

Motivation

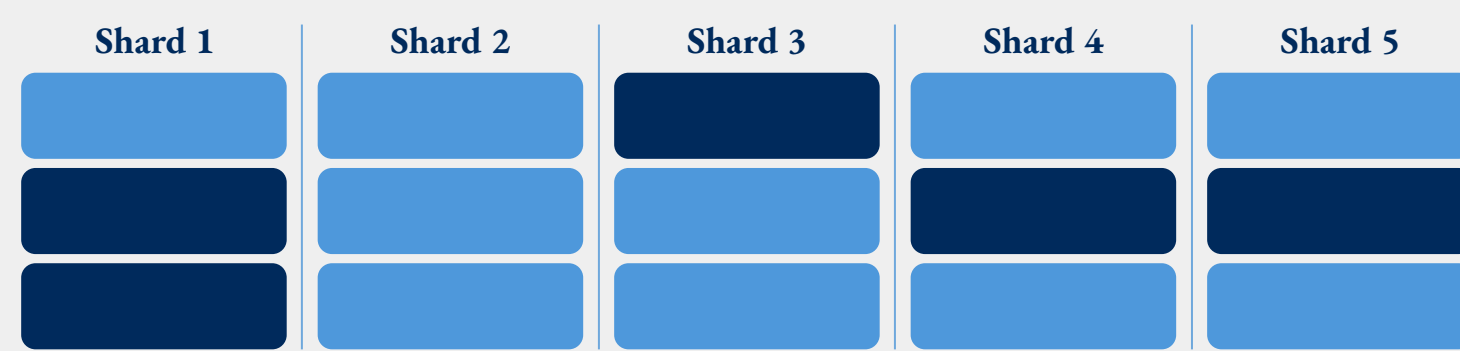
- The speed of unlearning techniques like SISA^[1] can be improved by treating samples with a high unlearning likelihood differently, typically resulting in a noticeable, but acceptable decrease in performance
- Common models may perform worse for minorities, which has been studied in further depth on the example of facial classifiers and racial minorities^[2]
- There is evidence that unlearning likelihoods correlate with belonging to a protected minority

Question: Is model unfairness amplified in SISA unlearning strategies using a-priori estimates?

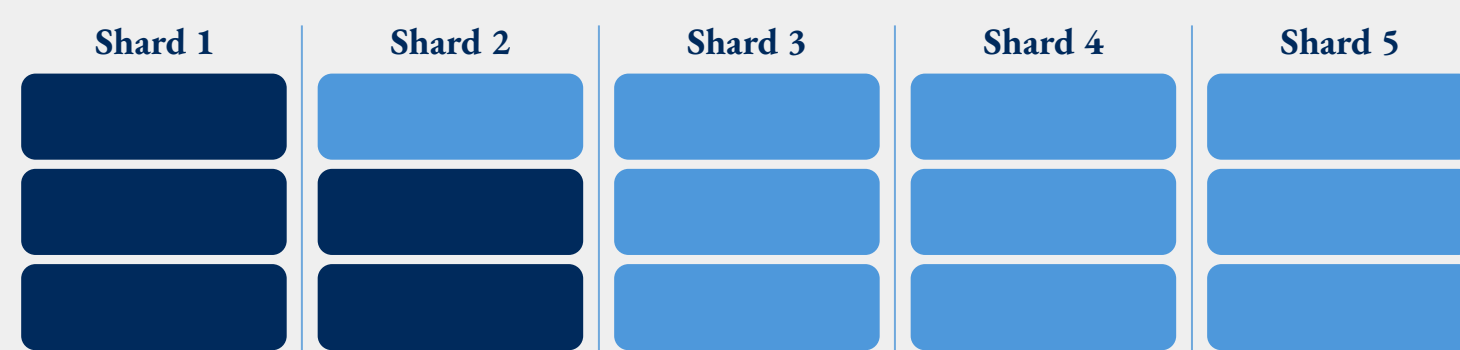
SISA^[1]

- Sharded, Isolated, Sliced and Aggregated** learning trains an ensemble of models on different subsets of the data
- Knowledge about a users individual likelihood to submit an unlearning request allows an adaptive placing of samples in specific shards or slices

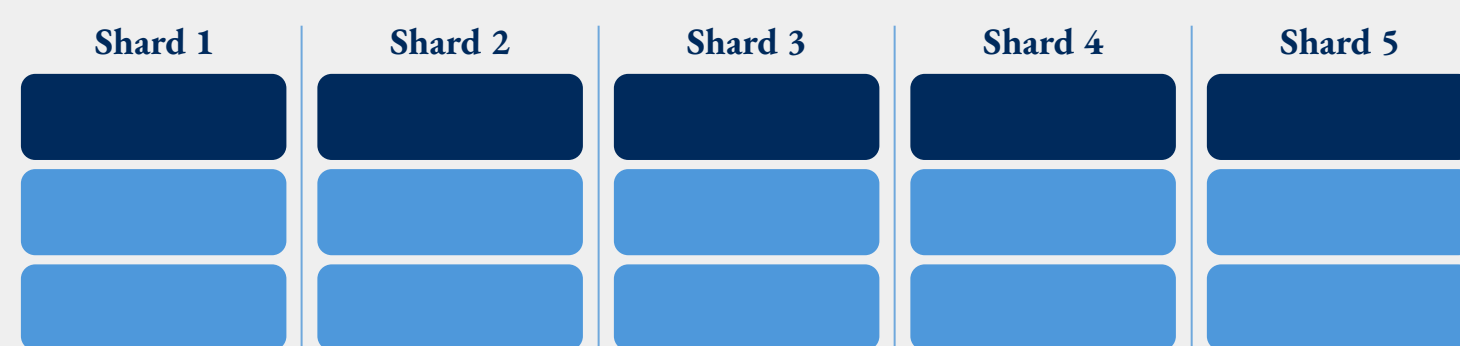
● = low unlearning likelihood ● = high unlearning likelihood



random assignment to shards and slices



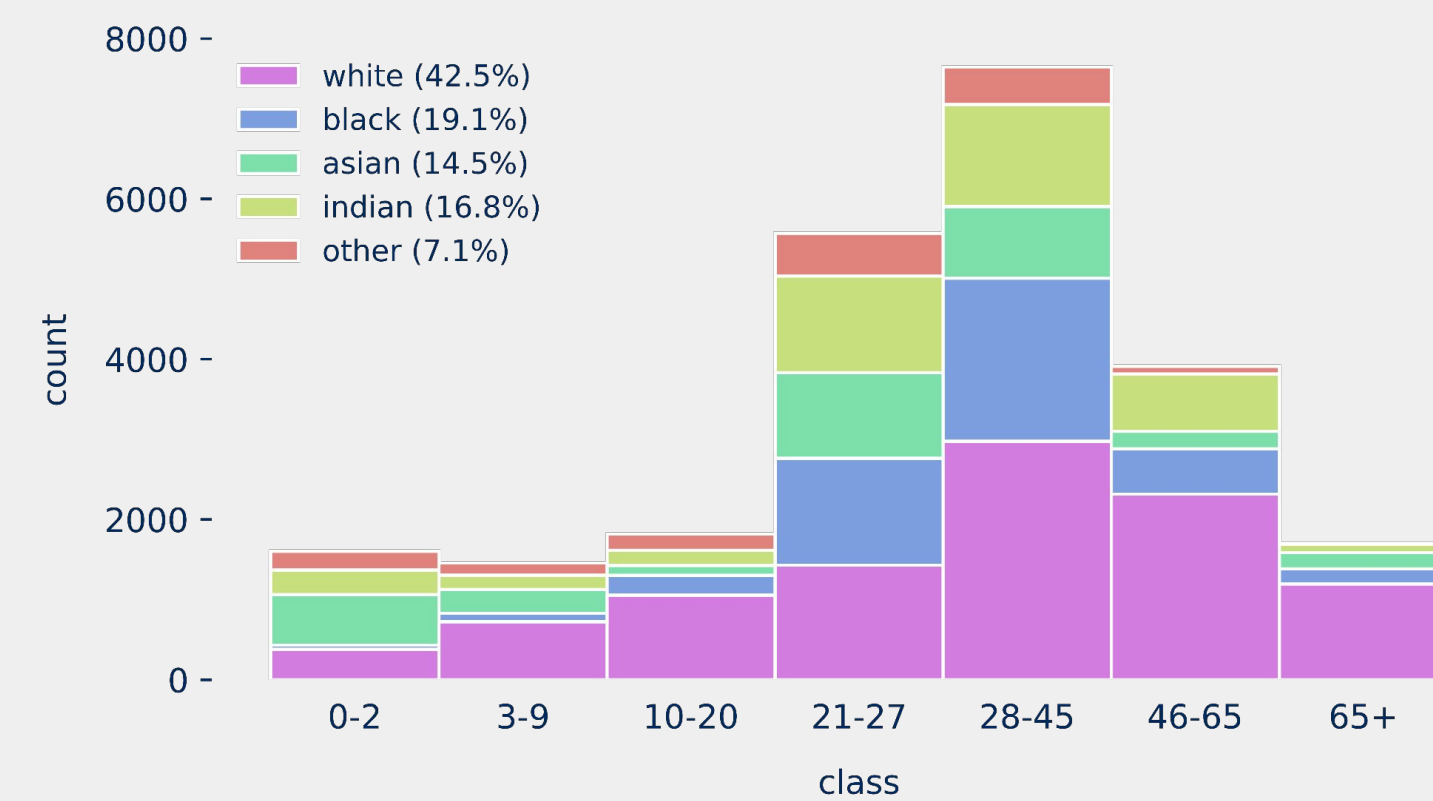
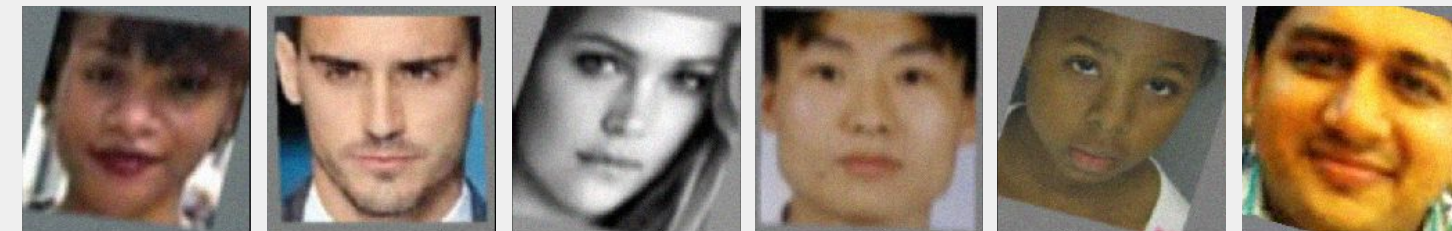
adaptive assignment into fewest possible shards



adaptive assignment into last trained slices

UTKFace^[3]

The **UTKFace** dataset contains 23,708 images of human faces with labels for **age** and **race** (and gender)

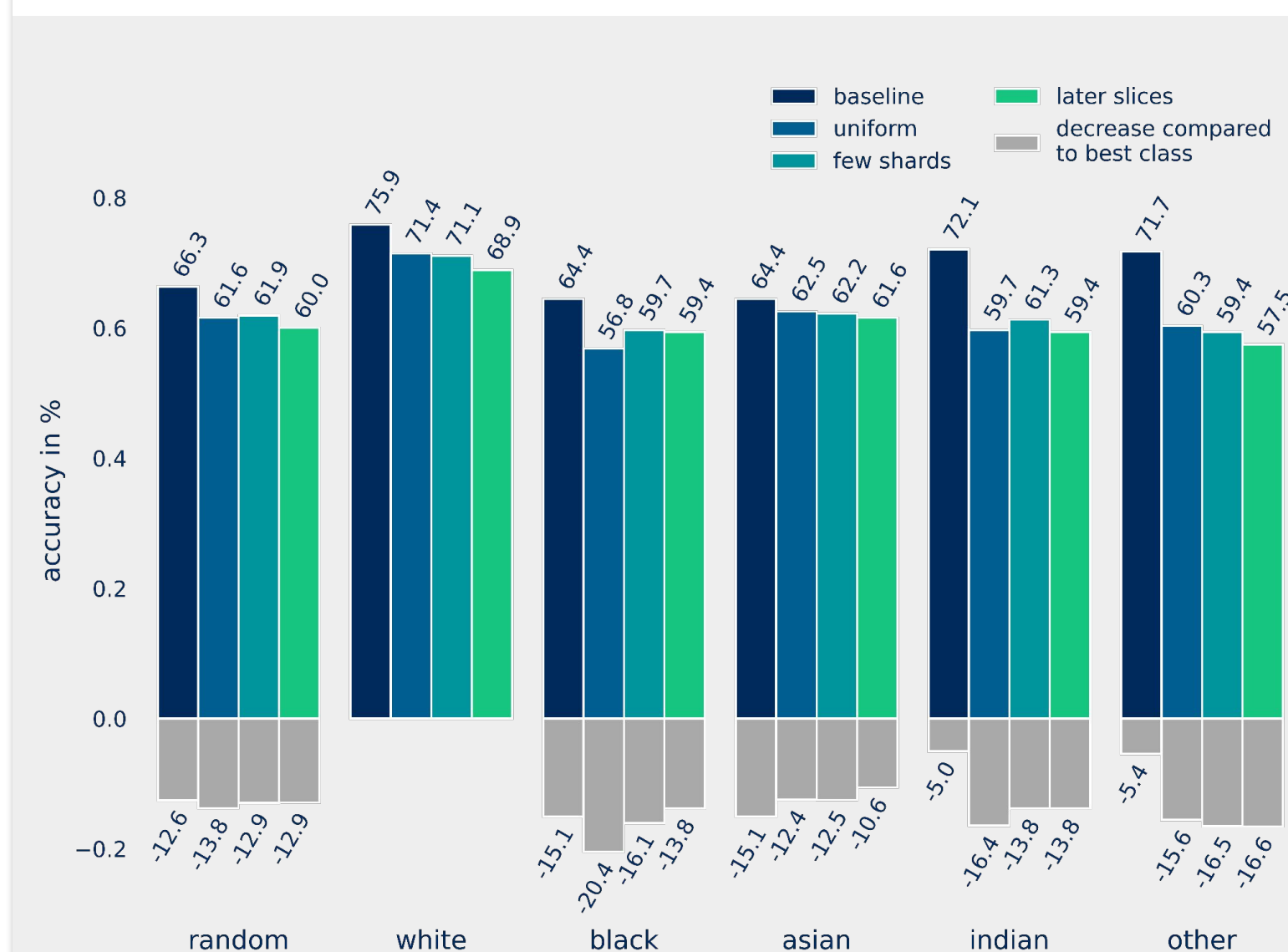


Experiments

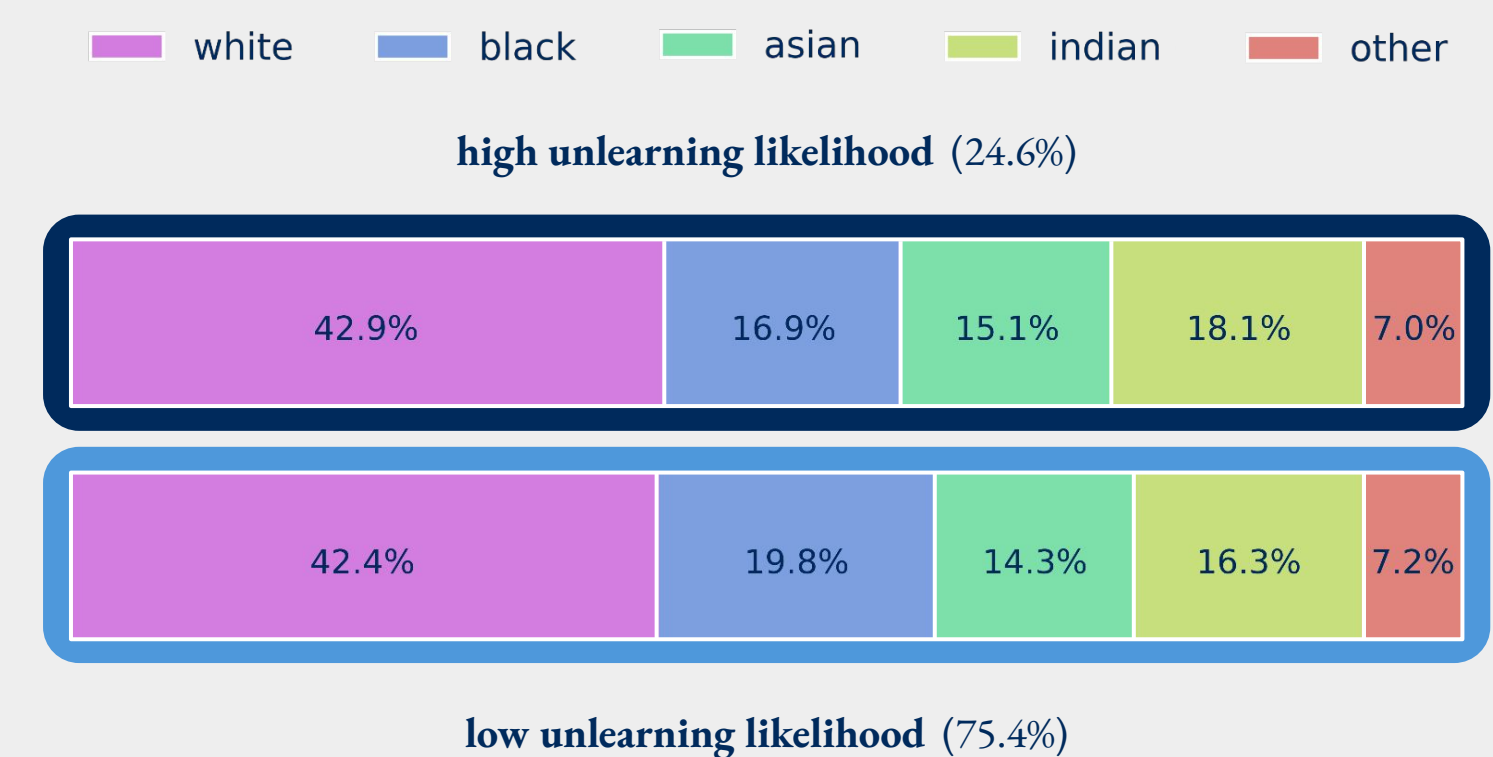
Baseline: monolithic model
SISA: 5 Shards \times 3 Slices (\approx 1,500 samples/slice) aggregation via summation of logits
Methods: uniform, few shards, later slices
Test Set: 9 random samples per class + 9 random samples per class per race
Training: 15 Epochs, lr 6e-4, lr decay of \times 0.2 after 6, 9, 12 epochs

all results are reported as average over 5 runs

Top-1 accuracies averaged over all classes per race and model



Average composition of slices



Possible causes of performance differences

- A)** relative amount of training samples ($r^2 = \underline{0.51}, \underline{0.81}, \underline{0.93}, \underline{0.94}$)
B) uniformity of class distribution ($r^2 = \underline{-0.22}, \underline{-0.10}, \underline{0.15}, \underline{0.24}$)
C) likelihood of unlearning indicator ($r^2 = \underline{0.22}, \underline{0.45}, \underline{0.43}, \underline{0.41}$)

correlation \neq causation

Accuracy matrices for race and age



Conclusions

- There are considerable performance differences across races in both the baseline and all SISA models
- Realistic modelling of unlearning predictors results in only minor distribution shifts across slices
- Weaknesses of the baseline model are inherited, but not necessarily amplified by all SISA surrogates

References

- [1] Lucas Bourtole et al. "Machine Unlearning". In: *CoRR* abs/1912.03817 (2019). arXiv: 1912.03817. URL: <http://arxiv.org/abs/1912.03817>.
- [2] Joy Buolamwini and Timnit Gebru. "Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification". In: *Conference on Fairness, Accountability and Transparency, FAT 2018, 23-24 February 2018, New York, NY, USA*. Ed. by Sorelle A. Friedler and Christo Wilson. Vol. 81. Proceedings of Machine Learning Research. PMLR, 2018, pp. 77–91. URL: <http://proceedings.mlr.press/v81/buolamwini18a.html>
- [3] Zhiwei Zhang, Yang Song, and Hairong Qi. "Age Progression/Regression by Conditional Adversarial Autoencoder". In: *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*. IEEE Computer Society, 2017, pp. 4352–4360. DOI: 10.1109/CVPR.2017.463. URL: <https://doi.org/10.1109/CVPR.2017.463>.
- [4] *Special Eurobarometer 487a*. (2019). ISBN: 978-92-76-08384-9. DOI: 10.2838/579882. URL: <https://europa.eu/eurobarometer/surveys/detail/2222>.