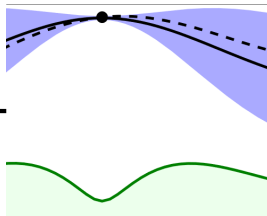


WWW.PHDCOMICS.COM



Data efficient
State of the art

select



Not easy to paral-
lelize

Requirements

- sklearnbot library
 - <https://github.com/openml/sklearn-bot>
- openmlcontrib library
 - Connects openml library with pandas and ConfigSpace library
 - <https://github.com/openml/openml-python-contrib>
- Both not on PyPI, installable from git
- Under active development (contributors welcome)

Plot Results

```
1 import sklearnbot # requires openml and openmlcontib and ConfigSpace
2
3 # given a study_id (int)
4 tasks = openml.study.get_study(study_id, 'tasks').tasks
5
6 configuration_space = sklearnbot.config_spaces.get_config_space('svc', None)
7 output_dir = os.path.join(args.output_dir, args.classifier_name)
8
9 n_executions = 10
10 for i in range(n_executions):
11     task_id = random.choice(tasks)
12     success, run_id, folder = sklearnbot.bot.run_bot_on_task(task_id,
13                                                                configuration_space,
14                                                                output_dir,
15                                                                args.upload_result)
```

Obtain Results

```
1 import openmlcontrib
2 import sklearnbot # requires openml and openmlcontrib and ConfigSpace
3
4 # given a study_id (int)
5 suite = openml.study.get_study(study_id, 'tasks')
6 configuration_space = sklearnbot.config_spaces.get_config_space('svc', None)
7 # given a flow_id (int)
8 performance_data = openmlcontrib.meta.get_tasks_result_as_dataframe(
9     task_ids=suite.tasks,
10     flow_id=flow_id,
11     num_runs=100,
12     per_fold=False,
13     configuration_space=config_space,
14     evaluation_measures='predictive_accuracy',
15     normalize=False
16 )
17
18 meta_features = openmlcontrib.meta.get_tasks_qualities_as_dataframe(
19     study.tasks, False, -1, True)
20 setup_data_with_meta_features = performance_data.join(
21     meta_features, on='task_id', how='inner')
```


Auto-sklearn

```
1 import autosklearn.classification
2 import sklearn.model_selection
3 import sklearn.datasets
4 import sklearn.metrics
5
6 # Load data and split into train and test data
7 X, y = sklearn.datasets.load_digits(return_X_y=True)
8 X_train, X_test, y_train, y_test = \
9     sklearn.model_selection.train_test_split(X, y)
10 # Let auto-sklearn run for 1h
11 autosklearn.classification.AutoSklearnClassifier(
12     time_left_for_this_task=3600)
13 automl.fit(X_train, y_train)
14 y_hat = automl.predict(X_test)
15 print("Accuracy score",
16       sklearn.metrics.accuracy_score(y_test, y_hat))
```

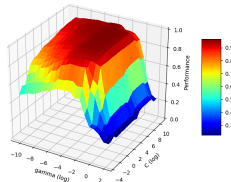
[M. Feurer, A. Klein, K. Eggenberger, J. Springenberg, M. Blum, and F. Hutter. [Efficient and robust automated machine learning](#).

In *Advances in neural information processing systems*, pages 2962–2970, 2015]



Hyperparameter Importance

- Typically unanswered questions:
 - What is a good Configuration Space?
 - What are important hyperparameters?
 - How to sample over these?
- Hyperparameter Importance using OpenML:
https://www.youtube.com/watch?v=mS4vL7_rSWQ



[J. N. van Rijn and F. Hutter. [Hyperparameter importance across datasets](#).

In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 2367–2376. ACM, 2018]

What have we learned

BENCHMARKING



suites

effect of a component

warm starting AutoML

- `openml`, `openmlcontrib` and `sklearnbot`
- Also available on Java and R
- Under active development, contributors welcome :)
- This afternoon: Hands on session

References

- B. Bischl, G. Casalicchio, M. Feurer, F. Hutter, M. Lang, R. G. Mantovani, J. N. van Rijn, and J. Vanschoren. Openml benchmarking suites and the openml100. *arXiv preprint arXiv:1708.03731*, 2017.
- M. Feurer, A. Klein, K. Eggenberger, J. Springenberg, M. Blum, and F. Hutter. Efficient and robust automated machine learning. In *Advances in neural information processing systems*, pages 2962–2970, 2015.
- M. J. Post, P. van der Putten, and J. N. van Rijn. Does Feature Selection Improve Classification? A Large Scale Experiment in OpenML. In *Advances in Intelligent Data Analysis XV*, pages 158–170. Springer, 2016.
- B. Strang, P. van der Putten, J. N. van Rijn, and F. Hutter. Don't rule out simple models prematurely: A large scale benchmark comparing linear and non-linear classifiers in openml. In *International Symposium on Intelligent Data Analysis*, pages 303–315. Springer, 2018.
- J. N. van Rijn and F. Hutter. Hyperparameter importance across datasets. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 2367–2376. ACM, 2018.