

# Introducing a transparency-driven platform for creating, comparing and explaining machine learning pipelines

Tom Klaver<sup>\*1</sup>, Erik Tjong Kim Sang<sup>1</sup>, Aysenur Bilgin<sup>2</sup>, Kim Smeenk<sup>3</sup>,  
Laura Hollink<sup>2</sup>, Jacco van Ossenbruggen<sup>2</sup>, Frank Harbers<sup>3</sup>, and Marcel Broersma<sup>3</sup>

<sup>1</sup>Netherlands eScience Center, Science Park 140, 1098 XG Amsterdam

<sup>2</sup>Centrum Wiskunde & Informatica, Science Park 123, 1098 XG Amsterdam

<sup>3</sup>University of Groningen, 9712 CP Groningen

**Abstract**—Machine learning offers great promises for many tasks relevant to science and industry. Before users accept machine-generated solutions, they need to be convinced of the quality of the system. However, the internal operation of a machine learning system is often complicated. In this demonstration, we present a transparency-driven platform that can create, compare and explain machine learning pipelines by means of graphical and tabular visualizations. The platform displays a genre classification task in digital newspaper archives, however it can be employed for any other machine learning task.

**Index Terms**—Machine learning, transparency, genre classification, digital humanities

## I. INTRODUCTION

Increasing amounts of data have facilitated the development of data-driven approaches that rely on complex machine learning algorithms. With respect to extracting meaningful insights, large amounts of data usage for automated decision making has been developing concerns [3] [7]. Inspired by the initiatives such as the Responsible Data Science (RDS<sup>1</sup>) and Explainable AI (XAI) [4] that aim to tackle these concerns by increasing awareness on transparency, we have built a web-based platform to encourage the trusted use of machine learning pipelines. The demonstrated application targets not only data scientists (as in many existing architectures such as TensorBoard<sup>2</sup>, Spark ML<sup>3</sup> and Skater [2]) but also lay users and can be applied to any other machine learning task.

In this demonstration, we present the functionality of the platform using a genre classification task (for further details about this research, see [1]). We address the following questions, which are central to most of the existing machine learning tasks from both developer and lay user perspectives:

- 1) How can we optimize the machine learning pipelines for task-specific accuracy in a user-friendly way?
- 2) What form of representations facilitate communicating the important criteria driving the automated decision making process? In particular, a user should be able to understand the type of errors being made and gauge to what extent the model can be trusted for a given task.

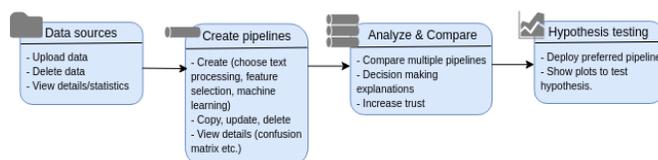


Fig. 1. Major platform functionality

- 3) Can we provide a user-friendly interface that allows a lay user to utilize complex machine learning pipelines by uploading data sets, comparing pipelines and visualizing the results as well as allowing them to understand the decisions given by the pipelines?

In the next section, we provide details regarding the live actions and interactions offered by the platform.

## II. DEMONSTRATION OF THE PLATFORM

The platform<sup>4</sup> is developed using Flask<sup>5</sup>, MongoDB<sup>6</sup>, scikit-learn [5], Frog<sup>7</sup> natural language processing suite and LIME [6]. It offers two interfaces: 1) web-based graphical user interface (GUI) and 2) Jupyter Lab<sup>8</sup> interface. Both interfaces allow for uploading data sets, constructing custom pipelines, analyzing results of individual pipelines as well as comparing multiple pipelines, providing explanations on local and global levels and deploying the pipelines to test hypotheses.

As summarized in Fig. 1, the web-based GUI has the following page views for demonstration:

- 1) Data view (Fig. 2): Users can upload, view and delete data sets. Clicking on a data set shows the task-specific statistics of the data set (e.g. size of the data set, distribution of features, etc.).
- 2) Pipelines view: Users can customize pipelines by choosing from a set of predefined text processing and feature selection steps as well as by configuring the desired machine learning algorithm (currently a limited number

\* E: t.klaver@esciencecenter.nl. T: +31651894359

<sup>1</sup><http://www.responsibledata-science.org/>

<sup>2</sup>[https://www.tensorflow.org/guide/summaries\\_and\\_tensorboard](https://www.tensorflow.org/guide/summaries_and_tensorboard)

<sup>3</sup><https://spark.apache.org/docs/1.2.2/ml-guide.html>

<sup>4</sup>The code is available at: <https://github.com/news-gac/platform>

<sup>5</sup><http://flask.pocoo.org/>

<sup>6</sup><https://www.mongodb.com>

<sup>7</sup><https://languagemachines.github.io/frog/>

<sup>8</sup><https://jupyterlab.readthedocs.io/en/stable/>

Data sources added by you

Title	Processed	Description	Time of execution	Number of instances	
N2B2G Training	Yes	The latest dataset for training	20-11-2018 09:27	617	<a href="#">Delete</a>
N2B2G Testing	Yes	The latest dataset for testing	20-11-2018 09:27	91	<a href="#">Delete</a>
Linked NRC (2018)	Yes	Linked NRC articles from Proton	22-11-2018 12:22	2390	<a href="#">Delete</a>
ROW + 5 features (N2B2G train)	Yes	ROW + 5 features (N2B2G train)	22-11-2018 17:27	617	<a href="#">Delete</a>
ROW + 5 features (N2B2G test)	Yes	ROW + 5 features (N2B2G test)	22-11-2018 17:34	91	<a href="#">Delete</a>
row-1050-1054-01	Yes	row-1050-1054-01	22-11-2018 23:43	3000	<a href="#">Delete</a>
row-1050-1054-01	Yes	row-1050-1054-01	22-11-2018 23:45	3000	<a href="#">Delete</a>
row-1050-1054-01	Yes	row-1050-1054-01	22-11-2018 23:48	3000	<a href="#">Delete</a>
row-1050-1054-01	Yes	row-1050-1054-01	22-11-2018 23:50	3000	<a href="#">Delete</a>
row-1050-1054-01	Yes	row-1050-1054-01	22-11-2018 23:53	2883	<a href="#">Delete</a>
integrated-1050-1054-01	Yes	integrated-1050-1054-01	22-11-2018 23:53	8000	<a href="#">Delete</a>
integrated-1050-1054-01	Yes	integrated-1050-1054-01	22-11-2018 00:05	8000	<a href="#">Delete</a>
integrated-1050-1054-01	Yes	integrated-1050-1054-01	22-11-2018 00:04	8000	<a href="#">Delete</a>
integrated-1050-1054-01	Yes	integrated-1050-1054-01	22-11-2018 00:06	8000	<a href="#">Delete</a>
integrated-1050-1054-01	Yes	integrated-1050-1054-01	25-11-2018 09:07	7104	<a href="#">Delete</a>

Fig. 2. Data view: users can upload, view and delete data sets. <https://bit.ly/2AN7cZs>

Display title:  The display title of the pre-processing configuration.

Data source:  The data source for this pipeline.

Preprocessing

Drop words removal:  The most common words that are not specific to the different classes are removed.

Levenshtein:  The process of removing additional codings of a word in order to return the base or dictionary form, which is known as the lemma.

Remove Quotes:  Remove quotes.

NLP tools

Natural language processing tool for feature extraction

Eng -> Frug -> Tr0D -> Tr0D

loading:  Transformers features by using embeddings that are robust to outliers.

Machine learning

Machine learning: The machine learning algorithm to use.

Classifier:  Loss:  Metric:  Model:  Hyperparameter:  Support Vector:  Softmax:

Parameters:

Fig. 3. Pipeline view: Users can customize pipelines by choosing from a set of predefined text processing and feature selection steps as well as by configuring the desired machine learning algorithm. <https://bit.ly/2sq37FZ>

of classifiers from scikit-learn [5] library are implemented) (Fig. 3). The option to use optimized algorithm parameters by means of grid search<sup>9</sup> is also available to the users. Once a pipeline is constructed, users can track the training status and view the run-times of training upon completion. Once successfully trained, pipelines can be copied, analyzed and deleted through the interface. The analysis view allows the users to examine several performance metrics tailored for the task (e.g. accuracy, recall, confusion matrix, etc. for classification tasks) as well as global-level explanations for the pipeline depending on the type of the machine learning algorithm selected (e.g. feature importance ranking plots). The trained pipelines are stored in database and can be deployed on other data sets.

- 3) ACE (Analyze-Compare-Explain) view (Fig. 4): Users can critically analyze and understand the decision making processes of not only a single pipeline but multiple pipelines. For this purpose, we make use of several graphical views along with tabular views that allow the user either explore the local-level explanation concerning one instance (i.e. one data point) or determine the majority vote (agreed prediction) of multiple pipelines. Using this

<sup>9</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.GridSearchCV.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html)

#### Analyse - Compare - Explain Pipelines

Show 10 entries

Article #	Pipeline	Predicted Genre	True Genre	Explanation
602	RF FR0G	LOS	SER	<a href="#">Click for explanation</a>
602	SVC FR0G	SER	SER	<a href="#">Click for explanation</a>
602	RF FR0G	LOS	SER	<a href="#">Click for explanation</a>
602	RF FR0G	MED	SER	<a href="#">Click for explanation</a>
602	NLP FR0G	LOS	SER	<a href="#">Click for explanation</a>
602	RF ROW	LOS	SER	<a href="#">Click for explanation</a>
602	SVC ROW	LOS	SER	<a href="#">Click for explanation</a>
602	RF ROW	LOS	SER	<a href="#">Click for explanation</a>
602	RF ROW	RIC	SER	<a href="#">Click for explanation</a>
602	NLP ROW	LOS	SER	<a href="#">Click for explanation</a>

Showing 1 to 10 of 10,239 entries

Search Article #  Search Pipeline  Search Predicted Genre  Search True Genre  Search Explanation

Showing 10 entries

Mutually Agreeing Pipelines # Predicted Genre Number of Articles Articles

RF FR0G	MED	22	<a href="#">[1]</a> <a href="#">[2]</a> <a href="#">[3]</a> <a href="#">[4]</a> <a href="#">[5]</a> <a href="#">[6]</a> <a href="#">[7]</a> <a href="#">[8]</a> <a href="#">[9]</a> <a href="#">[10]</a> <a href="#">[11]</a> <a href="#">[12]</a> <a href="#">[13]</a> <a href="#">[14]</a> <a href="#">[15]</a> <a href="#">[16]</a> <a href="#">[17]</a> <a href="#">[18]</a> <a href="#">[19]</a> <a href="#">[20]</a> <a href="#">[21]</a> <a href="#">[22]</a>
RF FR0G	MED	49	<a href="#">[1]</a> <a href="#">[2]</a> <a href="#">[3]</a> <a href="#">[4]</a> <a href="#">[5]</a> <a href="#">[6]</a> <a href="#">[7]</a> <a href="#">[8]</a> <a href="#">[9]</a> <a href="#">[10]</a> <a href="#">[11]</a> <a href="#">[12]</a> <a href="#">[13]</a> <a href="#">[14]</a> <a href="#">[15]</a> <a href="#">[16]</a> <a href="#">[17]</a> <a href="#">[18]</a> <a href="#">[19]</a> <a href="#">[20]</a> <a href="#">[21]</a> <a href="#">[22]</a> <a href="#">[23]</a> <a href="#">[24]</a> <a href="#">[25]</a> <a href="#">[26]</a> <a href="#">[27]</a> <a href="#">[28]</a> <a href="#">[29]</a> <a href="#">[30]</a> <a href="#">[31]</a> <a href="#">[32]</a> <a href="#">[33]</a> <a href="#">[34]</a> <a href="#">[35]</a> <a href="#">[36]</a> <a href="#">[37]</a> <a href="#">[38]</a> <a href="#">[39]</a> <a href="#">[40]</a> <a href="#">[41]</a> <a href="#">[42]</a> <a href="#">[43]</a> <a href="#">[44]</a> <a href="#">[45]</a> <a href="#">[46]</a> <a href="#">[47]</a> <a href="#">[48]</a> <a href="#">[49]</a>

Fig. 4. ACE view: Users can critically analyze and understand the decision making processes of multiple pipelines. <https://bit.ly/2FwaDa8>

view, the users can gather meaningful insights regarding the decision making processes of the complex machine learning pipelines and develop trust in them accordingly.

- 4) Hypothesis testing: Users can deploy the trusted pipeline in other data sets for the purpose of retrieving statistics that can help them test their hypotheses. In this view, the platform displays plots for a visual analysis of the predictions.

### III. SUMMARY

We present a ready-to-use platform for creating, comparing and explaining machine learning pipelines and demonstrate the practical impact of transparency for a genre classification task.

### ACKNOWLEDGMENTS

The study described in this paper was executed as a part of the NEWSGAC project which is funded by the Netherlands eScience Center<sup>10</sup> and CLARIAH<sup>11</sup>.

### REFERENCES

- [1] A. Bilgin, E. Tjong Kim Sang, K. Smeenk, L. Hollink, J. van Ossendruppen, F. Harbers, and M. Broersma. Utilizing a transparency-driven environment toward trusted automatic genre classification: A case study in journalism history. In *2018 IEEE 14th International Conference on e-Science (e-Science)*, pages 486–496, Oct 2018.
- [2] Primit Choudhary, Aaron Kramer, and contributors datascience.com team. Skater: Model Interpretation Library, March 2018.
- [3] Natasha Duarte, Emma Llanso, and Anna Loup. Mixed Messages? The Limits of Automated Social Media Content Analysis. In *Conference on Fairness, Accountability and Transparency*, volume 81 of *Proceedings of Machine Learning Research*. PMLR, 2018.
- [4] David Gunning. Explainable artificial intelligence (xai). *Defense Advanced Research Projects Agency (DARPA), nd Web*, 2017.
- [5] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [6] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Why should I trust you?: Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1135–1144. ACM, 2016.
- [7] Wil M. P. van der Aalst, Martin Bichler, and Armin Heinzl. Responsible data science. *Business & Information Systems Engineering*, 59(5):311–313, Oct 2017.

<sup>10</sup>esciencecenter.nl

<sup>11</sup>clariah.nl