

# Machine Unlearning via Algorithmic Stability <sup>\*</sup>

Enayat Ullah <sup>†</sup>    Tung Mai <sup>‡</sup>    Anup Rao <sup>§</sup>    Ryan Rossi <sup>¶</sup>    Raman Arora <sup>||</sup>

We study the problem of machine unlearning and identify a notion of algorithmic stability, Total Variation (TV) stability, which we argue, is suitable for the goal of *exact* unlearning. For convex risk minimization problems, we design TV-stable algorithms based on noisy Stochastic Gradient Descent (SGD). Our key contribution is the design of corresponding *efficient* unlearning algorithms, which are based on constructing a (maximal) coupling of Markov chains for the noisy SGD procedure. To understand the trade-offs between accuracy and unlearning efficiency, we give upper and lower bounds on excess empirical and populations risk of TV stable algorithms for convex risk minimization. Our techniques generalize to arbitrary non-convex functions, and our algorithms are differentially private as well.

---

<sup>\*</sup>The full version of the paper can be found at <https://arxiv.org/abs/2102.13179>

<sup>†</sup>Johns Hopkins University. Email: [enayat@jhu.edu](mailto:enayat@jhu.edu)

<sup>‡</sup>Adobe Research. Email: [tumai@adobe.com](mailto:tumai@adobe.com)

<sup>§</sup>Adobe Research. Email: [anuprao@adobe.com](mailto:anuprao@adobe.com)

<sup>¶</sup>Adobe Research. Email: [rrossi@adobe.com](mailto:rrossi@adobe.com)

<sup>||</sup>Johns Hopkins University. Email: [arora@cs.jhu.edu](mailto:arora@cs.jhu.edu)