

1 fairret: a Framework for Differentiable Fairness 2 Regularization Terms

3 MaryBeth Defrance¹, Maarten Buyl¹ and Tijl De Bie¹

4 ¹Ghent University, Belgium

5 Abstract

6 Current fairness toolkits in machine learning only admit a limited range of fairness definitions and have
7 seen little integration with automatic differentiation libraries, despite the central role these libraries
8 play in modern machine learning pipelines. We present a framework of fairness regularization terms
9 (FAIRRETS) which quantify bias as modular, flexible objectives that are easily integrated in automatic
10 differentiation pipelines. By employing a general definition of fairness through linear-fractional statistics,
11 many group fairness definitions can be enforced. Experiments show minimal loss of predictive power
12 compared to baselines. Our contribution includes a PyTorch implementation of the FAIRRET library.

13 Keywords

fairness, machine learning, library, automatic differentiation, fairness definitions

14 1. Introduction

15 The field of AI fairness has been concerned with formalizing ethical concepts of discrimination
16 and bias in technical definitions that can be assessed and pursued in AI systems [1]. A popular
17 paradigm for this formalization in binary classification is to use *group fairness* definitions [2],
18 which require the model's predictions to treat people from different sensitive groups similarly.

19 Despite ample research on group fairness definitions and methods to achieve them, an
20 easy-to-use and flexible implementation has not yet been realized. Popular fairness toolkits
21 such as Fairlearn [3] and AIF360 [4] expect the underlying model in the form of *scikit-learn*
22 *Estimators* [5] that can be retrained at-will in fairness meta-algorithms, but this aligns poorly
23 with the paradigm of automatic differentiation libraries like PyTorch [6], which have become
24 the bedrock of modern machine learning pipelines. These toolkits only integrate with automatic
25 differentiation in their implementations of adversarial fairness [7], but these still require full
26 control over the training process and lack generality in the fairness notions they can enforce.

27 We formally propose the FAIRRET framework in an effort to resolve these issues. At its
28 core, the framework uses *fairness regularization terms* (FAIRRETS) that can be easily integrated
29 into PyTorch-based pipelines (an example is given in Appendix A). They pursue any fairness
30 definition expressed as a parity between statistics in a linear-fractional notation, which covers
31 all group fairness definitions considered by Verma and Rubin [2]. Thus, all these definitions are
32 fully compatible with any FAIRRET in any differentiable model.


EWAF'24: European Workshop on Algorithmic Fairness, July 01–03, 2024, Mainz, Germany

✉ marybeth.defrance@ugent.be (M. Defrance); maarten.buyl@ugent.be (M. Buyl); tijl.debie@ugent.be (T. D. Bie)

🆔 0000-0002-6570-8857 (M. Defrance); 0000-0002-5434-2386 (M. Buyl); 0000-0002-2692-7504 (T. D. Bie)



© 2024 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

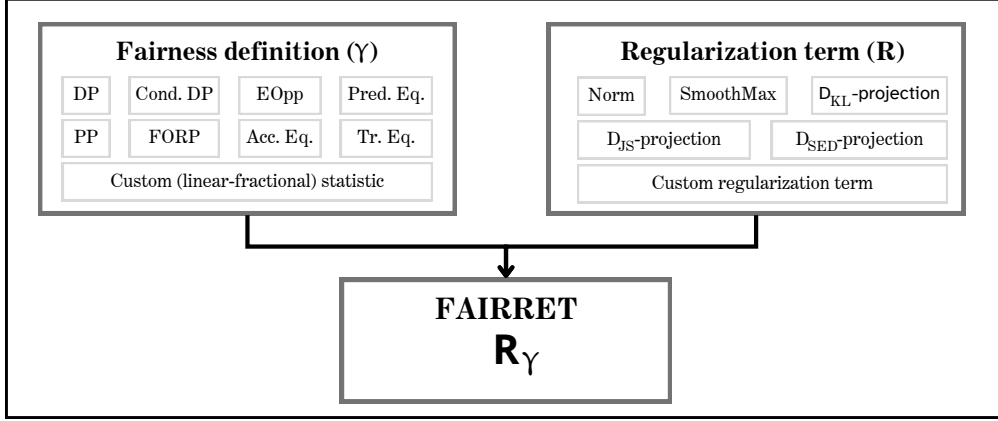


Figure 1: This diagram shows the modular nature of FAIRRET and provides an overview of the fairness definitions and methods present in the framework and notes the flexibility to implement novel ones.

33 In contrast to Fairlearn and AIF360, our proposed FAIRRETS act as a loss term that can simply
 34 be added *within* a training step. Two PyTorch-specific projects with similar goals as our paper
 35 are FairTorch [8] and the Fair Fairness Benchmark (FFB) [9]. However, neither present a formal
 36 framework and both only support a limited range of fairness definitions.

37 This work is an extended abstract of a full paper [10] presented at the *International Conference*
 38 *on Learning Representations (ICLR) 2024*. The implementation of our framework is available at
 39 <https://github.com/aida-ugent/fairret>, which we are currently extending into a full library.

40 2. How to build your FAIRRET

41 A FAIRRET is defined by two elements. First is the fairness definition it aims to satisfy. Second is
 42 the method used to evaluate the model with regard to that fairness definition. Figure 1 illustrates
 43 this combination and lists the definitions and methods already integrated into the framework.

44 2.1. Fairness definitions

45 Let $\mathbf{X} \in \mathbb{R}^{d_x}$ denote the feature vector of an individual, $\mathbf{S} \in \mathbb{R}^{d_s}$ their *sensitive* feature vector
 46 and $Y \in \{0, 1\}$ a binary output label. We want to learn a probabilistic classifier f such that
 47 its predictions $f(\mathbf{X})$ match Y while minimizing disparities over different \mathbf{S} . Our definition of
 48 sensitive features \mathbf{S} as real-valued, d_s -dimensional vectors allows us to take a mix of multiple
 49 sensitive traits into account, both discrete and continuous. Categorical sensitive features are
 50 one-hot encoded, e.g. by encoding ‘white’ or ‘non-white’ as the vectors $\mathbf{S} = (1, 0)^\top$ and
 51 $\mathbf{S} = (0, 1)^\top$ respectively. The variable S_k denotes the k th sensitive feature.

52 We use a simplified version of the solution from Celis et al. [11] to translate fairness definitions
 53 as a parity between linear-fractional statistics γ :

$$54 \quad \gamma(k; f) = \frac{\mathbb{E}[S_k(\alpha_0(\mathbf{X}, Y) + f(\mathbf{X})\beta_0(\mathbf{X}, Y))]}{\mathbb{E}[S_k(\alpha_1(\mathbf{X}, Y) + f(\mathbf{X})\beta_1(\mathbf{X}, Y))]} \quad (1)$$

Table 1

Fairness definitions and their α and β functions. Conditional Demographic Parity encompasses many notions with an arbitrary function ζ conditioned on the input \mathbf{X} .

| Fairness Definition | α_0 | β_0 | α_1 | β_1 |
|-------------------------------------|------------|---------------------|---------------------|-----------|
| Demographic Parity [12] | 0 | 1 | 1 | 0 |
| Conditional Demographic Parity [13] | 0 | $\zeta(\mathbf{X})$ | $\zeta(\mathbf{X})$ | 0 |
| Equal Opportunity [14] | 0 | Y | Y | 0 |
| False Positive Parity [14] | 0 | 1 - Y | 1 - Y | 0 |
| Predictive Parity [15] | 0 | Y | 0 | 1 |
| False Omission Parity | Y | -Y | 1 | -1 |
| Accuracy Equality [16] | 1 - Y | 2Y - 1 | 1 | 0 |
| Treatment Equality [16] | Y | -Y | 0 | 1 - Y |

55 with α_0 , α_1 , β_0 , and β_1 functions that do not depend on \mathbf{S} or f . Table 1 shows the statistic γ
 56 for a range of fairness definitions, defined through their α and β functions.

57 The set \mathcal{F}_γ of probabilistic classifiers f that adhere to the fairness definition is expressed as

$$58 \quad \mathcal{F}_\gamma \triangleq \{f : \mathbb{R}^{d_x} \rightarrow \{0, 1\} \mid \forall k \in [d_s] : \gamma(k; f) = \bar{\gamma}(f)\}. \quad (2)$$

59 In other words, the statistic $\gamma(k; f)$ for each sensitive attribute S_k should equal the overall
 60 statistic $\bar{\gamma}(f) \triangleq \frac{\mathbb{E}[\alpha_0(\mathbf{X}, Y) + f(\mathbf{X})\beta_0(\mathbf{X}, Y)]}{\mathbb{E}[\alpha_1(\mathbf{X}, Y) + f(\mathbf{X})\beta_1(\mathbf{X}, Y)]}$ computed independently of the sensitive attributes. By
 61 fixing $\bar{\gamma}$ to a constant $c \in \mathbb{R}$, any fairness definition can be enforced with a **linear** constraint:

$$62 \quad \gamma(k; f) = c \iff \mathbb{E}[S_k(\alpha_0(\mathbf{X}, Y) - c\alpha_1(\mathbf{X}, Y) + f(\mathbf{X})(\beta_0(\mathbf{X}, Y) - c\beta_1(\mathbf{X}, Y)))] = 0 \quad (3)$$

63 2.2. Regularization terms

64 The bias of a parameterized, probabilistic classifier h is quantified as a FAIRRET that can be
 65 minimized through automatic differentiation, in addition to any existing loss function \mathcal{L}_Y :

$$66 \quad \min_h \mathcal{L}_Y(h) + \lambda R_\gamma(h) \quad (4)$$

67 where $R_\gamma(h)$ is the FAIRRET for the fairness definition with statistic γ and strength $\lambda \in \mathbb{R}_{>0}$.

68 The FAIRRET framework admits many kinds of regularizers, due to the practical form of the
 69 statistics γ . Two types are currently integrated, namely *violation* and *projection* FAIRRETS.

70 We first discuss the **Norm** FAIRRET, a type of violation FAIRRET:

$$71 \quad R_\gamma(h) \triangleq \left\| \frac{\gamma(k; h)}{\bar{\gamma}(h)} - 1 \right\| \quad (5)$$

72 with $\|\cdot\|$ a norm over \mathbb{R}^{d_s} . Such a regularization term has been proposed several times [17, 18, 19],
 73 though without the same degree of modularity with respect to γ .

74 Second, an example of a projection FAIRRET is the **D_{KL} -projection**:

$$75 \quad R_\gamma(h) \triangleq \min_{f \in \mathcal{F}_{\gamma(\bar{\gamma}(h))}} \mathbb{E}[D_{KL}(f(\mathbf{X}) || h(\mathbf{X}))] \quad (6)$$

76 with D_{KL} the *Kullback-Leibler* divergence. The FAIRRET maps h onto the *closest* fair model
 77 $f \in \mathcal{F}_{\gamma(\bar{\gamma}(h))}$. Projection FAIRRETS generalize some prior work [20, 21, 22] to all definitions with
 78 linear-fractional statistics, as they are enforced with linear constraints using Eq. (3).

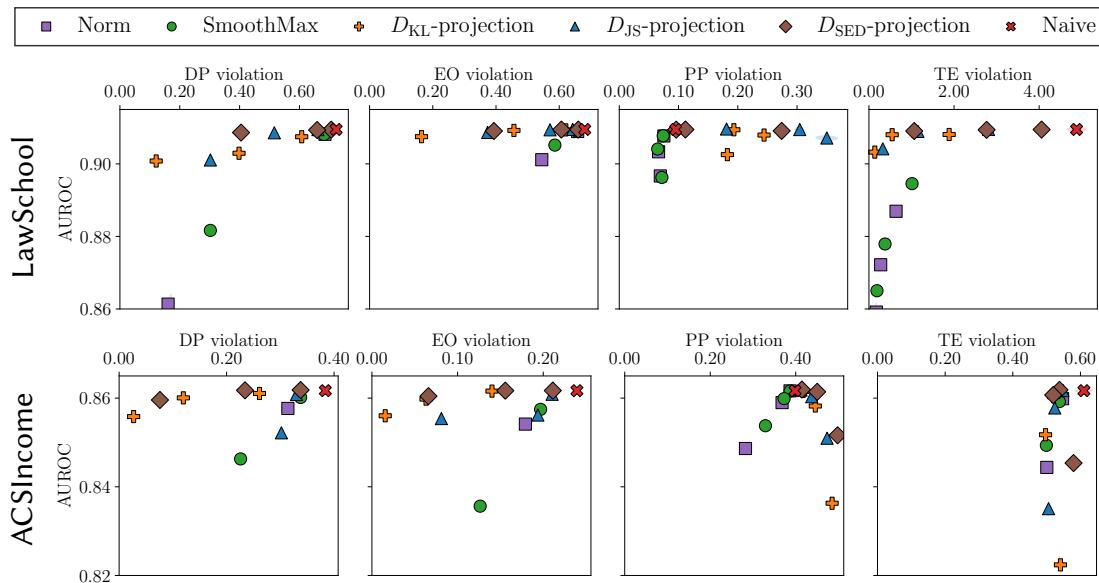


Figure 2: Mean test set results with confidence ellipse for the standard error (see Appendix B). Each marker is a separate combination of dataset, FAIRRET, FAIRRET strength, and statistic. Results in the top left are optimal. Failed runs (with an AUROC far worse than the rest) are omitted.

79 3. Experiments

80 Experiments were conducted on the *LawSchool*¹, and *ACSIncome* [23] datasets. Each dataset
 81 has multiple sensitive features, including some continuous. Figure 2 shows the results for the
 82 experiments. Each point represents a specific FAIRRET optimized for that statistic with a certain
 83 strength λ . A *Naive* baseline with $\lambda = 0$ is also included. In the full paper, the evaluation is
 84 done on two additional datasets and the FAIRRETS are compared to existing methods [10].

85 The results in Figure 2 show that the performance of a FAIRRET is dependent on the dataset
 86 itself and the fairness definition it aims to satisfy. The non-linear (yet still linear-fractional)
 87 fairness statistics like predictive parity and treatment equality seem more difficult to minimize.
 88 This leads us to conclude that not one FAIRRET can be chosen as the optimal solution, but rather
 89 that the best FAIRRET is dependent on the fairness definition and the dataset.

90 4. Conclusion

91 The FAIRRET framework allows for a wide range of fairness definitions by comparing linear-
 92 fractional statistics for each sensitive feature. We implement several FAIRRETS and show how
 93 they are easily integrated in existing machine learning pipelines utilizing automatic differentia-
 94 tion. More details can be found in the full paper [10].

¹Curated and published by the SEAPHE project

95 Acknowledgments

96 The research leading to these results has received funding from the Special Research Fund (BOF)
97 of Ghent University (BOF20/DOC/144 and BOF20/IBF/117), from the Flemish Government under
98 the “Onderzoeksprogramma Artificiële Intelligentie (AI) Vlaanderen” programme, and from the
99 FWO (project no. G0F9816N, 3G042220, G073924N).

100 References

- 101 [1] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, A. Galstyan, A Survey on Bias and
102 Fairness in Machine Learning, *ACM Computing Surveys* 54 (2021) 115:1–115:35. doi:10.
103 1145/3457607.
- 104 [2] S. Verma, J. Rubin, Fairness definitions explained, in: *Proceedings of the International*
105 *Workshop on Software Fairness*, ACM, Gothenburg Sweden, 2018, pp. 1–7. doi:10.1145/
106 3194770.3194776.
- 107 [3] S. Bird, M. Dudík, R. Edgar, B. Horn, R. Lutz, V. Milan, M. Sameki, H. Wallach, K. Walker,
108 Fairlearn: A toolkit for assessing and improving fairness in AI, Technical Report MSR-TR-
109 2020-32, Microsoft, 2020. URL: [https://www.microsoft.com/en-us/research/publication/
110 fairlearn-a-toolkit-for-assessing-and-improving-fairness-in-ai/](https://www.microsoft.com/en-us/research/publication/fairlearn-a-toolkit-for-assessing-and-improving-fairness-in-ai/).
- 111 [4] R. K. E. Bellamy, K. Dey, M. Hind, S. C. Hoffman, S. Houde, K. Kannan, P. Lohia, J. Martino,
112 S. Mehta, A. Mojsilovic, S. Nagar, K. N. Ramamurthy, J. Richards, D. Saha, P. Sattigeri,
113 M. Singh, K. R. Varshney, Y. Zhang, *AI Fairness 360: An extensible toolkit for detecting,*
114 *understanding, and mitigating unwanted algorithmic bias*, 2018. URL: [https://arxiv.org/
115 abs/1810.01943](https://arxiv.org/abs/1810.01943).
- 116 [5] L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae, P. Pret-
117 tenhofer, A. Gramfort, J. Grobler, R. Layton, J. VanderPlas, A. Joly, B. Holt, G. Varoquaux,
118 *API design for machine learning software: experiences from the scikit-learn project*, in:
119 *ECML PKDD Workshop: Languages for Data Mining and Machine Learning*, 2013, pp.
120 108–122.
- 121 [6] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin,
122 N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani,
123 S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, S. Chintala, *PyTorch: An Imperative Style,*
124 *High-Performance Deep Learning Library*, in: *Advances in Neural Information Processing*
125 *Systems*, volume 32, Curran Associates, Inc., 2019.
- 126 [7] B. H. Zhang, B. Lemoine, M. Mitchell, *Mitigating Unwanted Biases with Adversarial*
127 *Learning*, in: *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society,*
128 *AIES '18*, Association for Computing Machinery, New York, NY, USA, 2018, pp. 335–340.
129 doi:10.1145/3278721.3278779.
- 130 [8] S. Masashi, *Fairtorch*, <https://github.com/wbawakate/fairtorch>, 2020. Version 0.1.2.
- 131 [9] X. Han, J. Chi, Y. Chen, Q. Wang, H. Zhao, N. Zou, X. Hu, *FFB: A Fair Fairness Benchmark*
132 *for In-Processing Group Fairness Methods*, 2023. doi:10.48550/arXiv.2306.09468.
133 arXiv:2306.09468.
- 134 [10] M. Buyl, M. Defrance, T. D. Bie, *fairret: a framework for differentiable fairness regulariza-*

- 135 tion terms, in: The Twelfth International Conference on Learning Representations, 2024.
136 URL: <https://openreview.net/forum?id=NnyD0Rjx2B>.
- 137 [11] L. E. Celis, L. Huang, V. Keswani, N. K. Vishnoi, Classification with Fairness Constraints:
138 A Meta-Algorithm with Provable Guarantees, in: Proceedings of the Conference on
139 Fairness, Accountability, and Transparency, ACM, Atlanta GA USA, 2019, pp. 319–328.
140 doi:10.1145/3287560.3287586.
- 141 [12] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, R. Zemel, Fairness through awareness,
142 in: Proceedings of the 3rd Innovations in Theoretical Computer Science Conference,
143 ITCS '12, Association for Computing Machinery, New York, NY, USA, 2012, pp. 214–226.
144 doi:10.1145/2090236.2090255.
- 145 [13] S. Wachter, B. Mittelstadt, C. Russell, Why fairness cannot be automated: Bridging the
146 gap between eu non-discrimination law and ai, *Computer Law and Security Review* 41
147 (2020). URL: <https://papers.ssrn.com/abstract=3547922>. doi:10.2139/ssrn.3547922.
- 148 [14] M. Hardt, E. Price, E. Price, N. Srebro, Equality of Opportunity in Supervised Learning, in:
149 *Advances in Neural Information Processing Systems*, volume 29, Curran Associates, Inc.,
150 2016.
- 151 [15] A. Chouldechova, Fair Prediction with Disparate Impact: A Study of Bias in Recidivism
152 Prediction Instruments, *Big Data* 5 (2017) 153–163. doi:10.1089/big.2016.0047.
- 153 [16] R. Berk, H. Heidari, S. Jabbari, M. Kearns, A. Roth, Fairness in criminal justice risk
154 assessments: The state of the art, *Sociological Methods & Research* 50 (2021) 3–
155 44. URL: <https://doi.org/10.1177/0049124118782533>. doi:10.1177/0049124118782533.
156 arXiv:<https://doi.org/10.1177/0049124118782533>.
- 157 [17] R. Zemel, Y. Wu, K. Swersky, T. Pitassi, C. Dwork, Learning Fair Representations, in:
158 *Proceedings of the 30th International Conference on Machine Learning*, PMLR, 2013, pp.
159 325–333.
- 160 [18] M. Padala, S. Gujar, FNNC: Achieving fairness through neural networks, in: *Proceedings*
161 *of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI'20*,
162 *Yokohama, Yokohama, Japan, 2021*, pp. 2277–2283.
- 163 [19] M. Wick, s. panda, J.-B. Tristan, Unlocking Fairness: A Trade-off Revisited, in: *Advances*
164 *in Neural Information Processing Systems*, volume 32, Curran Associates, Inc., 2019.
- 165 [20] W. Alghamdi, S. Asoodeh, H. Wang, F. P. Calmon, D. Wei, K. N. Ramamurthy, Model
166 Projection: Theory and Applications to Fair Machine Learning, in: *2020 IEEE Interna-*
167 *tional Symposium on Information Theory (ISIT)*, IEEE, 2020, pp. 2711–2716. doi:10.1109/
168 *ISIT44484.2020.9173988*.
- 169 [21] D. Wei, K. N. Ramamurthy, F. Calmon, Optimized Score Transformation for Fair Clas-
170 sification, in: *Proceedings of the Twenty Third International Conference on Artificial*
171 *Intelligence and Statistics*, PMLR, 2020, pp. 1673–1683.
- 172 [22] M. Buyl, T. De Bie, The KL-Divergence Between a Graph Model and its Fair I-Projection
173 as a Fairness Regularizer, in: *Machine Learning and Knowledge Discovery in Databases*,
174 Springer International Publishing, 2021, pp. 351–366.
- 175 [23] F. Ding, M. Hardt, J. Miller, L. Schmidt, Retiring adult: New datasets for fair machine
176 learning, *Advances in Neural Information Processing Systems* 34 (2021).
- 177 [24] P. Schubert, M. Kirchner, Ellipse area calculations and their applicability in posturography,
178 *Gait & Posture* 39 (2014) 518–522. doi:10.1016/j.gaitpost.2013.09.001.

179 A. Code Use Examples

```
180 1 import torch
181 2 import torch.nn.functional as F
182 3 from fairret.statistic import TruePositiveRate
183 4 from fairret.loss.violation import NormLoss
184 5
185 6 # The TruePositiveRate class is a subclass of LinearFractionalStatistic.
186 7 statistic = TruePositiveRate()
187 8
188 9 # The fairret modules accept any LinearFractionalStatistic instance.
18910 fairret = NormLoss(statistic)
19011 fairret_strength = 1.0
19112
19213 def train_epoch(train_loader, model, optimizer):
19314     for feat, sens, target in train_loader:
19415         optimizer.zero_grad()
19516
19617         logit = model(feat)
19718         bce_loss = F.binary_cross_entropy_with_logits(logit, target)
19819         fairret_loss = fairret(logit, feat, sens, target)
19920         loss = bce_loss + fairret_strength * fairret_loss
20021         loss.backward()
20122
20223     optimizer.step()
```

Listing 1: Example use of the FAIRRET library in a simple PyTorch setup.

203 Listing 1 displays a code example of how the FAIRRET can easily be deployed in a typical
204 PyTorch [6] setup. It suffices to simply load a subclass of `LinearFractionalStatistic` and
205 pass it on to a FAIRRET implementation instance such as `NormLoss` (as defined in Def. 7). The
206 FAIRRET is then used to compute the quantification of unfairness as a loss like any other in
207 PyTorch. In this case, we use the true positive rate statistic to pursue the fairness notion of
208 equalized opportunity (EO).

209 B. Confidence Ellipses

210 The confidence ellipses we use in Fig. 2 are uncommon in machine learning literature. Yet,
211 they work well for our purpose of comparing trade-offs between metrics that may be noisy
212 depending on randomness during training and dataset split selection.

213 Recall that 1-dimensional confidence intervals typically assume a mean estimator to be
214 normally distributed. The confidence interval then denotes the uncertainty of the sample
215 mean using the standard error. Similarly, confidence ellipses assume a 2-dimensional point,
216 i.e. the 2-dimensional mean estimator, to have a multivariate normal distribution that can be
217 characterized through the sample mean and standard error statistics.

218 Our implementation of the confidence ellipses follows a featured implementation on `matplotlib`².
219 However, a crucial difference is that this implementation computes a confidence interval for a
220 2-dimensional random variable based on the covariance matrix for the standard *deviation* of

²https://matplotlib.org/3.7.0/gallery/statistics/confidence_ellipse.html.

221 samples of that variable. Following observations by Schubert and Kirchner [24], we instead
222 want to show the uncertainty of the mean estimator, which should use the standard deviation
223 of that estimator, i.e. the covariance for the standard *error*. This is accomplished by dividing the
224 covariance matrix in the `matplotlib` implementation by the number of seeds (5) we use in
225 our experiments.