

Deep Linear Networks

Rathindra Nath Karmakar
Session 2 : Implicit Acceleration

References

- On the Optimization of Deep Networks: Implicit Acceleration by Overparameterization, Sanjeev Arora, Nadav Cohen, Elad Hazan (2018)
- The geometry of the deep linear network, Govind Menon (2024)

Recap: Key Questions for Gradient Flow

For the gradient flow dynamics on a loss surface $\mathcal{L}(\mathbf{W})$:

$$\frac{d}{dt}\mathbf{W}(t) = -\nabla_{\mathbf{W}}\mathcal{L}(\mathbf{W}(t))$$

We want to understand:

- Convergence guarantees? (Yes, for *balanced* cases)
- Convergence rate?**
- Characterization of the minimizer reached?
- Effect of noise and discretization?

Today, we address the second question: can overparameterization improve the **rate of convergence**?

The Counter-Intuitive Message

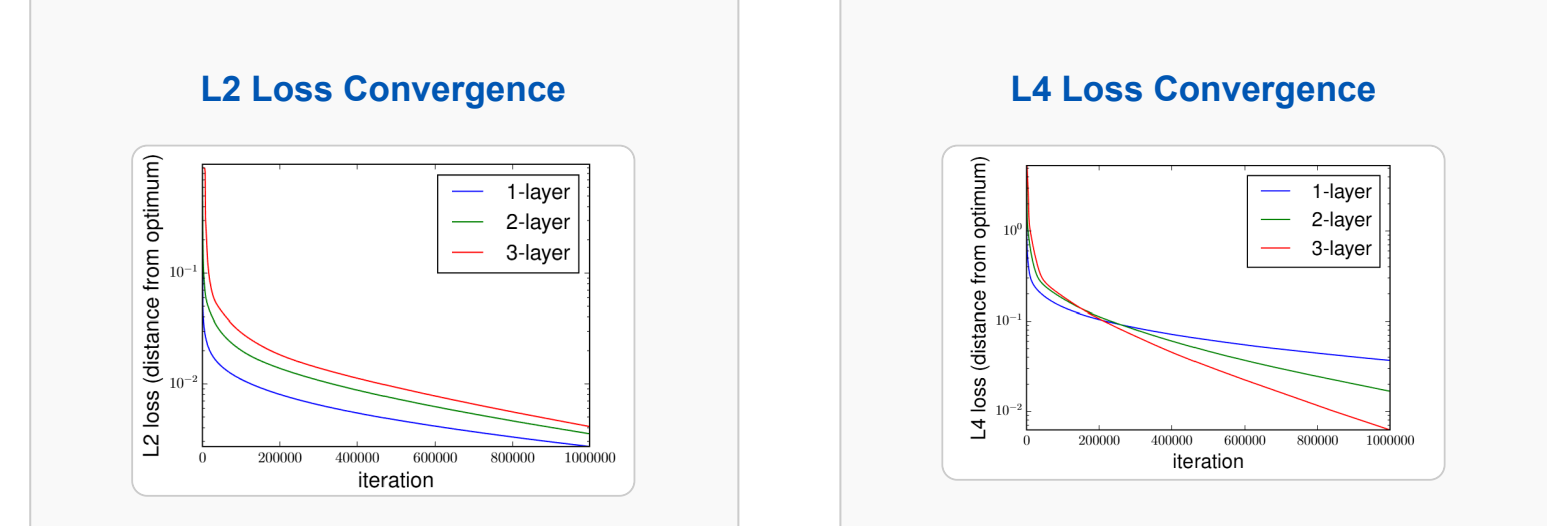
The conventional wisdom is that increasing depth improves expressiveness but complicates optimization. Today's message, based on Arora, Cohen, & Hazan (2018), is that sometimes:

Increasing depth can *accelerate* optimization.

We will see that overparameterization via depth can act as an implicit preconditioner, combining the effects of momentum and adaptive learning rates.

Empirical Evidence: Acceleration with Depth

Comparison of gradient descent on a linear regression task with ℓ_p loss. (Fig 3 from Arora et al. 2018)



- Left (ℓ_2 loss):** Increasing depth slightly "slows down" convergence, consistent with prior work.
- Right (ℓ_4 loss):** Increasing depth provides a "massive acceleration", outperforming the shallow model significantly.

The End-to-End Update Rule

Let $\mathbf{W}_e = \mathbf{W}_N \mathbf{W}_{N-1} \cdots \mathbf{W}_1$ be the end-to-end matrix. The paper shows that its gradient flow dynamics can be written purely in terms of \mathbf{W}_e itself.

Theorem 1 (Arora et al. 2018). Under the balanced initialization assumption, $\mathbf{W}_{j+1}(t_0)^T \mathbf{W}_{j+1}(t_0) = \mathbf{W}_j(t_0) \mathbf{W}_j(t_0)^T$, the gradient flow for \mathbf{W}_e is:

$$\dot{\mathbf{W}}_e(t) = -\eta \sum_{j=1}^N [\mathbf{W}_e(t) \mathbf{W}_e(t)^T]^{\frac{N-j}{N}} \cdot \nabla \mathcal{L}(\mathbf{W}_e(t)) \cdot [\mathbf{W}_e(t)^T \mathbf{W}_e(t)]^{\frac{j-1}{N}}$$

The standard gradient $\nabla \mathcal{L}(\mathbf{W}_e)$ is pre-multiplied and post-multiplied by fractional matrix powers of $\mathbf{W}_e \mathbf{W}_e^T$ and $\mathbf{W}_e^T \mathbf{W}_e$.

A Geometric Interpretation

The complex update rule from Theorem 1 seems arbitrary. But it has a beautiful geometric meaning. Let's define the linear preconditioning operator from Theorem 1 as $\mathbf{A}_{N, \mathbf{W}_e}$:

$$\mathbf{A}_{N, \mathbf{W}_e}(Z) = \sum_{k=1}^N (\mathbf{W}_e \mathbf{W}_e^T)^{\frac{N-k}{N}} Z (\mathbf{W}_e^T \mathbf{W}_e)^{\frac{k-1}{N}}$$

So the dynamics are $\dot{\mathbf{W}}_e = -\mathbf{A}_{N, \mathbf{W}_e}(E'(\mathbf{W}_e))$, where $E'(\mathbf{W}_e) = \nabla \mathcal{L}(\mathbf{W}_e)$ is the standard Euclidean gradient.

Following Menon (Sec 4.3), we define a position-dependent Riemannian metric g^N on the tangent space at \mathbf{W}_e via the inner product:

Definition: A New Metric.

$$g^N(\mathbf{W}_e)(Z_1, Z_2) := \text{Tr}(Z_1^T (\mathbf{A}_{N, \mathbf{W}_e})^{-1} Z_2)$$

(Assuming \mathbf{W}_e has full rank, $\mathbf{A}_{N, \mathbf{W}_e}$ is invertible.)

Under this specific metric, the complicated dynamics of Theorem 1 become a simple and natural **Riemannian gradient flow**:

$$\dot{\mathbf{W}}_e(t) = -\text{grad}_{g^N} E(\mathbf{W}_e(t))$$

Interpreting the Dynamics (Single Output Case)

For the special case of a single output ($k = 1$), the preconditioning scheme simplifies to a more intuitive form.

Claim 2 (Arora et al. 2018). For a single output network, the end-to-end dynamics are equivalent to:

$$\dot{\mathbf{W}}_e = \underbrace{-\eta \|\mathbf{W}_e\|^{2-\frac{2}{N}}}_{\text{Adaptive LR}} \left(\nabla \mathcal{L}(\mathbf{W}_e) + \underbrace{(N-1) \text{Pr}_{\mathbf{W}_e}(\nabla \mathcal{L}(\mathbf{W}_e))}_{\text{Momentum-like term}} \right)$$

where $\text{Pr}_{\mathbf{W}_e}(V)$ is the orthogonal projection of vector V onto the direction of vector \mathbf{W}_e .

- Adaptive Learning Rate:** As the weight vector $\|\mathbf{W}_e\|$ grows (moves away from zero init), the effective learning rate increases.
- Momentum:** The gradient is amplified along the direction of the current weight vector \mathbf{W}_e . Since \mathbf{W}_e is the integral of past updates, this promotes movement along the historical trajectory.

Detailed Derivations

Proof of Theorem 1 (N=2 case)

Let's derive the end-to-end dynamics for $\mathbf{W}_e = \mathbf{W}_2 \mathbf{W}_1$. The time derivative is $\dot{\mathbf{W}}_e = \dot{\mathbf{W}}_2 \mathbf{W}_1 + \mathbf{W}_2 \dot{\mathbf{W}}_1$.

- Substitute the gradient flow dynamics for each layer:

$$\dot{\mathbf{W}}_1 = -\eta \mathbf{W}_1^T \nabla \mathcal{L}(\mathbf{W}_e) \quad , \quad \dot{\mathbf{W}}_2 = -\eta \nabla \mathcal{L}(\mathbf{W}_e) \mathbf{W}_1^T$$

- Plug these into the expression for $\dot{\mathbf{W}}_e$:

$$\dot{\mathbf{W}}_e = (-\eta \nabla \mathcal{L}(\mathbf{W}_e) \mathbf{W}_1^T) \mathbf{W}_1 + \mathbf{W}_2 (-\eta \mathbf{W}_2^T \nabla \mathcal{L}(\mathbf{W}_e))$$

- Rearrange the terms:

$$\dot{\mathbf{W}}_e = -\eta ((\mathbf{W}_2 \mathbf{W}_2^T) \nabla \mathcal{L}(\mathbf{W}_e) + \nabla \mathcal{L}(\mathbf{W}_e) (\mathbf{W}_1^T \mathbf{W}_1))$$

- Using the identities for the N=2 case, $\mathbf{W}_2 \mathbf{W}_2^T = (\mathbf{W}_e \mathbf{W}_e^T)^{1/2}$ and $\mathbf{W}_1^T \mathbf{W}_1 = (\mathbf{W}_e^T \mathbf{W}_e)^{1/2}$:

$$\dot{\mathbf{W}}_e = -\eta \left((\mathbf{W}_e \mathbf{W}_e^T)^{1/2} \nabla \mathcal{L}(\mathbf{W}_e) + \nabla \mathcal{L}(\mathbf{W}_e) (\mathbf{W}_e^T \mathbf{W}_e)^{1/2} \right)$$

This matches the general formula from Theorem 1 for $N = 2$.

Deriving the Preconditioner

SVD Setup (N=3 case)

Let's analyze $N = 3$. Let the SVD of each layer be $\mathbf{W}_j = \mathbf{U}_j \Sigma_j \mathbf{V}_j^T$.

- Balance Invariants:** $\mathbf{W}_2^T \mathbf{W}_2 = \mathbf{W}_1 \mathbf{W}_1^T$. Hence, their spectra are equal as well $\implies \mathbf{W}_1$ and \mathbf{W}_2 have the same set of singular values.
- In terms of Polar decomposition, the balance equation implies $\mathbf{P}_2 = \mathbf{W}_2^T \mathbf{W}_2 = \mathbf{W}_1 \mathbf{W}_1^T = \mathbf{R}_1$. Hence, \mathbf{V}_2 and \mathbf{U}_1 can be chosen to be equal. Therefore, we get the SVD decompositions: $\mathbf{W}_2 = \mathbf{U}_2 \Sigma \mathbf{V}_2^T$ and $\mathbf{W}_1 = \mathbf{V}_2 \Sigma \mathbf{V}_1^T$.
- Similarly using the balance equation for \mathbf{W}_2 and \mathbf{W}_3 , we get the SVDs: $\mathbf{W}_3 = \mathbf{U}_3 \Sigma \mathbf{V}_3^T$, $\mathbf{W}_2 = \mathbf{V}_3 \Sigma \mathbf{V}_2^T$ and $\mathbf{W}_1 = \mathbf{V}_2 \Sigma \mathbf{V}_1^T$.

Simplifying the Product

Let's express the end-to-end matrix $\mathbf{W}_e = \mathbf{W}_3 \mathbf{W}_2 \mathbf{W}_1$ using the SVDs and the relationships we found. The expression simplifies because the intermediate orthogonal matrices cancel out:

$$\mathbf{W}_e = \mathbf{U}_3 \Sigma^3 \mathbf{V}_1^T$$

This looks like an SVD for \mathbf{W}_e , with singular value matrix Σ^3 and orthogonal matrices \mathbf{U}_3 and \mathbf{V}_1^T .

Final Identities

From $\mathbf{W}_e = \mathbf{U}_3 \Sigma^3 \mathbf{V}_1^T$, we can identify the SVD components of $\mathbf{W}_e = \mathbf{U}_e \Sigma_e \mathbf{V}_e^T$ as $\mathbf{U}_e = \mathbf{U}_3$, $\Sigma_e = \Sigma^3$, $\mathbf{V}_e = \mathbf{V}_1$.

Now we can express the individual layer terms using the end-to-end SVD components.

- For the last layer, $\mathbf{W}_3 \mathbf{W}_3^T = \mathbf{U}_3 \Sigma^2 \mathbf{U}_3^T$. Since $\mathbf{U}_3 = \mathbf{U}_e$ and $\Sigma = \Sigma_e^{1/3}$:
$$\mathbf{W}_3 \mathbf{W}_3^T = \mathbf{U}_e (\Sigma_e^{1/3})^2 \mathbf{U}_e^T = (\mathbf{U}_e \Sigma_e \mathbf{U}_e^T)^{2/3} = (\mathbf{W}_e \mathbf{W}_e^T)^{2/3}$$
- For the first layer, $\mathbf{W}_1^T \mathbf{W}_1 = \mathbf{V}_1 \Sigma^2 \mathbf{V}_1^T$. Using $\Sigma = \Sigma_e^{1/3}$ and $\mathbf{V}_1 = \mathbf{V}_e$:
$$\mathbf{W}_1^T \mathbf{W}_1 = \mathbf{V}_1 (\Sigma_e^{1/3})^2 \mathbf{V}_1^T = (\mathbf{V}_e^T \Sigma_e \mathbf{V}_e)^{2/3} = (\mathbf{W}_e^T \mathbf{W}_e)^{2/3}$$

How This Leads to Acceleration

Setup

We analyze a simple, ill-conditioned linear regression problem with ℓ_4 loss and $N = 2$ overparameterization.

- Data:** Two points, $([1, 0], y_1)$ and $([0, 1], y_2)$.
- Loss:** $L(w_1, w_2) = \frac{1}{4}(w_1 - y_1)^4 + \frac{1}{4}(w_2 - y_2)^4$.
- Ill-Conditioning:** Assume $|y_1| \gg |y_2| \approx 1$.
- Initialization:** Near-zero weights, $w_1^{(0)} = \epsilon_1$, $w_2^{(0)} = \epsilon_2$, with $\epsilon_2/\epsilon_1 \approx y_2/y_1$.

Standard GD

The standard gradient descent update for each coordinate is decoupled:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} - \eta (w_i^{(t)} - y_i)^3$$

- To avoid divergence, the learning rate η is limited by the coordinate with the largest gradient, which is w_1 . The stability condition requires: $\eta < \frac{2}{(w_1 - y_1)^3} \approx \frac{2}{y_1^3}$.
- This very small learning rate, dictated by y_1 , is then applied to the update for w_2 .
- The convergence for w_2 is extremely slow. The error $\Delta_2 = w_2 - y_2$ shrinks by a factor of approximately $(1 - \eta y_2^3)$ at each step, which is very close to 1.

Overparameterized GD

For $N = 2$, the discrete update rule for a single output network is:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \left(\|\mathbf{w}^{(t)}\| \nabla \mathcal{L}(\mathbf{w}^{(t)}) + \text{Pr}_{\mathbf{w}^{(t)}}(\nabla \mathcal{L}(\mathbf{w}^{(t)})) \right)$$

Using the condition $\epsilon_2/\epsilon_1 \approx y_2/y_1 \ll 1$, we have $\|\mathbf{w}^{(0)}\| \approx \epsilon_1$. The updates simplify to:

$$w_1^{(1)} \approx \epsilon_1 + \eta(2\epsilon_1 y_1^3) \quad , \quad w_2^{(1)} \approx \epsilon_2 + \eta(\epsilon_1 y_2^3 + \epsilon_2 y_1^3)$$

The Punchline

- Choose η :** We can now choose a large learning rate for w_1 to make it converge in one step. Let $\eta = \frac{1}{2\epsilon_1 y_1^3}$. This gives $w_1^{(1)} \approx y_1$.
- Analyze w_2 update:** With this η , the update for w_2 becomes approximately $\frac{y_2}{2}$.
- Compare Rates:** In one step, w_2 has moved halfway to its target y_2 . The effective learning rate for w_2 after w_1 converges is $\eta_{OP} \approx \frac{1}{2\epsilon_1 y_1^3}$. The speedup is:

$$\frac{\eta_{OP}}{\eta_{GD}} > \frac{1/(2\epsilon_1 y_1^3)}{2/y_1^3} = \frac{y_1}{4\epsilon_1} \gg 1$$

Overparameterization allows the large coordinate to "lend" its scale to accelerate the small coordinate.

Next Time...

Characterizing the Minimizer

Thank You!

Any Questions?

