

Figure F2. Comparison of convergence rates for optimization algorithms. The figure illustrates typical error decay behaviors: linear convergence (exponential decay), sublinear convergence (e.g., $1/k$), and superlinear convergence (faster than exponential). Additionally, three finite-step convergence scenarios are shown, where error drops to zero suddenly after a fixed number of iterations (at step 10, 20, and 30 respectively), representing idealized convergence in a finite number of steps. The y -axis is shown on a logarithmic scale to highlight differences in decay speed.

If $\delta_{K+1} \neq \delta^\Delta$, then it satisfies the initialization condition of Theorem 3.1, and the subsequent sequence δ_{K+i} for $i = 1, 2, \dots$ converges at least at a linear rate.

If $\delta_{K+1} = \delta^\Delta$, the sequence δ_k exhibits finite-step convergence. Prior works [1–5] show that such convergence is effectively instantaneous and can be considered faster than linear (See Figure F2). This does not contradict our conclusion, as our paper emphasizes that PGD converges at least at a linear rate, and the lower bound applies under the specific initialization assumed in Theorem 3.1.

In summary, PGD achieves global convergence and eventually at least local linear convergence even under general initialization.

Reference:

- [1] Nocedal, Jorge, and Stephen J. Wright, eds. Numerical optimization. New York, NY: Springer New York, 1999.
- [2] Wright S J. Primal-dual interior-point methods[M]. Society for Industrial and Applied Mathematics, 1997.
- [3] Trefethen L N, Bau D. Numerical linear algebra[M]. Society for Industrial and Applied Mathematics, 2022.
- [4] Dennis Jr J E, Schnabel R B. Numerical methods for unconstrained optimization and nonlinear equations[M]. Society for Industrial and Applied Mathematics, 1996.
- [5] Han, Ningning, Jian Lu, and Shidong Li. "The finite steps of convergence of the fast thresholding algorithms with f-feedbacks in compressed sensing." Numerical Algorithms 90.3 (2022): 1197-1223.

E.2. Experiments on Logistic Regression Loss for Q3

Logistic Regression Loss (used for binary classification):

$$\Phi(\delta) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \mathbf{w}^\top (\mathbf{x}_i + \delta))), \quad (3)$$

where \mathbf{w} is given vectors, \mathbf{x}_i denotes the input features, and $y_i \in \{-1, 1\}$ are the labels. \mathbf{w} , \mathbf{x}_i and y_i are generated for simulation. Then, PGD is applied to solve the constrained optimization problem:

$$\min_{\delta \in \mathbb{B}(\mathbf{0}, \epsilon)} \Phi(\delta),$$

where $\Phi(\delta)$ represents the Logistic Regression Loss (3). This function is convex and smooth in δ , as it is a composition of the convex and smooth function $\log(1 + \exp(\cdot))$ with a linear transformation. The gradient does not vanish in the small ball. Therefore, there is no critical point inside the interior of the constraint set. No prior knowledge of the optimal solution δ^Δ is assumed, and the initialization is randomly sampled from the constraint ball $\mathbb{B}(\mathbf{0}, \epsilon)$.

We consider various dimensions $d \in \{10, 50, 100, 200\}$ and constraint radius $\epsilon \in \{0.1, 0.5, 1, 5\}$. The experimental results for Logistic Regression Loss are shown in Figures F3. As illustrated, PGD demonstrates both global convergence and local linear convergence, which is consistent with our theoretical analysis.

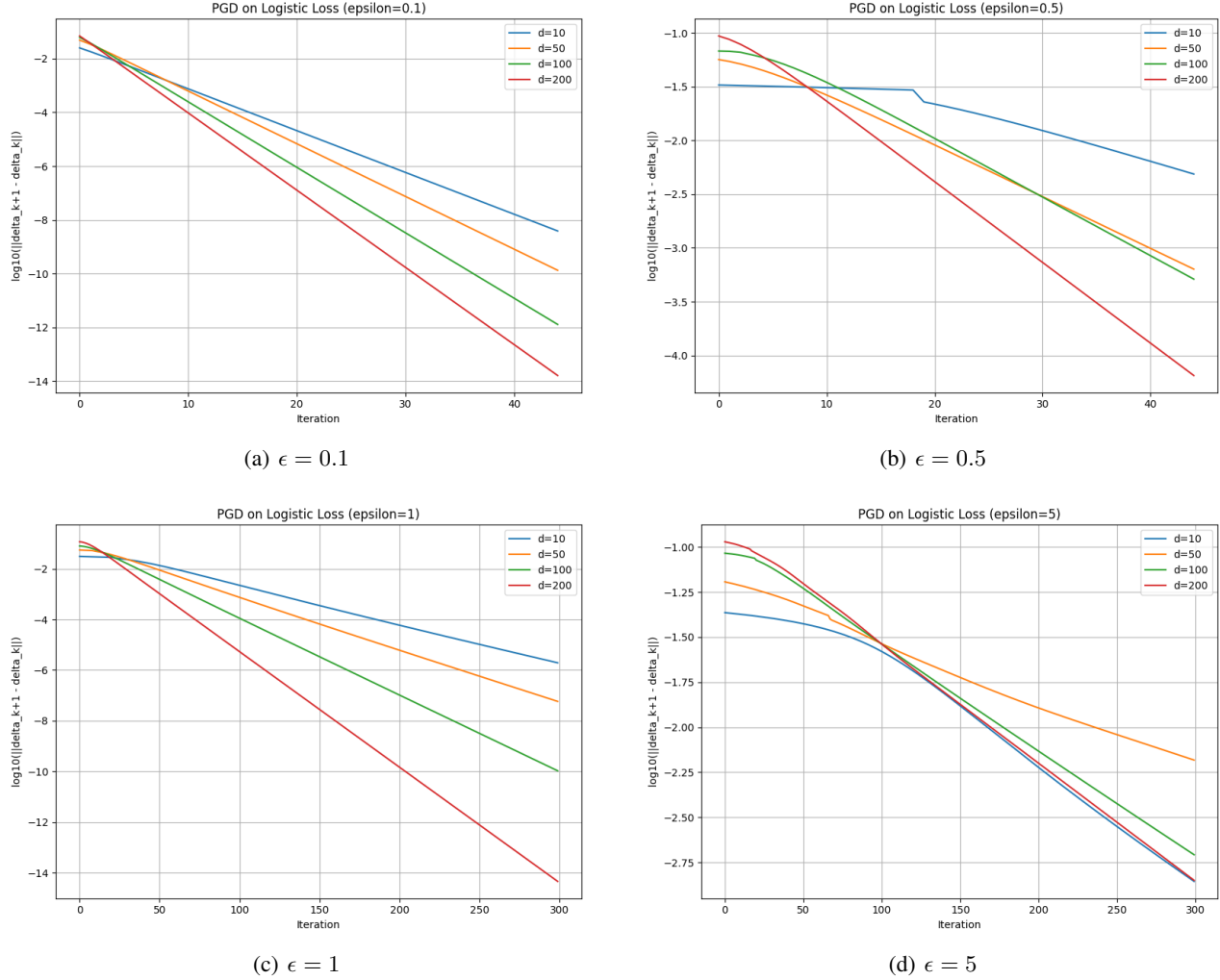


Figure F3. The convergence of δ_k produced by PGD for solving optimization Logistic Regression Loss (3) with different dimensions d . We use a logarithmic coordinate system where the horizontal axis represents the number of iterations k and the vertical axis represents $\log_{10}(\|\delta_{k+1} - \delta_k\|)$.