

Tools that help you get your experiments under control

Katharina Rasch
PyData Berlin 2019

Reproducibility in theory and in practice

Katharina Rasch
PyData Berlin 2019

Katharina Rasch

`kat@krasch.io`

<https://github.com/krasch/presentations>

PhD Computer Science

Previously: Data science / Computer Vision at zalando

Now: Freelance data science + teaching





20190921T1114.
json



20190921T1114.
pkl



20190921T1803.
json



20190921T1803.
pkl



20190922T1117.
json



20190924T1815.
json



20190924T1815.
pkl



20190926T1733.
json



20190926T1733.
pkl



20190927T1043.
json



20190927T1043.
pkl



20190927T1344.
json



20190927T1344.
pkl



20190927T1346.
json



20190927T1346.
pkl



20190927T1411.
json



20190927T1411.
pkl

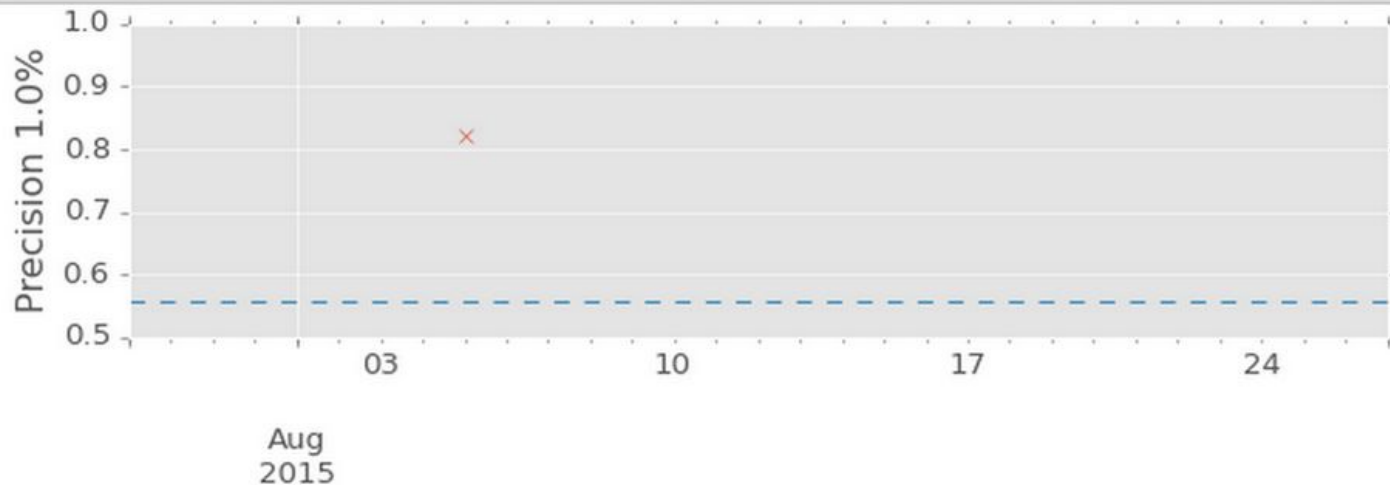


20190927T1607.
json



20190927T1607.
pkl

File Edit View Insert Format Styles Sheet Data Tools Window Help									
Liberation Sans 10									
A32									
	A	B	C	D	E	F	G	H	I
1	Timestamp	Experiment name	Git commit hash	Dataset	Model type	<u>Hyperparameters</u>	Precision	Recall	<u>AUC ROC</u>
2									
3									
4									
5									
6									
7									
8									
9									
10									
11									
12									
13									
14									
15									
16									
17									
18									
19									
20									
21									
22									
23									
24									
25									
26									
27									
28									
29									



Time	Features	Model	Precision 1.0%
------	----------	-------	----------------

05 Aug 2015

12:49	all	RandomForest	0.82
-------	-----	--------------	------

```
{'max_features': 'auto',
'n_estimators': 50, 'depth': 3,
'criterion': 'gini'}
```

12:52	all	RandomForest	0.82
-------	-----	--------------	------

```
{'max_features': 'auto'
```

1. Motivation

Pride

Worry

Trust

Teamwork

Reproducibility

1. Motivation
2. Scope
3. Theory
4. Practice

IMHO

2. Scope

Model exploration / development

many models, features, parameters
(mostly) fixed dataset

errors expected to happen

flexibility needed

Model usage / deployment (e.g. nightly training on new data)

fixed model type, features, parameters
frequent data updates

should be reliable, monitored

rigidity helpful

Model exploration / development

many models, features, parameters
(mostly) fixed dataset

errors expected to happen

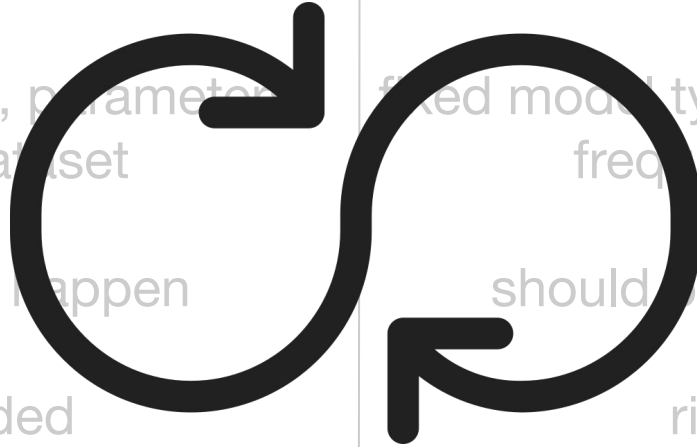
flexibility needed

Model usage / deployment (e.g. nightly training on new data)

fixed model type, features, parameters
frequent data updates

should be reliable, monitored

rigidity helpful



Model exploration / development

many models, features, parameters
(mostly) fixed dataset

errors expected to happen

flexibility needed

--> today

Model usage / deployment (e.g. nightly training on new data)

fixed model type, features, parameters
frequent data updates

should be reliable, monitored

rigidity helpful

--> not today

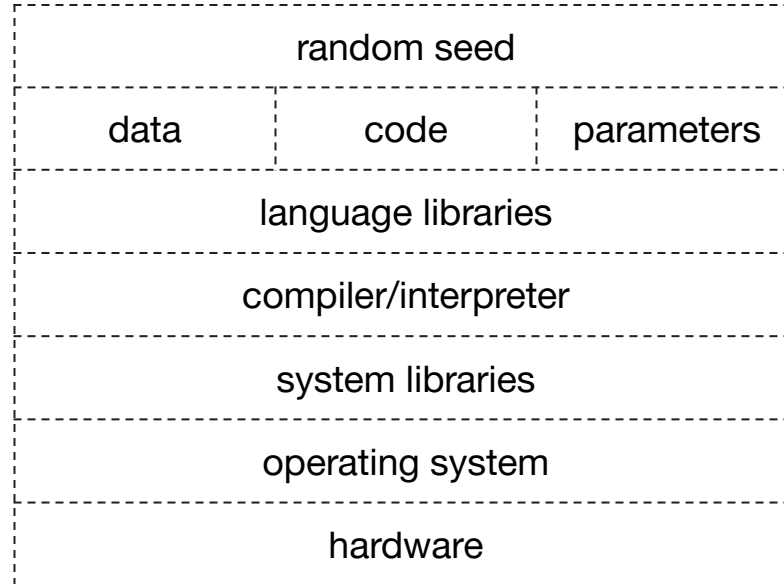
Reproducibility vs. Replicability:
A Brief History of a Confused Terminology
HE Plesser, Frontiers in Neuroinformatics, 2018

This talk:

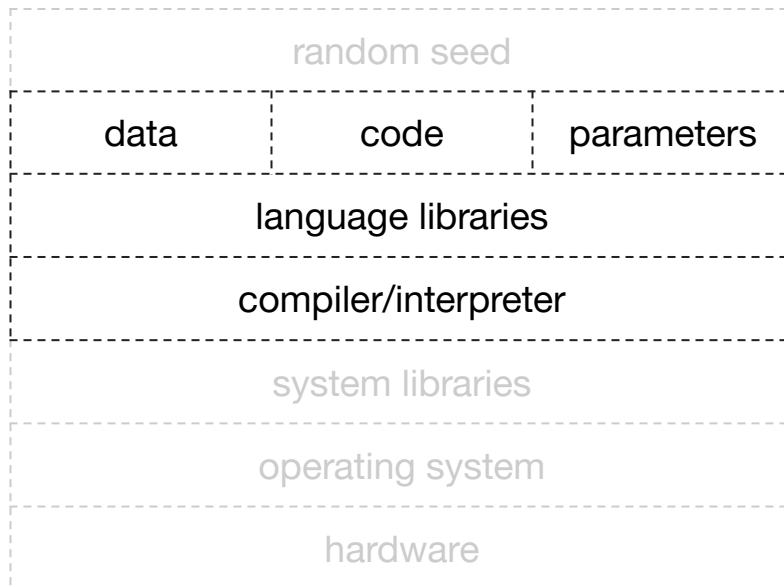
reproducibility = ability to obtain the
same* model with reasonable* efforts

3. Theory

Exact same model



Same* model



Provenance (?) (can compare models)

data hash / version	code hash / version	parameter names and values
language libraries name and version		
compiler/interpreter name and version		

Reproducibility (can reproduce models)

actual data	actual code	parameter names and values
actual language libraries		
actual compiler/interpreter		

Meta		Provenance				Metrics	
Time-stamp ▲▼	Data version ▲▼	Code version ▲▼	Model type ▲▼	Over-sample? ▲▼	Hyper-params	Precision ▲▼	Recall ▲▼

Meta		Provenance					Metrics	
Time-stamp ▲▼	Name	Data version ▲▼	Code version ▲▼	Model type ▲▼	Over-sample? ▲▼	Hyper-params	Precision ▲▼	Recall ▲▼
	test							
	test							
	baseline							
	baseline_fixed							

Meta		Provenance				Metrics	
Time-stamp ▲▼	Data version ▲▼	Code version ▲▼	Model type ▲▼	Over-sample? ▲▼	Hyper-params	Precision ▲▼	Recall ▲▼

Parameters							
Time-stamp ▲▼	Data version ▲▼	Code version ▲▼	Model type ▲▼	Over-sample? ▲▼	Hyper-params	Precision ▲▼	Recall ▲▼

Project-dependent

Time-stamp ▲▼	Data version ▲▼	Code version ▲▼	Model type ▲▼	Over-sample? ▲▼	Hyper-params	Precision ▲▼	Recall ▲▼

Provenance

Reproducibility

data hash / version	code hash / version	parameter names and values	actual data	actual code	parameter names and values
language libraries name and version			actual language libraries		
compiler/interpreter name and version			actual compiler/interpreter		



Out-of-sync error

Provenance

Git commit hash



data hash / version	code hash / version	parameter names and values
language libraries name and version		
compiler/interpreter name and version		

Reproducibility

Developer's local source code directory



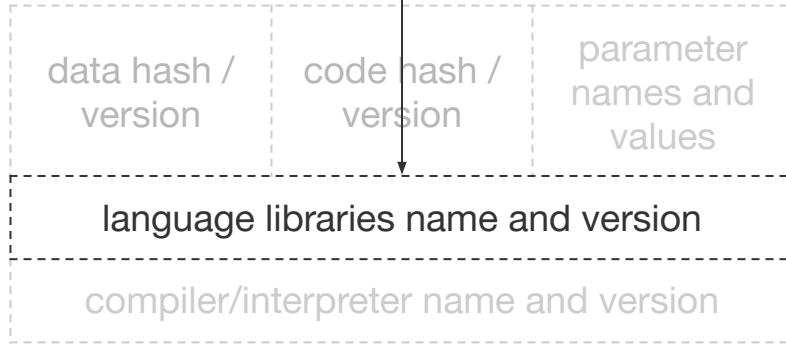
actual data	actual code	parameter names and values
actual language libraries		
actual compiler/interpreter		



Forgot to commit before starting experiment

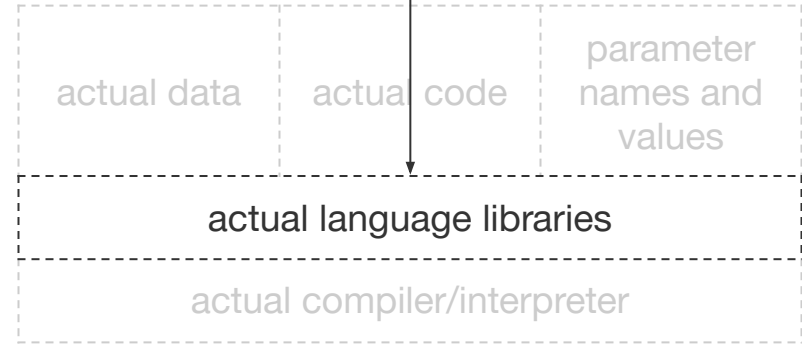
Provenance

requirements.txt
conda.yaml
Dockerfile



Reproducibility

Developer's local environment



Forgot to update requirements file after installing new library

Provenance

Hash calculated over data file(s)



Reproducibility

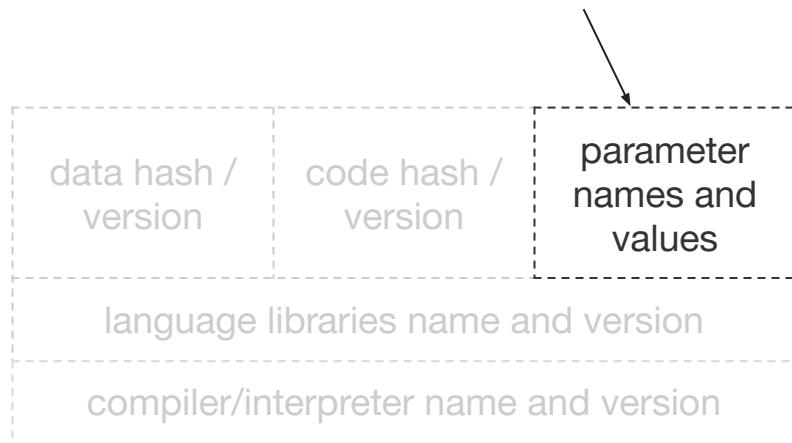
File on developer machine, S3, etc



Altered data without calculating new hash

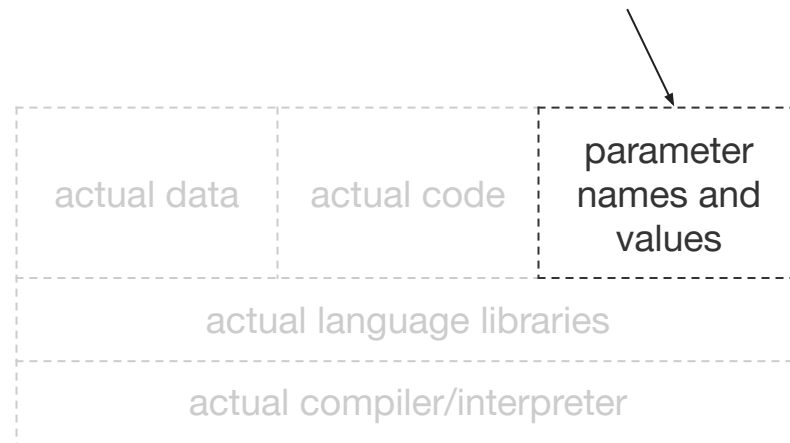
Provenance

command line parameters, config file, config in code



Reproducibility

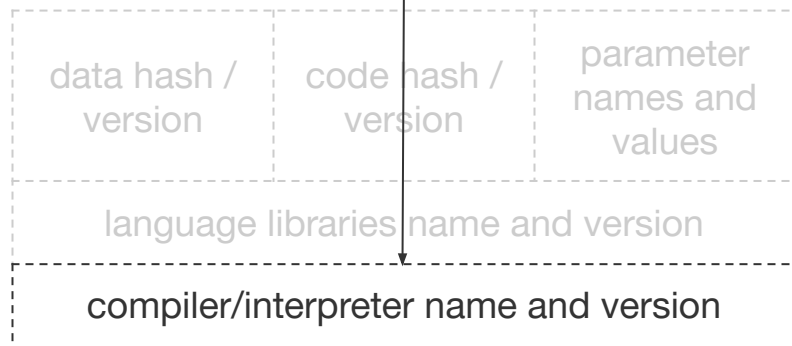
command line parameters, config file, config in code



Not actually passing parameters to model -> code review, tests

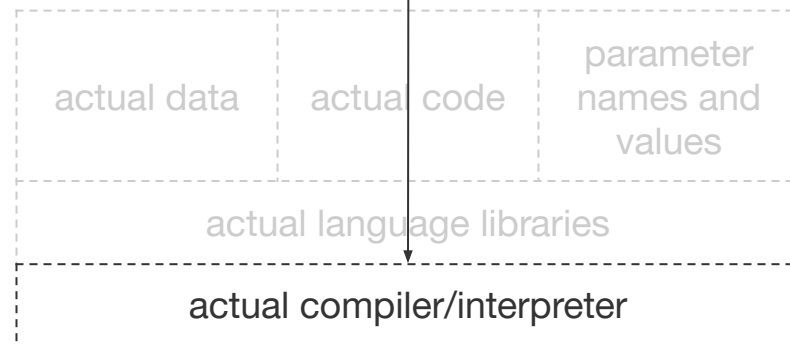
Provenance

Readme.md
Dockerfile



Reproducibility

Developer's local python installation



Using wrong python version locally

For every experiment

- ensure using correct data

- ensure local git repo clean

- ensure requirements up-to-date

- ensure using correct python version

Too much responsibility!

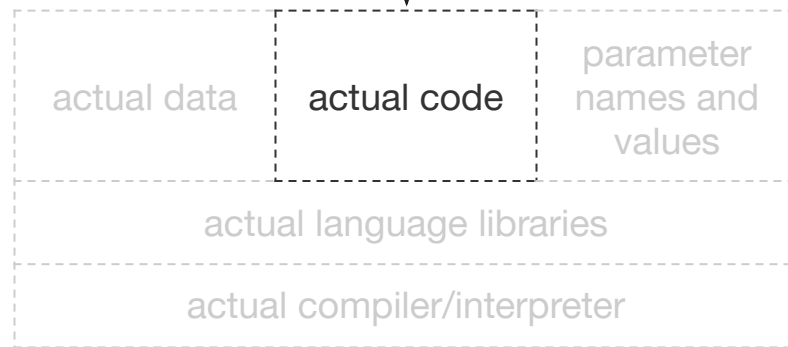
Too much hassle!

Provenance



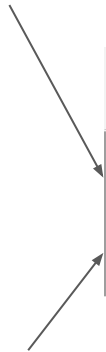
Reproducibility

upload zip archive of local directory



Duplicating git functionality
Less easy to compare








Development mode: best-effort reproducibility



Code	Where executed	Who triggers
-/zip	local	user
git	local	user

Reproducibility mode: permanently recording successful experiments

Code	Where executed	Who triggers
-/zip	local	user
git	local	user

<input type="checkbox"/>	#21 Train model 	<u>9 minutes ago</u>	discokugel	<u>46.2254</u>	<u>43.8166</u>	✔ Complete	<u>39 seconds</u>	Baseline	lag_one_week
<input type="checkbox"/>	#18 Train model  	<u>38 minutes ago</u>	discokugel			❗ Error	<u>52 seconds</u>	Baseline	lag_one_week
<input type="checkbox"/>	#17 Train model  	<u>38 minutes ago</u>	discokugel			❗ Error	<u>5 seconds</u>	LinearRegression	all_lag_date
<input type="checkbox"/>	#16 Train model  	<u>42 minutes ago</u>	discokugel	<u>37.6169</u>	<u>36.8913</u>	✔ Complete	<u>6 seconds</u>	LinearRegression	lag_features

Screenshot from valohai.com

Code	Where executed	Who triggers
-/zip	local	user
git	local	user
git	remote	user
git	remote	git push / merge request

Open krasch wants to merge 2 commits into master from random_forest

Conversation 0 Commits 2 Checks 0 Files changed 3 +33 -10

krasch commented 40 seconds ago
No description provided.

krasch added 2 commits 13 minutes ago

- Added random forest be87b40
- Plot prediction 3cd8dc8

Add more commits by pushing to the random_forest branch on krasch/gradient-bikes.

All checks have failed
1 errored check

gradientci

This branch has no conflicts with the base branch
Merging can be performed automatically.

Merge pull request or view command line instructions.

Write Preview AA B i “ < > ☁ ☰ ☷ ☸ @ 📎 ↶

Leave a comment

Attach files by dragging & dropping, selecting or pasting them.

Close pull request Comment

Reviewers
No reviews

Assignees
No one—assign yourself

Labels
None yet

Projects
None yet

Milestone
No milestone

Notifications
Customize
Unsubscribe
You're receiving notifications because you're watching this repository.

1 participant

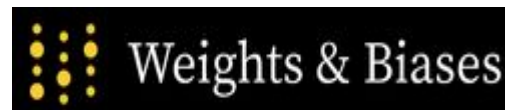
Lock conversation

Screenshot from github.com, using <https://gradient.paperspace.com/> CI

How do you want to work?

4. Practice

Software as a Service solutions



Open source tools

 MODELCHIMP

mlflow

DEPLOY-ML

Pachyderm

Sumatra

Sacred

FGLab

 Polyaxon

VisTrails

DVC

noWorkflow

 StreamSets

Open source tools



Pachyderm

Sacred



Open source tools



Sacred
[MIT license]



Sacred

```
from sacred import Experiment
ex = Experiment("bikes")

ex.add_resource("data/hourly_counts.csv")

ex.log_scalar("train_rmse", train_rmse)
ex.add_artifact("results/predictions.png")
```

```
@ex.named_config
def baseline():
    split = "2018-01-01"
    feature_functions = [features.lag_one_week]
```

```
@ex.automain
def run_experiment(split, feature_functions):
    # do data science stuff
```

```
> python train.py with baseline
```

Status: 7 selected

Filters:

Column Name

==

Enter Value...

Add Filter

Compare

+/- Metric Columns

Show/Hide Columns

Auto Refresh

Last Update: September 26th, 09:24:59 PM

Reload

24 experiments

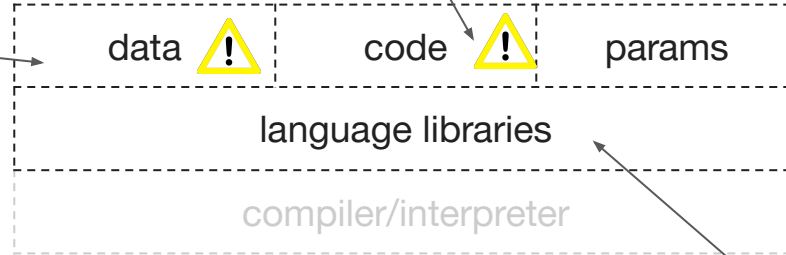
		Id	Experiment Name	Start Time	Status	Notes	Resources	Duration	Result	Stop Time	Seed	RMSE Train	RMSE Test
<input type="checkbox"/>	▶	24	● bikes	2019-09-26T16:44:44	FAILED	Enter Notes	["/home/arbe/tembrace-/data/hourly_counts.csv"]	356ms		2019-09-26T16:44:44	965283310		
<input type="checkbox"/>	▶	23	● bikes	2019-09-26T16:44:14	FAILED	Enter Notes	[]	18ms		2019-09-26T16:44:14	549376400		
<input type="checkbox"/>	▶	21	● bikes	2019-09-22T22:27:44	COMPLETED	Enter Notes	["/home/arbe/tembrace-/data/hourly_counts.csv"]	53.9s		2019-09-22T22:28:38	152215597	107.66954716456115	107.66954716456115
<input type="checkbox"/>	▶	20	● bikes	2019-09-22T22:11:34	COMPLETED	Enter Notes	["/home/arbe/tembrace-/data/hourly_counts.csv"]	1.7s		2019-09-22T22:11:36	403214227	34.45518373155321	34.45518373155321
<input type="checkbox"/>	▶	19	● bikes	2019-09-22T22:11:21	COMPLETED	Enter Notes	["/home/arbe/tembrace-/data/hourly_counts.csv"]	1.7s		2019-09-22T22:11:23	9631558	34.24920057236309	34.24920057236309
<input type="checkbox"/>	▶	18	● bikes	2019-09-22T22:11:07	FAILED	Enter Notes	[]	12ms		2019-09-22T22:11:07	209327342		
<input type="checkbox"/>	▶	17	● bikes	2019-09-22T22:10:52	FAILED	Enter Notes	[]	15ms		2019-09-22T22:10:52	66310965		
<input type="checkbox"/>	▶	16	● bikes	2019-09-22T21:13:56	COMPLETED	Enter Notes	["/home/arbe/tembrace-/data/hourly_counts.csv"]	1.6s		2019-09-22T21:13:58	136782595	34.42437733721705	34.42437733721705
<input type="checkbox"/>	▶	15	● bikes	2019-09-22T21:13:41	COMPLETED	Enter Notes	["/home/arbe/tembrace-/data/hourly_counts.csv"]			2019-09-22T21:13:44	178673814		
<input type="checkbox"/>	▶	14	● bikes	2019-09-22T21:05:27	COMPLETED	Enter Notes	["/home/arbe/tembrace-/data/hourly_counts.csv"]	1.6s		2019-09-22T21:05:28	241019292	34.31922926598029	34.31922926598029

Code	Where executed	Who triggers
zip	local	user
git (with local changes, flagged as “dirty” in UI)	local	user
git (no local changes)	local	user
git	remote	user
git	remote	git push / merge request

always: zip uploaded

optional: reproducibility mode, must remember to use `--enforce_clean`

md5 hash
also stores data??
how to get data out ??



auto-detects
used libraries

Steps to actually reproduce

1. Obtain the data (somehow)
2. Obtain the code
 - a. download zip from UI
 - b. git checkout (if clean commit)
3. Find list of libraries in UI and install
4. Find config values in UI
5. `> python train.py with param1=value1 ...`



```
with mlflow.start_run():  
    mlflow.log_param("model", args.model)  
    mlflow.log_param("features", args.features)  
  
    # do data science stuff  
    ...  
  
    mlflow.log_metric("RMSE_train", rmse_train)  
    mlflow.log_artifact("results/predictions.png")
```

conda.yaml

```
name: bikes
channels:
  - defaults
dependencies:
  - scikit-learn=0.21.1
  - matplotlib
  - keras
  - pip:
    - mlflow
```

```
> mlflow run -e hyper_train .
```

```
> mlflow run git@gitlab.com:org/mlflow_bikes.git
```

MLProject

```
name: bikes

conda_env: conda.yaml
entry_points:

  hyper_train:
    parameters:
      max_depth: int
      n_estimators: int
      bootstrap: boolean
    command: "python -m
hyper.train ..."
```

Experiments



Default

Default

Experiment ID: 0

Artifact Location: jmlruns/0

Description:

Search Runs: metrics.rmse < 1 and params.model = "tree"



State:

Active ▾

Search

Filter Params: alpha, lr

Filter Metrics: rmse, r2

Clear

Showing 36 matching runs

Compare

Delete

Download CSV

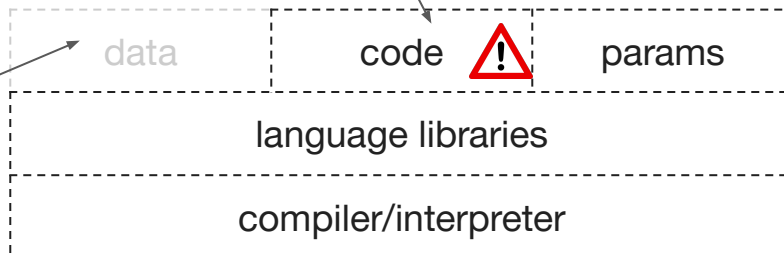
<input type="checkbox"/>	Date	User	Run Name	Source	Version	Tags	model	Parameters	Metrics
<input type="checkbox"/>	2019-09-20 11:06:03	arbeit		bikes	eaac67		LinearRegression	features: all_lag_date_w... log: true	RMSE_test: 35.6768504256... RMSE_train: 36.5240951492...
<input type="checkbox"/>	2019-09-20 11:05:43	kat		bikes	eaac67		LinearRegression	features: all_lag_date_w... log: true	
<input type="checkbox"/>	2019-09-20 11:05:17	kat		bikes	eaac67		LinearRegression	features: all_lag_date_w... log: true	
<input type="checkbox"/>	2019-09-20 11:02:05	kat		bikes	eaac67		LinearRegression	features: all_lag_date_w... log: true	
<input type="checkbox"/>	2019-09-20 11:01:43	kat		bikes	eaac67		LinearRegression	features: all_lag_date_w... log: true	RMSE_test: 35.6768504256... RMSE_train: 36.5240951492...
<input type="checkbox"/>	2019-09-20 10:59:15	kat		bikes	eaac67		LinearRegression	features: all_lag_date_w... log: true	RMSE_test: 35.6768504256... RMSE_train: 36.5240951492...
<input type="checkbox"/>	2019-09-20 10:59:01	kat		bikes	eaac67		LinearRegression	features: all_lag_date_w... log: true	
<input type="checkbox"/>	2019-09-20 10:56:54	kat		mlflow ...	eaac67		LinearRegression	features: all_lag_date_w... log: true	RMSE_test: 35.6768504256... RMSE_train: 36.5240951492...
<input type="checkbox"/>	2019-09-20 10:55:46	arbeit		mlflow ...	eaac67		LinearRegression	features: all_lag_date_w... log: true	RMSE_test: 35.6768504256... RMSE_train: 36.5240951492...
<input type="checkbox"/>	2019-09-20 10:53:39	arbeit		mlflow ...	eaac67		LinearRegression	features: all_lag_date_w... log: true	RMSE_test: 35.6768504256... RMSE_train: 36.5240951492...

Code	Where executed	Who triggers
-	local	user
git (with local changes, NOT flagged as dirty in UI)	local	user
git (no local changes)	local	user
git	remote	user
git	remote	git push / merge request

```
> mlflow run .
```

does not check if local git repo clean, records stale hash

not mentioned in
documentation



Steps to actually reproduce

(if have been careful with `mlflow run .`)

1. Obtain the data (somehow)
2. `> mlflow run -v commit_hash git@repo...`



// add data and push to remote storage (e.g. S3)

```
> dvc add hourly_counts.csv  
> dvc push
```

// set up pipeline

```
> dvc run -d hourly_counts.csv \  
          -d train.py \  
          -M results/rmse.json \  
          python train.py
```

... change data or train.py

// re-run affected pipeline steps

```
> dvc repro
```

git tags

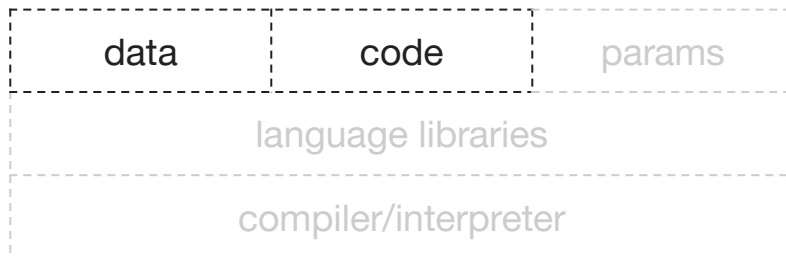
```
> dvc metrics show -T
```

```
baseline:  
  results/rmse.json: {"RMSE (train)": 46.22540681588555, "RMSE (test)": 43.81664782510664}  
linear_regression:  
  results/rmse.json: {"RMSE (train)": 37.616881728555924, "RMSE (test)": 36.89125150760678}  
linear_regression_with_lag_date_weather:  
  results/rmse.json: {"RMSE (train)": 36.52409514921099, "RMSE (test)": 35.67685042569101}
```

Code	Where executed	Who triggers
-/zip	local	user
git (with local changes)	local	user
git (no local changes)	local	user
git	remote	user
git	remote	git push / merge request



Need to ensure pipeline in-sync with reality



Need to ensure `dvc repro` is run before committing (commit hook?)

Steps to actually reproduce

1. `> git clone ...`
2. `> dvc pull`
3. `> dvc repro -f`

How do you want to work?

Envy

~~Envy~~ Inspiration


How do we want to work?

How do we want to work?

<https://github.com/krasch/presentations>
kat@krasch.io



From 01/01/2018

To 02/10/2018 

1. Oberbaumbrücke		2,923,788
2. Jannowitzbrücke		2,390,985
3. Berliner Straße		1,592,578

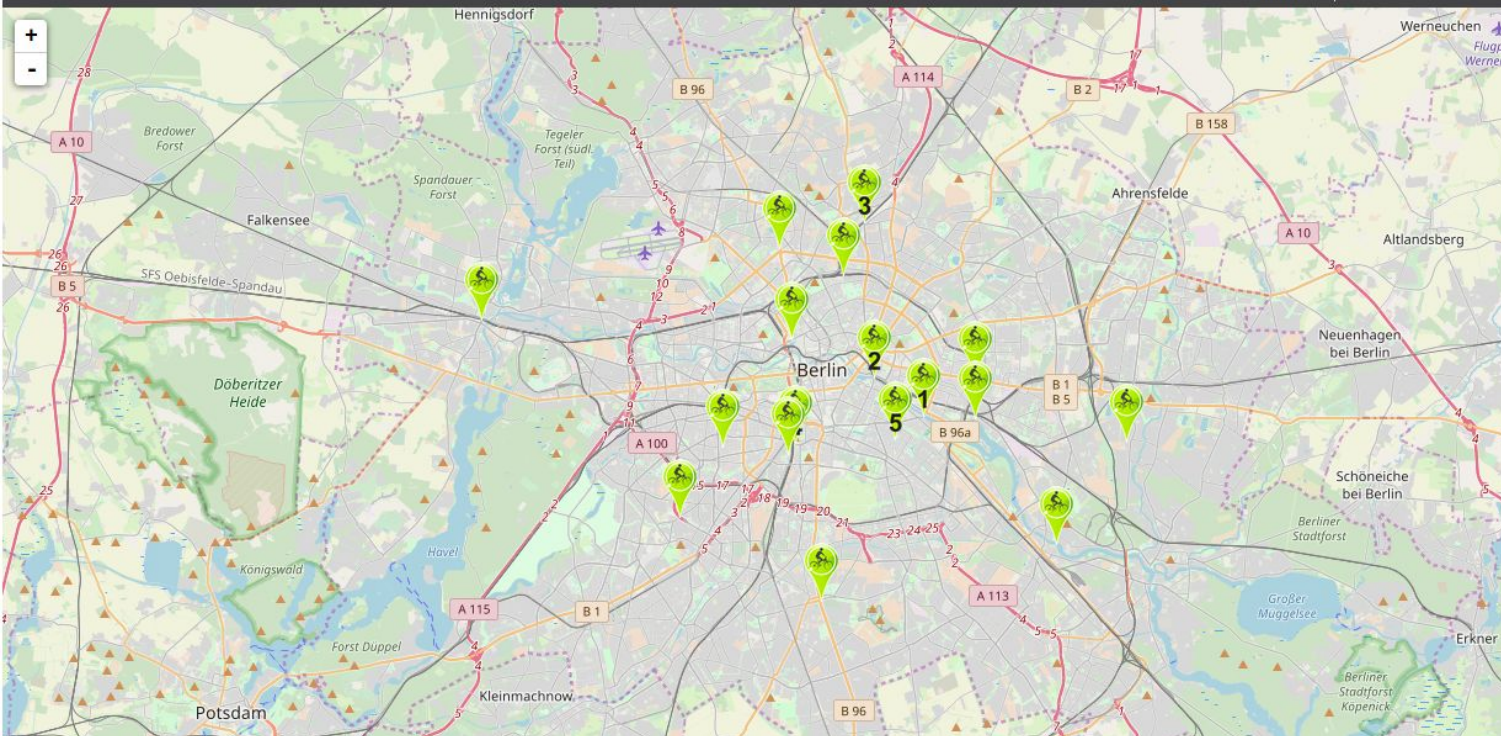
[more...](#)

Total counts:

17,059,860

Daily Average:

3,665



Oberbaumbrücke

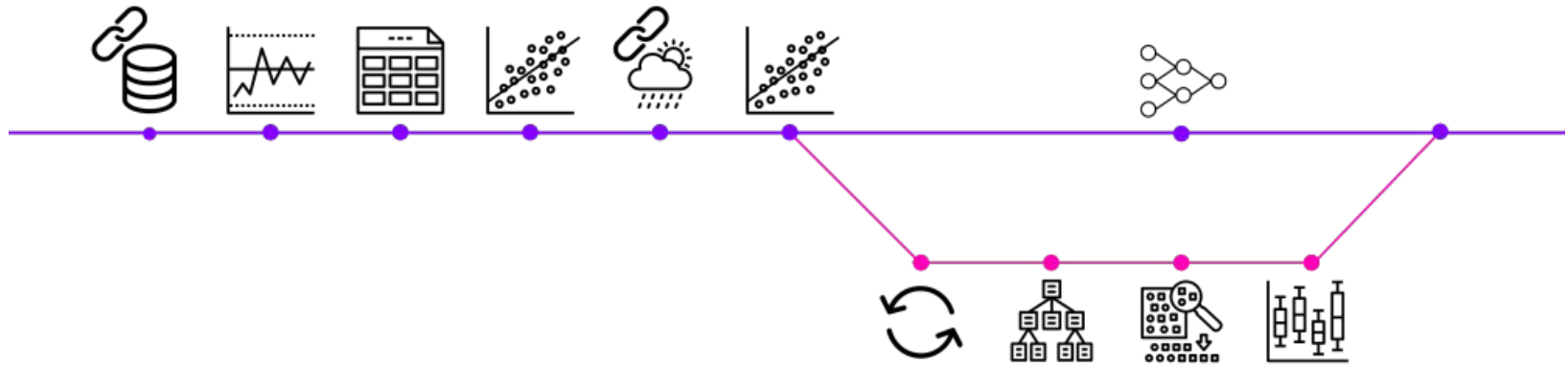


Public Web Page

Yesterday	Daily AVG	Total
11,916	10,671	2,923,788

All the Data


☒ Days
 ☐ Weeks
 ☐ Months

Icons by Vectorstall, Vaibhav Radhakrishnan,
Kimmi Studio, Becris; <https://thenounproject.com>