

NLP games и автоматизация пайплайнов с Kubeflow

Mikhail Rozhkov



web: ml-repa.ru

telegram: t.me/mlrepa

AI-hat game



Игра на загадывание и отгадывание слов.

Player 1:

- пытается объяснить вытянутое из шляпы слово с помощью набора неоднокоренных слов, которые называются по очереди

Other players:

- предлагают несколько вариантов отгадок

Scoring

- **guess score:** чем раньше игрок из Other players угадает вытянутое из шляпы слово, тем больше очков он получит.
- **explain score:** чем больше Other players угадают вытянутое из шляпы слово и чем раньше они это сделают, тем больше очков Player 1

AI Hat competition on Raiffeisen bootcamp 2019

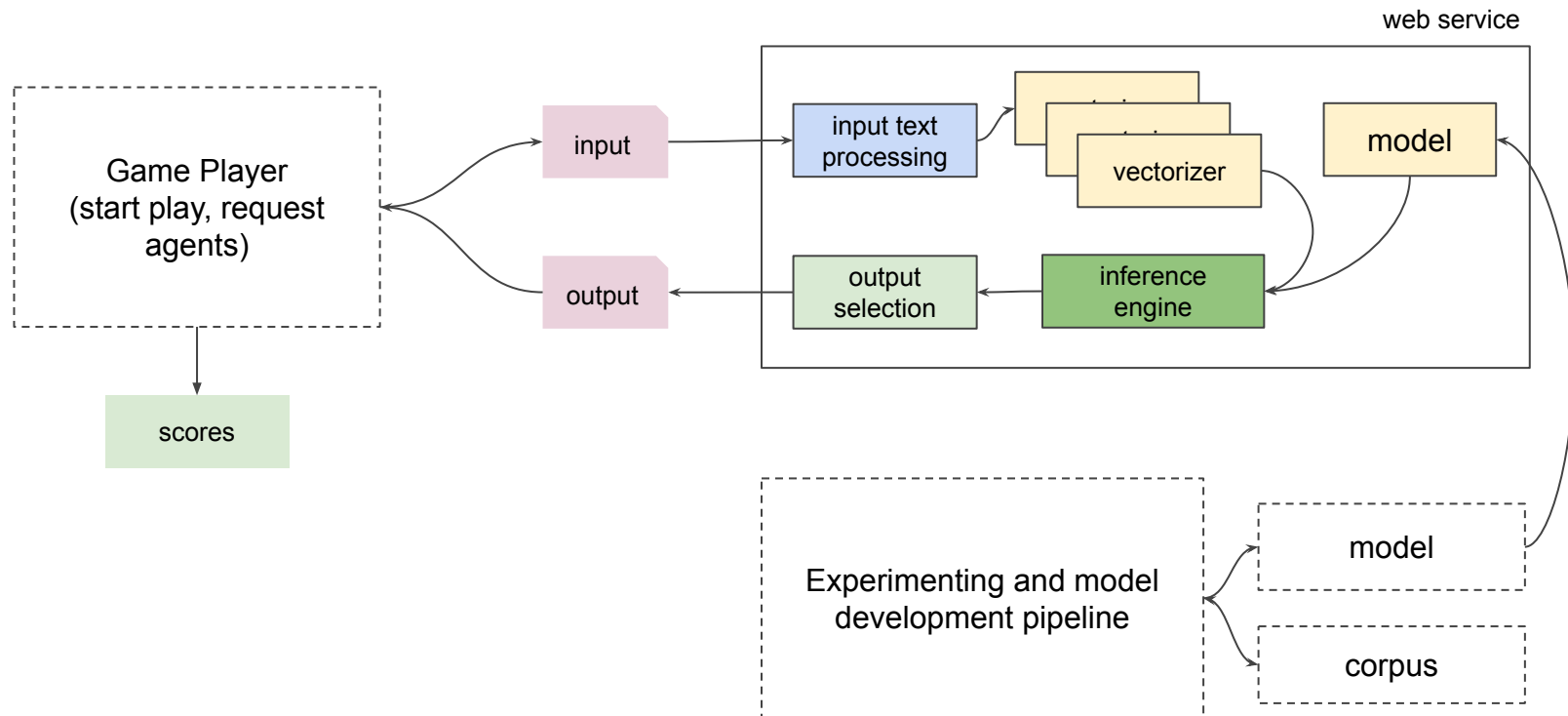
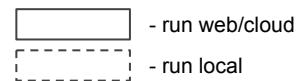
Task

- teams develop their own model, which implements the logic of guessing and explaining
- model serving as a web service with REST API endpoints
 - /guess
 - /explain



Text

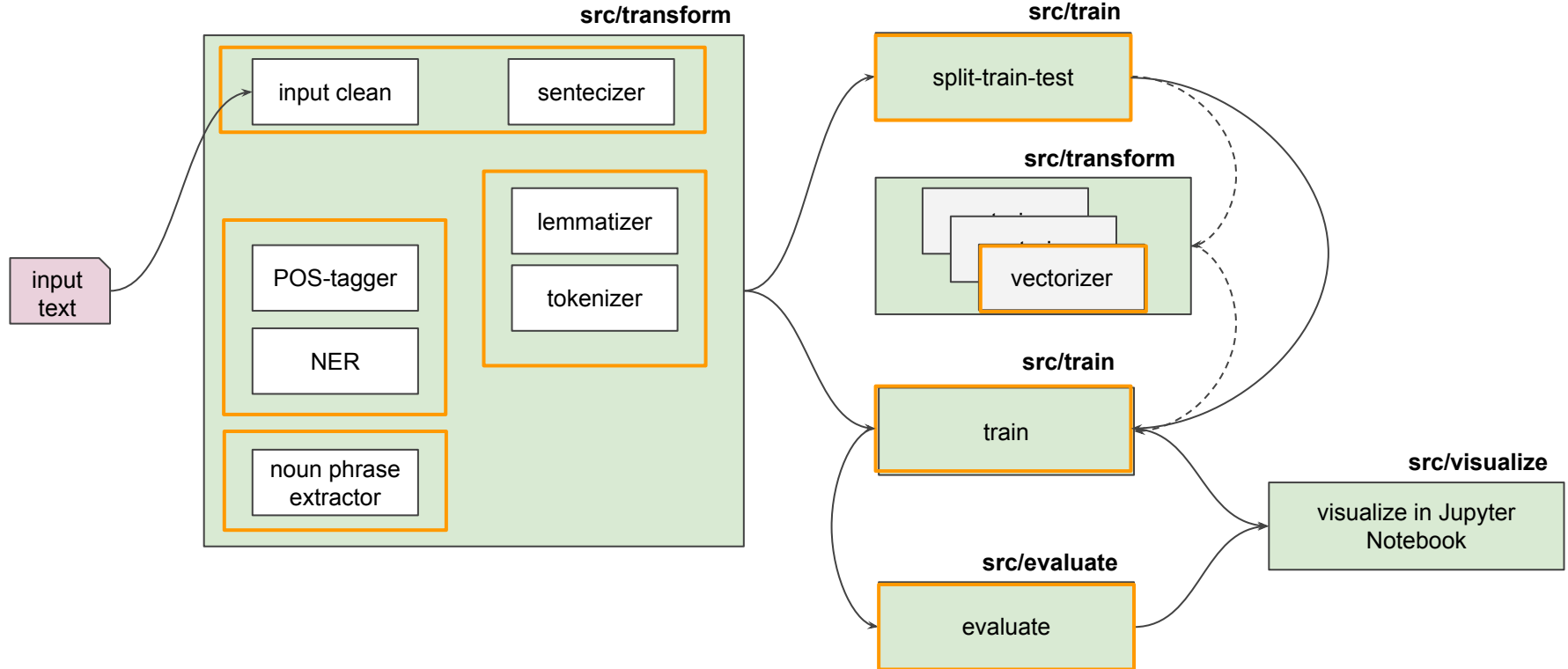
- financial news
- words: financial terms

Game train/serve architecture



Code structure

 - package in B src
 - stage in pipeline



Text processing

Методы:

- 1) Remove symbols and word reduction forms (- -> ' ', 'll -> will and etc.)
- 2) Remove links and tags (www, /, .com and etc.)
- 3) Lower-case
- 4) Lemmatization
- 5) Remove multiple spaces

Train embeddings (fasttext)

CBOW - for guessing

- dim=50, epoch=120, charNgram=(3, 4), wordNgram=3, window=5

SKIPGRAM - for explain

- dim=50, epoch=140, charNgram=(4, 6), wordNgram=3

Evaluation

Metrics

- total score = guess + explain
- separate validation set
- competition among models

Approach 1: dvc & mlflow

- dvc
 - pipelines automation
 - artifacts and models versioning
- mlflow
 - metrics tracking and experiment management

mlflow



Approach 2: kubeflow

- pipeline configuration
- metrics tracking
- artifacts



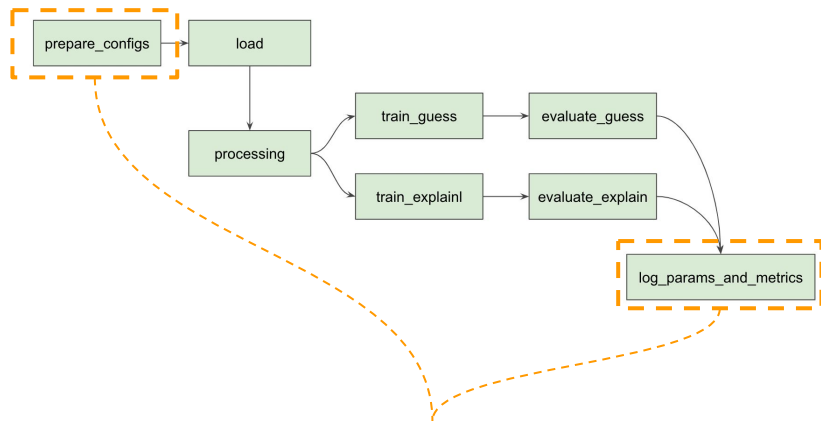
Kubeflow: a platform for building ML products

- *A curated set of compatible tools and artifacts that lays a foundation for running production ML apps*
 - Run containers on Kubernetes cluster
 - Kubernetes runs everywhere
 - Enterprises can adopt shared infrastructure and patterns for ML and non ML services
 - Key features
 - Easy, repeatable, portable deployments on a diverse infrastructure
 - Deploying and managing loosely-coupled microservices
 - Scaling based on demand
- Pipelines
 - Notebooks
 - TensorFlow model training
 - Model serving
 - Multi-framework



Pipeline

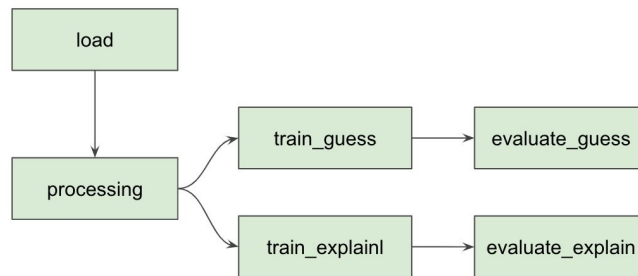
dvc + mlflow



DS решает вопросы с конфигом и логированием метрик/параметров*

* можно не выделять в отдельные этапы

kubeflow



kubeflow предлагает UI для конфига, логирует все inputs/outputs этапов**

** это не всегда удобно

Code structure

dvc + mlflow

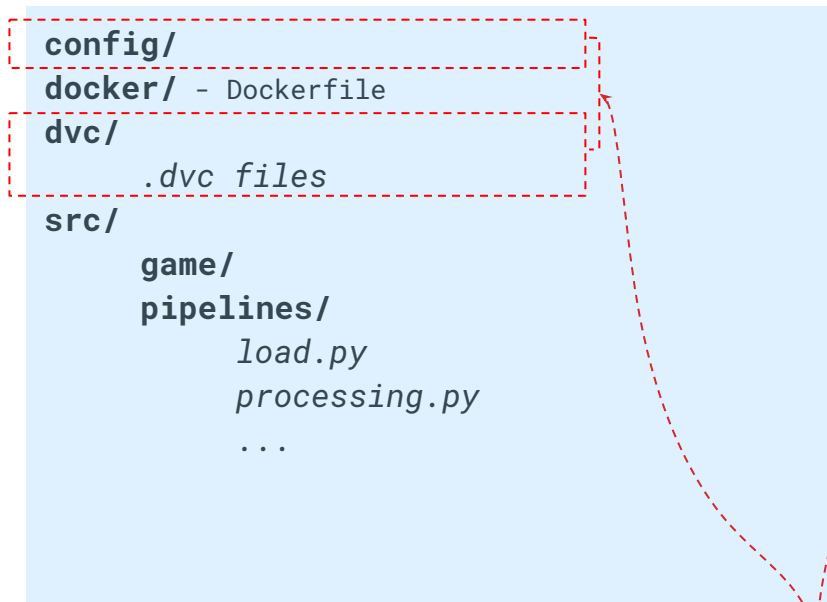
```
config/  
docker/ - Dockerfile  
dvc/  
src/  
    game/  
    pipelines/  
        load.py  
        processing.py  
        ...
```

kubeflow

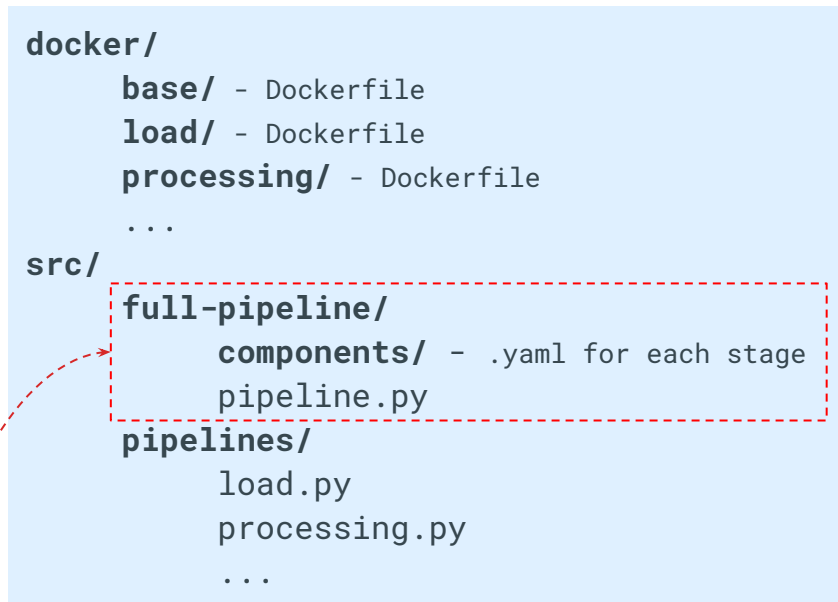
```
docker/  
    base/ - Dockerfile  
    load/ - Dockerfile  
    processing/ - Dockerfile  
    ...  
src/  
    full-pipeline/  
        components/ - .yaml for each stage  
        pipeline.py  
    pipelines/  
        load.py  
        processing.py  
        ...
```

Code structure

dvc + mlflow



kubeflow

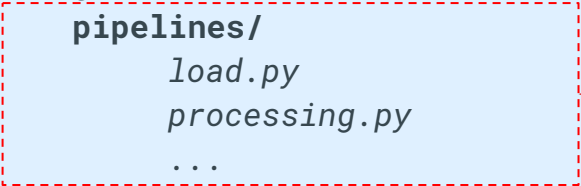


different approach to config pipeline

Code structure

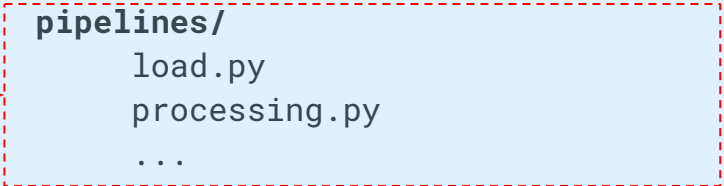
dvc + mlflow

```
config/  
docker/ - Dockerfile  
dvc/  
src/  
  game/  
  pipelines/  
    load.py  
    processing.py  
    ...
```



kubeflow

```
docker/  
  base/ - Dockerfile  
  load/ - Dockerfile  
  processing/ - Dockerfile  
  ...  
src/  
  full-pipeline/  
    components/ - .yaml for each stage  
    pipeline.py  
  pipelines/  
    load.py  
    processing.py  
    ...
```



~ same .py for stages

Find experiments / runs

dvc + mlflow

mlflow

Experiments

- Default
- GuessExperiment_1**
- ExplainExperiment_1

GuessExperiment_1

Experiment ID: 1 Artifact Location: file:///home/ai-hat/mlruns/1

Description: ☒

Search Runs:

Filter Params: Filter Metrics:

Showing 7 matching runs [Compare](#) [Delete](#) [Download CSV](#)

<input type="checkbox"/>	Date	User	Run Name	Source
<input type="checkbox"/>	2019-10-23 19:05:36	user		
<input type="checkbox"/>	2019-10-23 18:41:51	user		

kubeflow

All experiments All runs

Filter experiments

Experiment name	Description	Last 5 runs				
fasttext-two-models	Experiment uses FastText to train models. 2 separate models: CBOW - for ques...					
Run name	Status	Duration	Pipeline	Start time	explain-accuracy	guess-accuracy
<input type="checkbox"/> run6	✓	0:01:59	nlp_1_aihat_game...	10/23/2019, 7:00:...	3.203%	22.222%
<input type="checkbox"/> run5-tune-ngrams-texts	✓	0:01:30	nlp_1_aihat_game...	10/23/2019, 6:52:...	3.400%	17.172%
<input type="checkbox"/> run5-tune-15epoch	✓	0:01:30	nlp_1_aihat_game...	10/23/2019, 6:49:...	3.222%	19.192%
<input type="checkbox"/> run4-tune-15epoch	✓	0:01:31	nlp_1_aihat_game...	10/23/2019, 6:41:...	2.958%	17.172%
<input type="checkbox"/> run3-tune-params	✓	0:01:18	nlp_1_aihat_game...	10/23/2019, 6:38:...	2.714%	8.081%
<input type="checkbox"/> run2-add-texts	✓	0:02:21	nlp_1_aihat_game...	10/23/2019, 6:31:...	1.792%	1.010%
<input type="checkbox"/> run1-baseline	✓	0:01:23	nlp_1_aihat_game...	10/23/2019, 6:26:...	2.643%	0.000%
exp1-alex ✓						
Default All runs created without specifying an experiment will be grouped here.						

list of experiments

New experiment

dvc + mlflow

```
1 1
2 2
3 3
4 4
5 5
6 6
7 7
8 8
9 9
10 10
11 11
12 12
13 13
14 14
15 15
16 16
17 17
18 18
19 19
20 20
21 21
22 22
23 23
24 24
25 25
26 26
27 27
28 28
29 29
30 30
31 31
32 32
33 33
34 34
35 35
36 36
37 37
38 38
39 39
40 40
41 41
42 42
```

```
load:
  raw_data: data/raw/us-financial-news-articles
  docs_number: 100000
  min_words_per_sentence: 4
  cleaned_sentences_file: data/processed/cleaned_sentences.txt

processing:
  corpus_file: data/processed/corpus.txt

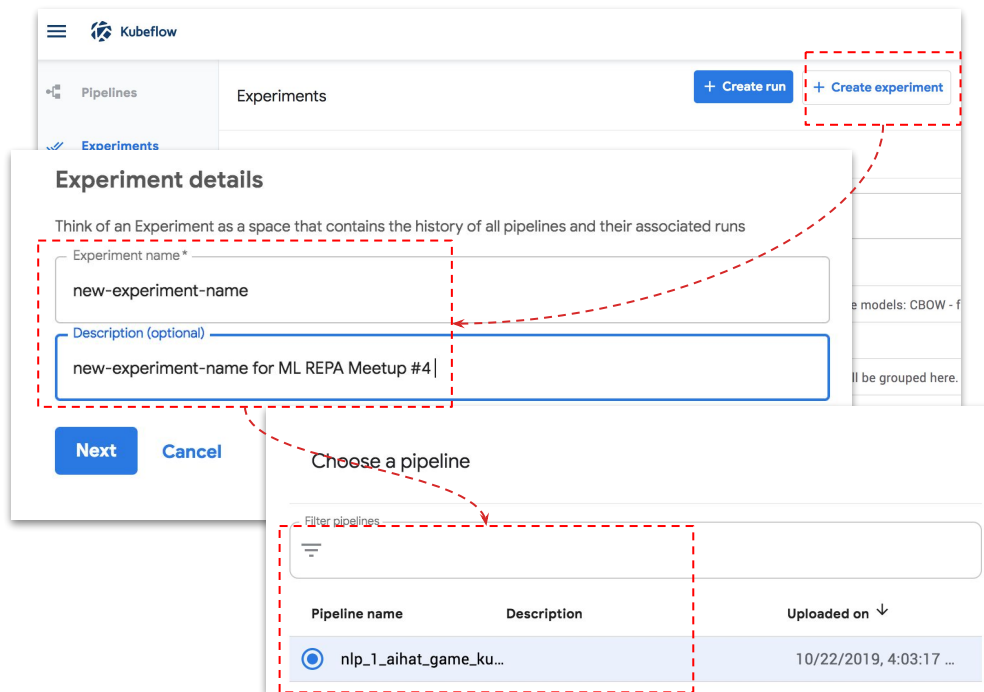
train_guess:
  type: cbow
  dim: 50
  thread: 4
  epoch: 15
  wordNgrams: 1
  model_path: models/guess.model

train_explain:
  type: skipgram
  dim: 50
  thread: 4
  epoch: 15
  wordNgrams: 1
  model_path: models/explain.model

evaluate_guess:
  synonyms_file: data/processed/synonyms.json
  accuracy_n_words: 20
  metrics_file: reports/guess_metrics.json
```

\$ dvc repro last_pipeline_stage.dvc

kubeflow



Run experiment (kubeflow)

The screenshot displays the Kubeflow Experiments interface. At the top, there are buttons for '+ Create run', '+ Create experiment', 'Compare runs', 'Clone run', 'Archive', and 'Refresh'. The 'Clone run' button is highlighted with a red dashed box. Below the buttons, there are tabs for 'All experiments' and 'All runs'. A filter input is present. The main table lists experiments with columns for 'Experiment name', 'Description', 'Run name', 'Status', 'Duration', and 'Pipeline'. The 'fasttext-two-models' experiment is expanded, showing a table of runs. The 'run6' row is highlighted with a red dashed box. A modal titled 'Run parameters' is open, showing a form for specifying pipeline parameters. The 'project' parameter is set to 'vision-230607'. The 'bucket' parameter is set to 'vision-storage1'. The 'raw-data' parameter is set to 'nlp-1-aihat-game-data/short_dataset.zip'. The 'docs-number' parameter is set to '100000'. The 'min-words-per-sentence' parameter is set to '4'. The 'cleaned-sentences-file' parameter is set to 'nlp-1-aihat-game-data/cleaned_sentences_file.txt'. The 'corpus-file' parameter is set to 'nlp-1-aihat-game-data/corpus.txt'. The 'synonyms-file' parameter is set to 'nlp-1-aihat-game-data/synonyms.json'. The 'guess-model-type' parameter is set to 'cbow'. The 'guess-model-dim' parameter is set to '70'. The 'guess-model-thread' parameter is set to '70'. The 'accuracy' parameter is set to '222%'. A red dashed line connects the 'Clone run' button to the 'Run parameters' modal.

Experiments

+ Create run + Create experiment Compare runs Clone run Archive Refresh

All experiments All runs

Filter experiments

Experiment name	Description	Run name	Status	Duration	Pipeline
new-experiment-name	new-experiment-name for ML REPA Me				
fasttext-two-models	Experiment uses FastText to train mode				
		run6	✓	0:01:59	nlp_1_aihat_g

Run parameters

Specify parameters required by the pipeline

project
vision-230607

bucket
vision-storage1

raw-data
nlp-1-aihat-game-data/short_dataset.zip

docs-number
100000

min-words-per-sentence
4

cleaned-sentences-file
nlp-1-aihat-game-data/cleaned_sentences_file.txt

corpus-file
nlp-1-aihat-game-data/corpus.txt

synonyms-file
nlp-1-aihat-game-data/synonyms.json

guess-model-type
cbow

guess-model-dim
70

guess-model-thread
70

accuracy
222%

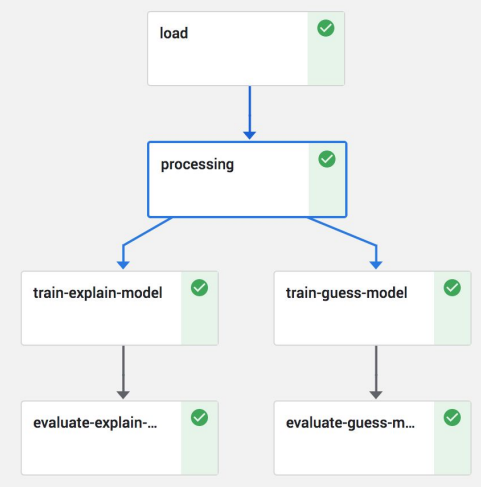
Metrics tracking for a Run (kubeflow)

Experiments > fasttext-two-models

Clone run Terminate Archive

← ✓ run1-baseline

Graph Run output Config



```
graph TD; load[load] --> processing[processing]; processing --> train-explain-model[train-explain-model]; processing --> train-guess-model[train-guess-model]; train-explain-model --> evaluate-explain-model[evaluate-explain-model]; train-guess-model --> evaluate-guess-model[evaluate-guess-model];
```

iris-pipeline-xm779-1271021395

Artifacts Input/Output Volumes Manifest Logs

Input parameters

bucket
corpus-
load-cl
project
Output
process

← ✓ new-run-to-win-game

Graph Run output Config

Metrics

	guess-accuracy	explain-accuracy
evaluate-guess-model	22.222%	
evaluate-explain-model		3.373%

Metrics tracking for a Run (mlflow)

GuessExperiment_1 > **Run 4f6051f36e98433d890888c362a7923f** ▾

Date: 2019-10-23 19:05:36

Run ID: 4f6051f36e98433d890888c362a7923f

Source:

User: user


▼ Notes 

None

▼ Parameters

Name	Value
docs_number	100000
min_words_per_sentence	4
model_dim	50
model_epoch	15
model_type	cbow
model_wordNgrams	1

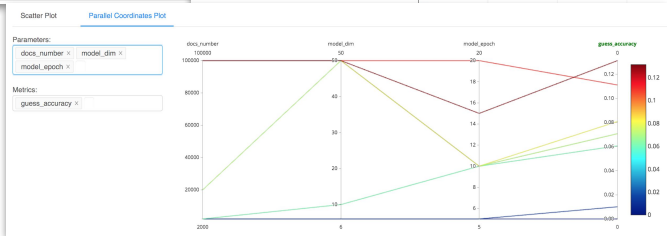
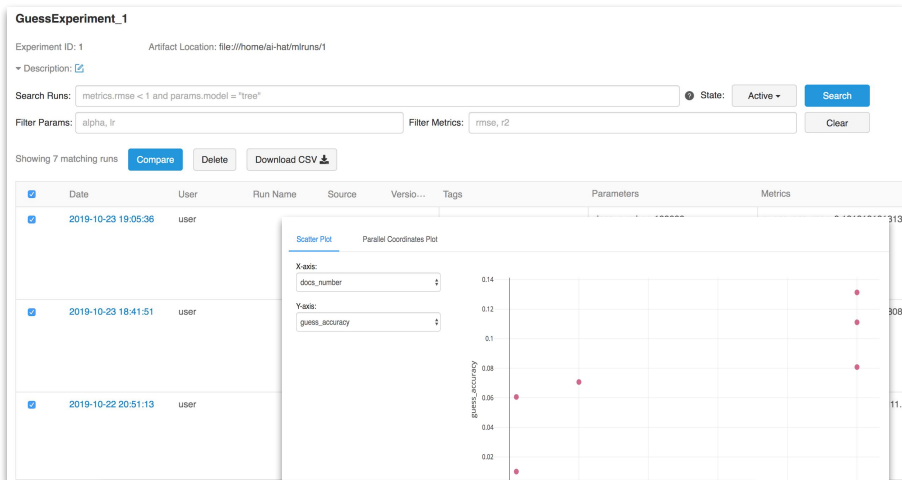
▼ Metrics

Name	Value
guess_accuracy 	0.131

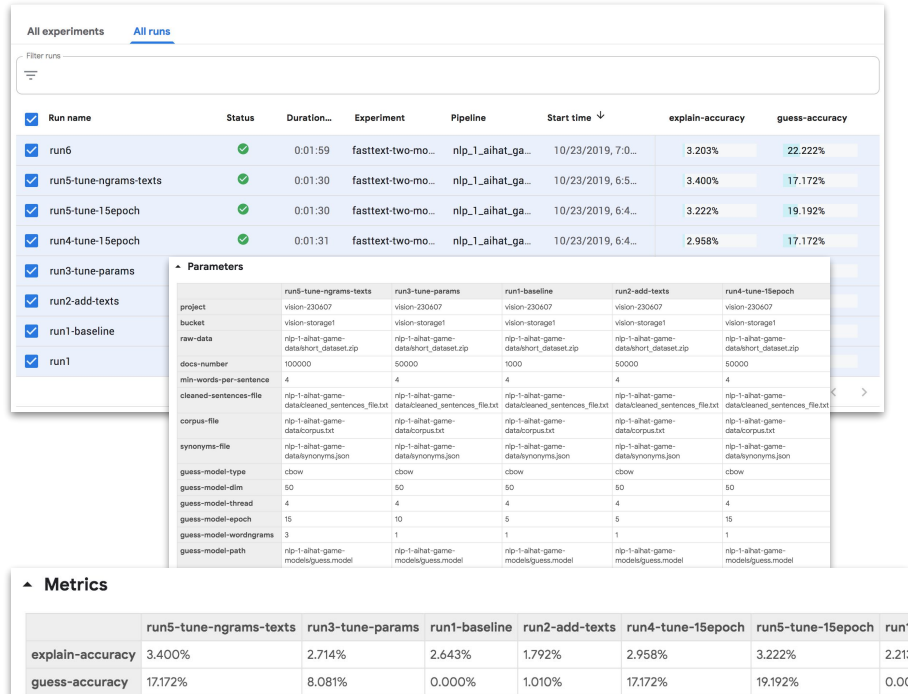
**guess_accuracy и explain_accuracy
можно логировать
в одном RUN**

Benchmark results

dvc + mlflow

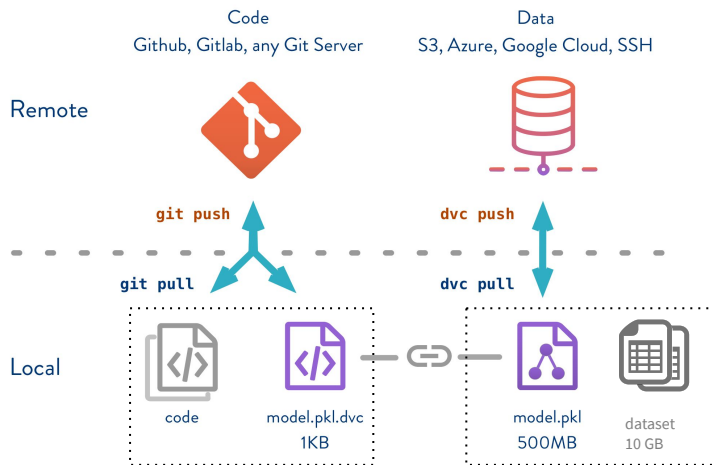


kubeflow



Models / artifacts

dvc + mlflow



models/artifact versioning with
DVC

kubeflow

Artifacts	Input/Output	Volumes	Manifest	Logs
Input parameters				
bucket	vision-storage1			
explain-model-dim	70			
explain-model-epoch	20			
explain-model-path	nlp-1-aihat-game-models/explain.model			
explain-model-thread	4			
explain-model-type	skipgram			
explain-model-wordngrams	2			
processing-corpus-file-output	gs://vision-storage1/nlp-1-aihat-game-data/corpus.txt			
project	vision-230607			
Output parameters				
train-explain-model-model-path-output	gs://vision-storage1/nlp-1-aihat-game-models/explain.model			

need manual edit for model path
each run

Benchmark approaches of DVC, MLflow and kubeflow

	DVC	MLflow	kubeflow
Artifacts version control (models, datasets, etc.)	yes dvc run args	yes log_artifact()	yes* via metadata API
Pipeline execution DAG	yes	no*	yes
Caching of intermediate results	yes	no	no
Experiment management (tracking metrics, comparison, visualization)	yes-no*	yes	yes
Metadata	.dvc files	params, metrics, artifacts meta	kfmd library
Deployment/serving	no	yes	yes
Works locally	yes	yes	no*

* not out of the box or not flexible enough but possible to do/use/hack

Benchmark approaches of DVC, MLflow and kubeflow

	DVC	MLflow	kubeflow
Learning rate level (for DS)	moderate (git, shell cmd)	low (python only)	high* *k8s, GCP
Efforts for start with	moderate	low	high*
Cost of money	free	free	k8s cluster
Maintainability / customization efforts	low	low	high 50+ microservices

* not out of the box or not flexible enough but possible to do/use/hack

Benchmark approaches of DVC, MLflow and kubeflow

	DVC	MLflow	kubeflow
pipelines	Complex pipeline with intermediate data saved into separate files. No duplicated computations and copy of artifacts.	Simple pipelines, one model. Serving model out of the box.	Pipelines with different resources requirements.
cool feature	Handful for experimentation local or collaboration (shared resources).	Nice UI for tracking metrics/params and experiment benchmark.	Reusable components, experiments benchmark and computation graph visualization
reproducibility	Reproducibility out of the box. Easy to checkout to previous version.	Need to save a copy for all data/code and artifacts to get reproducibility	Work in progress to versioning but still many drawbacks. Users' responsibility.

Thank you

Mikhail Rozhkov

ml-repa.ru

Some links

- [DVC tutorials](#)
- [MLflow tracking](#)
- [Kubeflow pipelines quickstart](#)
- [Reproducibility in Machine Learning](#)
- [Kubeflow v0.6: support for artifact tracking, data versioning & multi-user](#)
- [The Data Science Bill of Rights](#)
- [KFServing](#)