

Average-case acceleration for bilinear games and normal matrices

Carles Domingo-Enrich^a, Fabian Pedregosa^b, Damien Scieur^c

^aCourant Institute (NYU), ^bGoogle Research, ^cSamsung SAIT AI Lab & Mila

ICLR 2021

May, 2021

Background

- Traditional theory in optimization is worst-case analysis: Not representative of typical behavior.
- Optimal average-case methods have recently been developed for quadratic minimization problems [Berthier et al., 2020, Pedregosa and Scieur, 2020, Lacotte and Pilanci, 2020].
- Optimal methods for smooth games only exist for the worst-case analysis [Azizian et al., 2020].
- **Gap:** Average-case optimal methods for smooth games.

Contributions of the paper

We combine average-case analysis with smooth games.

1. We develop novel average-case optimal algorithms for finding the root of a linear system determined by a (potentially non-symmetric) normal matrix.
2. We show that solving the Hamiltonian using an average-case optimal method is optimal to find equilibria in bilinear games.

Framework (1/2)

- For $\mathbf{A} \in \mathbb{R}^{d \times d}$ and $\mathbf{x}^* \in \mathbb{R}^d$, we consider the non-symmetric operator problem (**NSO**):

$$\text{Find } \mathbf{x} : F(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{A}(\mathbf{x} - \mathbf{x}^*) = \mathbf{0}.$$

- We define

$$\text{dist}(\mathbf{x}, \mathcal{X}^*) \stackrel{\text{def}}{=} \min_{\mathbf{v} \in \mathcal{X}^*} \|\mathbf{x} - \mathbf{v}\|_2, \quad \text{with } \mathcal{X}^* = \{\mathbf{x} \in \mathbb{R}^d \mid \mathbf{A}(\mathbf{x} - \mathbf{x}^*) = \mathbf{0}\}.$$

Framework (2/2)

- **First-order methods:** For all $t \geq 0$,

$$\mathbf{x}_t = \mathbf{x}_0 + \sum_{i=0}^t \alpha_{t,i} F(\mathbf{x}_i),$$

for some coefficients $\alpha_{t,i}$.

- For some random \mathbf{A} , \mathbf{x}^* and initialization \mathbf{x}_0 , average-case first-order optimal algorithms solve:

$$\min_{\mathbf{x}_t \text{ first-order method}} \mathbb{E}_{(\mathbf{A}, \mathbf{x}^*, \mathbf{x}_0)} \text{dist}(\mathbf{x}_t, \mathcal{X}^*)$$

Sketch of the theory

- *Residual polynomial*: polynomial P that satisfies $P(0) = 1$.
- Known: if (\mathbf{x}_t) is the sequence generated by a first-order method, there exist residual polynomials P_t of degree at most t such that $\mathbf{x}_t - \mathbf{x}^* = P_t(\mathbf{A})(\mathbf{x}_0 - \mathbf{x}^*)$.
- Empirical spectral distribution of \mathbf{A} : $\hat{\mu}_{\mathbf{A}}(\lambda) = \frac{1}{d} \sum_{i=1}^d \delta_{\lambda_i}(\lambda)$, where $(\lambda_i)_{i=1}^d$ are the eigenvalues of \mathbf{A} . Expected spectral distribution: $\mu_{\mathbf{A}} = \mathbb{E}_{\mathbf{A}} \hat{\mu}_{\mathbf{A}}(\lambda)$.
- If \mathbf{A} is a random normal matrix and \mathbf{x}^* are sampled appropriately, we show that **for any first order method with associated polynomials (P_t) , we have**
$$\mathbb{E}[\text{dist}(\mathbf{x}_t, \mathcal{X}^*)] = R^2 \int_{\mathbb{C} \setminus \{0\}} |P_t|^2 d\mu_{\mathbf{A}}.$$
- For simple measures μ , we can compute the sequence of residual polynomials that optimize $\int_{\mathbb{C} \setminus \{0\}} |P_t|^2 d\mu_{\mathbf{A}}$, and their corresponding first-order methods.

Average-case optimal methods for bilinear games

We want to find a Nash equilibrium of the zero-sum minimax game given by

$$\min_{\theta_1} \max_{\theta_2} \ell(\theta_1, \theta_2) \stackrel{\text{def}}{=} (\theta_1 - \theta_1^*)^\top \mathbf{M} (\theta_2 - \theta_2^*).$$

where $\theta_1, \theta_1^* \in \mathbb{R}^{d_1}$, $\theta_2, \theta_2^* \in \mathbb{R}^{d_2}$, $\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$. Defining

$$\mathbf{A} = \begin{bmatrix} 0 & \mathbf{M} \\ -\mathbf{M}^\top & 0 \end{bmatrix},$$

we recast the problem as solving the NSO $F(\mathbf{x}) \stackrel{\text{def}}{=} \mathbf{A}(\mathbf{x} - \mathbf{x}^*) = 0$.

Example: M with i.i.d components

Setting: Each entry of M is sampled from iid from distribution with mean 0 and variance σ^2 , in the regime $d_1, d_2 \rightarrow \infty$, $d_1/d_2 = r$.

Optimal average-case algorithm.

Initialization. $\mathbf{x}_{-1} = \mathbf{x}_0 = (\boldsymbol{\theta}_{1,0}, \boldsymbol{\theta}_{2,0})$.

Main loop. For $t \geq 0$,

$$\mathbf{g}_t = F(\mathbf{x}_t - F(\mathbf{x}_t)) - F(\mathbf{x}_t) \quad (= \frac{1}{2} \nabla \|F(\mathbf{x}_t)\|^2)$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - h_{t+1} \mathbf{g}_t + m_{t+1} (\mathbf{x}_{t-1} - \mathbf{x}_t) \quad \text{where}$$

$$h_t = -\frac{\delta_t}{\sigma^2 \sqrt{r}}, \quad m_t = 1 + \rho \delta_t, \quad \rho = \frac{1+r}{\sqrt{r}}, \quad \delta_t = (-\rho - \delta_{t-1})^{-1}, \quad \delta_0 = 0.$$

Example: M with i.i.d components

Bilinear Problems

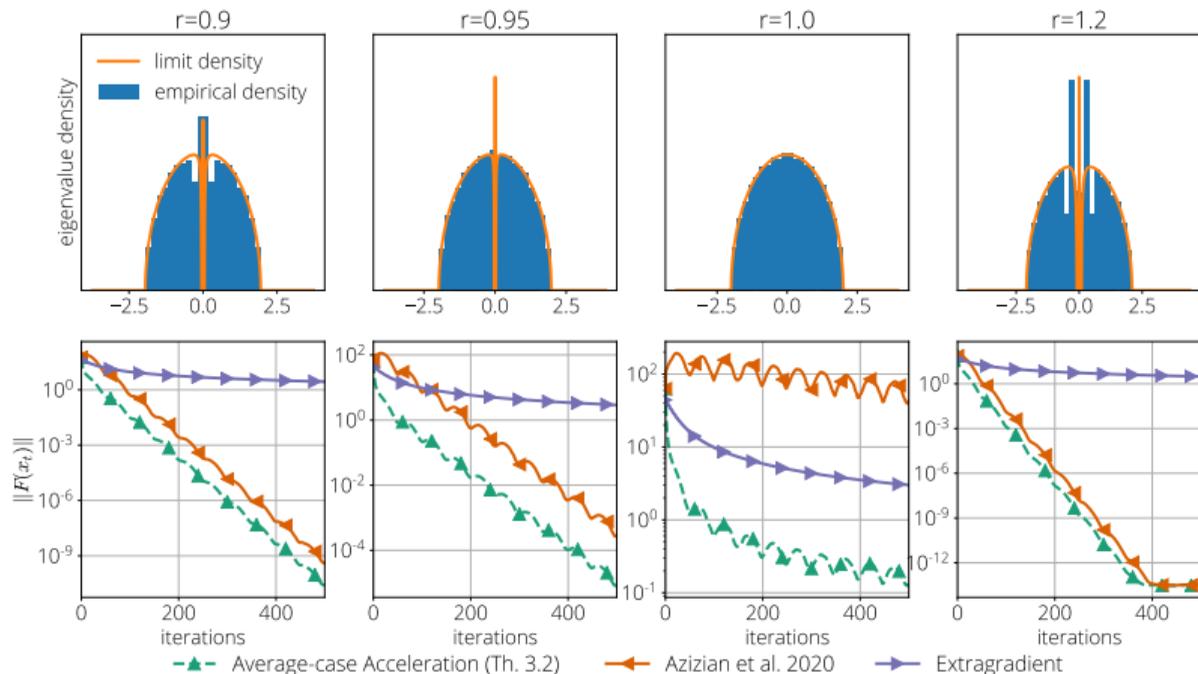


Figure: *First row:* spectral density associated with bilinear games for varying values of the ratio parameter $r = n/d$ (the x-axis represents the imaginary line). *Second row:* Comparison of gradient norm decay with benchmark. The largest gain is in the ill-conditioned regime ($r \approx 1$).

Normal matrices with circular spectral distribution

Setting: Assume that the expected spectral distribution $\mu_{\mathbf{A}}$ is the uniform probability measure on the complex disk of center $C \in \mathbb{R}$, $C > 0$ and radius $R < C$.

Optimal average-case algorithm.

Initialization. $y_{-1} = y_0 = x_0$.

Main loop. For $t \geq 0$,

$$y_t = y_{t-1} - \frac{1}{C} F(y_{t-1}), \quad \beta_t = \left(\frac{C}{R}\right)^{2t} (t+1), \quad B_t = B_{t-1} + \beta_{t-1},$$

$$x_t = \frac{B_t}{B_t + \beta_t} x_{t-1} + \frac{\beta_t}{B_t + \beta_t} y_t.$$

Moreover, $\mathbb{E}_{(\mathbf{A}, x^*, x_0)} \text{dist}(x_t, \mathcal{X}^*)$ converges to zero at rate $1/B_t$.

Uniform circular spectral distribution

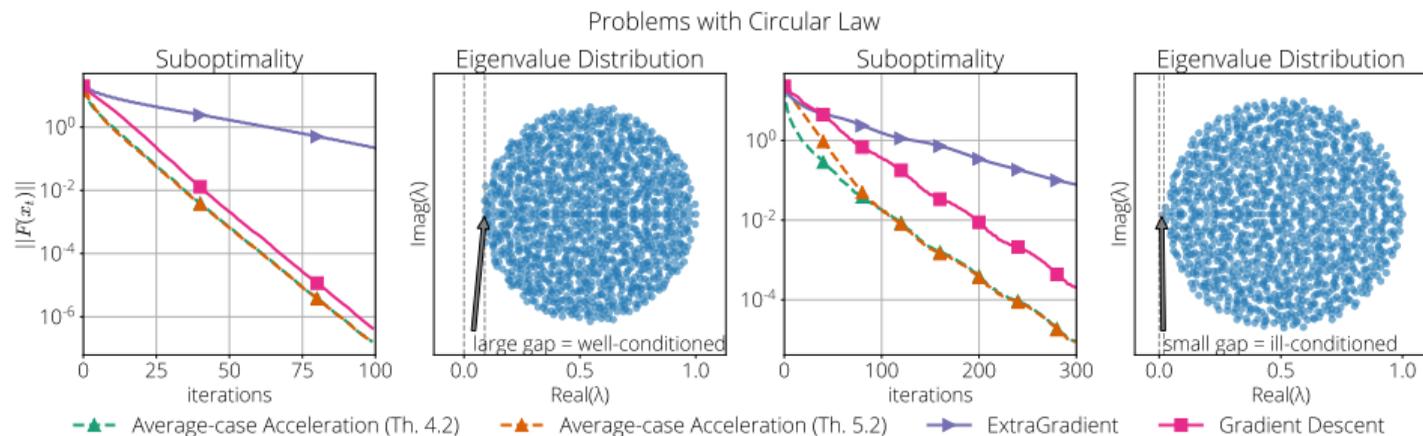


Figure: Benchmarks (columns 1 and 3) and eigenvalue distribution of a design matrix generated with iid entries for two different degrees of conditioning. Despite the normality assumption not being satisfied, we still observe an improvement of average-case optimal methods vs worst-case optimal ones.

References

-  Azizian, W., Scieur, D., Mitliagkas, I., Lacoste-Julien, S., and Gidel, G. (2020). Accelerating smooth games by manipulating spectral shapes. *In Proceedings of Machine Learning Research.*
-  Berthier, R., Bach, F., and Gaillard, P. (2020). Accelerated gossip in networks of given dimension using Jacobi polynomial iterations. *SIAM Journal on Mathematics of Data Science*, 2(1):24–47.
-  Lacotte, J. and Pilanci, M. (2020). Optimal randomized first-order methods for least-squares problems. *Proceedings of the 37th International Conference on Machine Learning.*
-  Pedregosa, F. and Scieur, D. (2020). Average-case acceleration through spectral density estimation. *In Proceedings of the 37th International Conference on Machine Learning.*