

# FAT Forensics: A Python Toolbox for Implementing and Deploying Fairness, Accountability and Transparency Algorithms in Predictive Systems

Kacper Sokol<sup>1</sup>, Alexander Hepburn<sup>2</sup>, Rafael Poyiadzi<sup>2</sup>, Matthew Clifford<sup>2</sup>, Raul Santos-Rodriguez<sup>2</sup>, and Peter Flach<sup>1</sup>

<sup>1</sup> Department of Computer Science, University of Bristol <sup>2</sup> Department of Engineering Mathematics, University of Bristol

DOI: [10.21105/joss.01904](https://doi.org/10.21105/joss.01904)

## Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Ariel Rokem](#) ↗

## Reviewers:

- [@bernease](#)
- [@osolari](#)

Submitted: 12 September 2019

Published: 29 January 2020

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC-BY](#)).

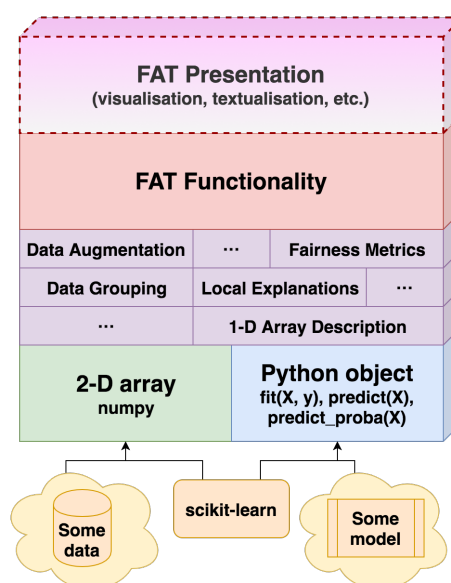
## Background

Predictive systems, in particular machine learning algorithms, can take important, and sometimes legally binding, decisions about our everyday life. In most cases, however, these systems and decisions are neither regulated nor certified. Given the potential harm that these algorithms can cause, their qualities such as **fairness**, **accountability** and **transparency** (FAT) are of paramount importance. To ensure high-quality, fair, transparent and reliable predictive systems, we developed an open source Python package called *FAT Forensics*. It can inspect important fairness, accountability and transparency aspects of predictive algorithms to automatically and objectively report them back to engineers and users of such systems. Our toolbox can evaluate all elements of a predictive pipeline: data (and their features), models and predictions. Published under the BSD 3-Clause open source licence, *FAT Forensics* is opened up for personal and commercial usage.

## Summary

*FAT Forensics* is designed as an interoperable framework for *implementing*, *testing* and *deploying* novel algorithms devised by the FAT research community. It facilitates their evaluation and comparison against the state of the art, thereby democratising access to these techniques. In addition to supporting research in this space, the toolbox is capable of analysing all components of a predictive pipeline – data, models and predictions – by considering their fairness, accountability (robustness, security, safety and privacy) and transparency (interpretability and explainability).

*FAT Forensics* collates all of these diverse tools and algorithms under a common application programming interface. This is achieved with a modular design that allows to share and reuse a collection of core algorithmic components – see Figure 1. This architecture makes the process of creating new algorithms as easy as connecting the right blocks, therefore supporting a range of diverse use cases.



**Figure 1:** Modular architecture of FAT Forensics.

The format requirements for data sets and predictive models are kept to a minimum, lowering any barriers for adoption of *FAT Forensics* in new and already well-established projects. In this abstraction a data set is assumed to be a two-dimensional NumPy array: either a classic or a structured array. The latter is a welcome addition given that some of the features may be categorical (string-based). A predictive model is assumed to be a plain Python object that has `fit`, `predict` and, optionally, `predict_proba` methods. This flexibility makes our package compatible with scikit-learn – the most popular Python machine learning toolbox – without introducing additional dependencies. Moreover, this approach makes *FAT Forensics* compatible with other packages for predictive modelling since their predictive functions can be easily wrapped inside a Python object with all the required methods.

Our package improves over existing solutions as it collates algorithms across the FAT domains, taking advantage of their shared functional building blocks. The common interface layer of the toolbox supports several *modes of operation*. The **research mode** (data in – visualisation out), where the tool can be loaded into an interactive Python session, e.g., a Jupyter Notebook, supports prototyping and exploratory analysis. This mode is intended for FAT researchers who may use it to propose new fairness metrics, compare them with the existing ones or use them to inspect a new system or a data set. The **deployment mode** (data in – data out) can be used as a part of a data processing pipeline to provide a (numerical) FAT analytics, supporting automated reporting and dashboarding. This mode is intended for machine learning engineers and data scientists who may use it to monitor or evaluate a predictive system during its development and deployment.

To encourage long-term maintainability, sustainability and extensibility, *FAT Forensics* has been developed employing software engineering best practice such as unit testing, continuous integration, well-defined package structure and consistent code formatting. Furthermore, our toolbox is supported by a thorough and beginner-friendly documentation that is based on four main pillars, which together build up the user’s confidence in using the package:

- narrative-driven **tutorials** designated for new users, which provide a step-by-step guidance through practical use cases of all the main aspects of the package;
- **how-to guides** created for relatively new users of the package, which showcase the flexibility of the toolbox and explain how to use it to solve user-specific FAT challenges; for example, how to build your own local surrogate explainer by pairing a data augementer and an inherently transparent local model;

- **API documentation** describing functional aspects of the algorithms implemented in the package and designated for a technical audience as a reference material; it is complemented by task-focused *code examples* that put the functions, objects and methods in context; and
- a **user guide** discussing theoretical aspects of the algorithms implemented in the package such as their restrictions, caveats, computational time and memory complexity, among others.

We hope that this effort will encourage the FAT community to contribute their algorithms to *FAT Forensics*. We offer it as an attractive alternative to releasing yet more standalone packages, keeping the toolbox at the frontiers of algorithmic fairness, accountability and transparency research. For a more detailed description of *FAT Forensics*, we point the reader to its documentation<sup>1</sup> and the paper (Sokol, Santos-Rodriguez, & Flach, 2019) describing its design, scope and usage examples.

## Related Work

A recent attempt to create a common framework for FAT algorithms is the *What-If* tool<sup>2</sup>, which implements various fairness and explainability approaches. A number of Python packages collating multiple state-of-the-art algorithms for either fairness, accountability or transparency also exist. Available algorithmic *transparency* packages include:

- Skater<sup>3</sup> (Kramer et al., 2018),
- ELI5<sup>4</sup>,
- Microsoft's Interpret<sup>5</sup> (Nori, Jenkins, Koch, & Caruana, 2019), and
- IBM's AI Explainability 360<sup>6</sup>.

Packages implementing individual algorithms are also popular. For example, LIME<sup>7</sup> for Local Interpretable Model-agnostic Explanations (Ribeiro, Singh, & Guestrin, 2016) and PyCEbox<sup>8</sup> for Partial Dependence (Friedman, 2001) and Individual Conditional Expectation (Goldstein, Kapelner, Bleich, & Pitkin, 2015) plots.

Algorithmic *fairness* packages are also ubiquitous, for example: Microsoft's fairlearn<sup>9</sup> (Agarwal, Beygelzimer, Dudik, Langford, & Wallach, 2018) and IBM's AI Fairness 360<sup>10</sup> (Bellamy et al., 2018). However, *accountability* is relatively underexplored. The most prominent software in this space deals with robustness of predictive systems against adversarial attacks, for example:

- FoolBox<sup>11</sup>,
- CleverHans<sup>12</sup> and
- IBM's Adversarial Robustness 360 Toolbox<sup>13</sup>.

---

<sup>1</sup><https://fat-forensics.org>

<sup>2</sup><https://pair-code.github.io/what-if-tool>

<sup>3</sup><https://github.com/oracle/Skater>

<sup>4</sup><https://github.com/TeamHG-Memex/eli5>

<sup>5</sup><https://github.com/interpretml/interpret>

<sup>6</sup><https://github.com/IBM/AIX360>

<sup>7</sup><https://github.com/marcotcr/lime>

<sup>8</sup><https://github.com/AustinRochford/PyCEbox>

<sup>9</sup><https://github.com/fairlearn/fairlearn>

<sup>10</sup><https://github.com/IBM/AIF360>

<sup>11</sup><https://github.com/bethgelab/foolbox>

<sup>12</sup><https://github.com/tensorflow/cleverhans>

<sup>13</sup><https://github.com/IBM/adversarial-robustness-toolbox>

*FAT Forensics* aims to bring together all of this functionality from across fairness, accountability and transparency domains with its modular implementation. This design principle enables the toolbox to support two modes of operation: research and deployment. Therefore, the package caters to a diverse audience and supports a range of tasks such as implementing, testing and deploying *FAT* solutions. Abstracting away from fixed data set and predictive model formats adds to its versatility. The development of the toolbox adheres to best practices for software engineering and the package is supported by a rich documentation, both of which make it stand out amongst its peers.

## Acknowledgements

This work was financially supported by Thales, and is the result of a collaborative research agreement between Thales and the University of Bristol.

## References

- Agarwal, A., Beygelzimer, A., Dudik, M., Langford, J., & Wallach, H. (2018). A reductions approach to fair classification. In J. Dy & A. Krause (Eds.), *Proceedings of the 35th international conference on machine learning*, Proceedings of machine learning research (Vol. 80, pp. 60–69). Stockholmsmässan, Stockholm Sweden: PMLR. Retrieved from <http://proceedings.mlr.press/v80/agarwal18a.html>
- Bellamy, R. K. E., Dey, K., Hind, M., Hoffman, S. C., Houde, S., Kannan, K., Lohia, P., et al. (2018). AI fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias. *arXiv preprint arXiv:1810.01943*.
- Friedman, J. H. (2001). Greedy function approximation: A gradient boosting machine. *Annals of statistics*, 1189–1232.
- Goldstein, A., Kapelner, A., Bleich, J., & Pitkin, E. (2015). Peeking inside the black box: Visualizing statistical learning with plots of individual conditional expectation. *Journal of Computational and Graphical Statistics*, 24(1), 44–65.
- Kramer, A., Choudhary, P., silversurfer84, Dyke, B. V., Thai, A., Pasumathy, N., Lemaitre, G., et al. (2018). *Datascienceinc/skater: 1.1.2*. Zenodo. doi:[10.5281/zenodo.1423046](https://doi.org/10.5281/zenodo.1423046)
- Nori, H., Jenkins, S., Koch, P., & Caruana, R. (2019). InterpretML: A unified framework for machine learning interpretability. *arXiv preprint arXiv:1909.09223*.
- Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "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, san francisco, ca, usa, august 13-17, 2016* (pp. 1135–1144).
- Sokol, K., Santos-Rodriguez, R., & Flach, P. (2019). *FAT forensics: A python toolbox for algorithmic fairness, accountability and transparency*. *arXiv preprint arXiv:1909.05167*.