

Replicability Crisis

nature

Published: 25 May 2016

1,500 scientists lift the lid on reproducibility

Monya Baker

- Over 70% of researchers failed to replicate others' work
- Over 50% failed to replicate their own work!



- 2019 NeurIPS / ICLR Reproducibility Challenge (github.com/reproducibility-challenge)
- Ongoing ML Reproducibility Challenge (paperswithcode.com/rc2022)

Goal: Mathematical Study of Replicability

- X data domain
- \mathcal{D} distribution over X
- $S_1, S_2 \sim_{i.i.d.} \mathcal{D}^n$ size n datasets
- ξ random binary string

Definition (Replicable Algorithm) [Impagliazzo, Lei, Pitassi, Sorrell '22]

A randomized algorithm $\mathcal{A} : X^n \rightarrow Y$ is ρ -replicable if

$$\Pr_{S_1, S_2, \xi} [\mathcal{A}(S_1; \xi) = \mathcal{A}(S_2; \xi)] \geq 1 - \rho.$$

PAC, Differentially Private, Online, and SQ Learning

- PAC Learning:** for all α, β , given a sufficiently large dataset S , the learner outputs a classifier $\hat{h} = \mathcal{A}(S)$ satisfying $\Pr_{(x,y)} [\hat{h}(x) \neq y] \leq \alpha$ with probability $1 - \beta$ over the draw of S .
- DP Learning:** PAC learning requirements and (ϵ, δ) -DP requirements, i.e. for all neighboring datasets S, S' and for all events E it holds that $\Pr[\mathcal{A}(S) \in E] \leq e^\epsilon \Pr[\mathcal{A}(S') \in E] + \delta$.
- Online Learning:** adversary picks "hard" function h^* and in every round t gives x_t to the learner; learner guesses a label \hat{y}_t and makes a mistake if $\hat{y}_t \neq h^*(x_t)$. The goal of the learner is to minimize # of mistakes.
- SQ Learning:** instead of getting labeled samples as input, the learner has access to a (noisy) oracle that can answer statistical queries about the target.

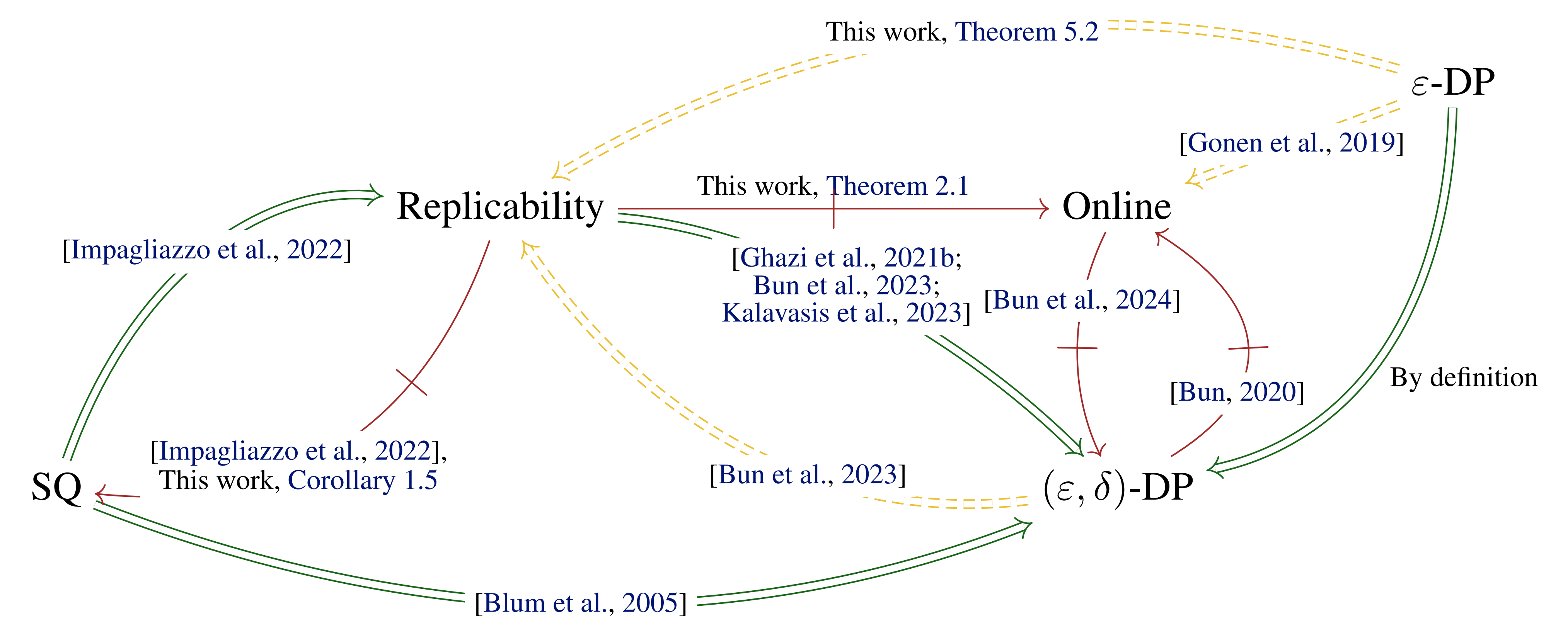
Main Results

Black-Box Transformations

- A pure DP learner can be efficiently* transformed into an online learner.

Separations under Cryptographic Assumptions

- There is a concept class that can be efficiently learned by a replicable PAC learner, but not an efficient online learner.
- There is a concept class that can be efficiently† learned by a replicable PAC learner, but not an efficient SQ learner.



Future Work

- Is there a computationally efficient transformation from online learners to replicable learners?
- *Can we derive replicable learners from pure DP learners which are efficient with respect to the complexity of the underlying concept class?
- †Can we design efficient replicable algorithms for every distribution?



SCAN ME