

Abstract

Recent theory shows that training wide neural net-Context. works amounts to doing regression with a positive-definite kernel.

Contributions. This *lazy training* phenomenon: • is not intrinsically due to width but to a degenerate relative scale \rightarrow depends on early stopping, initialization and normalization • removes some benefits of depth and may **hinder generalization**

_azy Training

Setting. Adjust parameters of a differentiable model $h : \mathbb{R}^p \to \mathcal{F}$ by minimizing a loss $R: \mathcal{F} \to \mathbb{R}_+$ using gradient flow on the objective

$$F(w) = R(\alpha h(w)) / \alpha^2.$$

- \mathcal{F} is a Hilbert space of predictors, R typically the empirical or population risk, h typically a neural network
- $\alpha > 0$ is a scale, often implicitly present
- gradient flows approximate (stochastic, accelerated) gradient descent

Training paths. For initialization w_0 and stopping time T, let

- $(w_{\alpha}(t))_{t \in [0,T]}$ be the *original* optimization path
- $(\overline{\boldsymbol{w}}_{\alpha}(t))_{t\in[0,T]}$ be the *tangent* optimization path, for the tangent model



$\bar{h}(w) = h(w_0) + Dh(w_0)(w - w_0)$

Lazy Training (definition)

When the *original* and *tangent* optimization paths are close

Consequences. Lazy training is a type of implicit bias for gradient descent that leads to strong guarantees:

- on optimization speed (theory of convex optimization)
- on generalization (theory of kernel regression)

Lénaïc Chizat³, Édouard Oyallon⁴ and Francis Bach^{1,2} ¹INRIA, ²ENS-PSL Paris, ³CNRS, ⁴Centrale-Supélec Is it desirable in practice? — Synthetic experiments. Two-layer ReLU neural network, square loss, initialized with variance τ , best predictor has 3 neurons. aradient flow (-) Lazy Training (au = 0.1) Non-Lazy Training $(\tau = 2)$ Trajectory of each "hidden" neuron during training (2-D input) 2.5 -S 2.0 -............. est - 1.5 end of training Over-parameterized Under-parameterized (GD on train loss until 0 loss) (SGD on population loss) Impact of laziness on performance (100-D input) Image recognition. Does lazy training explain deep learning? 90 •••• train accuracy 80 test accuracy -- stability of activations α (scale of the model) Effect on laziness (VGG11 model) Linear vs. lazy vs. deep models **Theoretical arguments.** Neural networks can be superior to kernel/fixed features methods, thanks to their adaptivity (Bach 2017). Main references • Jacot et al., Neural Tangent Kernel: Convergence and Generalization in Neural Networks. 2018. • Du et al., Gradient Descent Provably Optimizes Over-parameterized Neural Networks. 2018.



On Lazy Training in Differentiable Programming Lazy Training Theorems — When does it occur? $\kappa_{h_m}(w_0) \lesssim m^{-1/2} + (\alpha m)^{-1}$

• over-parameterization is not needed

• see paper for precise statements

 $h(w_0)$ is controlled by

 \rightarrow lazy if $h(w_0)$ small and α large

Finite horizon $\lim_{\alpha \to \infty} \sup_{t \in [0,T]} \|\alpha h(\boldsymbol{w}_{\alpha}(t)) - \alpha \bar{h}(\boldsymbol{\bar{w}}_{\alpha}(t))\| = 0.$ **Infinite horizon** $\lim_{\alpha \to \infty} \sup_{t>0} \|\alpha h(\boldsymbol{w}_{\alpha}(t)) - \alpha \bar{h}(\bar{\boldsymbol{w}}_{\alpha}(t))\| = 0.$ **Case 1: Rescaled models Case 2: Homogeneous models** \rightarrow lazy if $h(w_0)$ small and λ large Case 3: Wide neural networks \rightarrow lazy if $\lim_{m\to\infty} \alpha m = \infty$ (e.g. $\alpha = 1/\sqrt{m}$)

PSL 😥 If $h(w_0) = 0$ and R potentially non-convex then for any T > 0, If $h(w_0) = 0$, and R is strongly convex, then **A sufficient criterion.** For the square loss $R(y) = \frac{1}{2} ||y - y^{\star}||^2$ and $\alpha = 1$, the relative difference $\Delta \coloneqq \|h(\mathbf{w}(t)) - \bar{h}(\bar{\mathbf{w}}(t))\|/\|y^{\star} - \bar{h}(\bar{\mathbf{w}}(t))\|/\|y^{\star} - \bar{h}(\bar{\mathbf{w}}(t))\|/\|y^{\star}\|$ $\Delta \lesssim \tilde{t}^2 \cdot \kappa_h(w_0) \quad \text{where} \quad \kappa_h(w_0) \coloneqq \frac{\|h(w_0) - y^\star\| \|D^2 h(w_0)\|}{\|Dh(w_0)\|^2}$ where $\tilde{t} = t \|Dh(w_0)\|^2$ is the normalized time (\approx iteration number). For $\alpha > 0$, one has $\kappa_{\alpha h}(w_0) \lesssim \|h(w_0) - y^*/\alpha\|$ If $h(\lambda w) = \lambda^q h(w)$, one has $\kappa_h(\lambda w_0) \lesssim \|h(w_0) - y^*/\lambda^q\|$ If $h_m(w) = \alpha \sum_{i=1}^m \phi(\theta_i)$ where $w = (\theta_1, \dots, \theta_m)$ are i.i.d. and satisfy $\mathbb{E}\phi(\theta_i) = 0$ (two-layer neural network), then

 \rightarrow can be extended to deep networks (Jacot et al.)

Model	Train acc.	Test acc.
ResNet wide, linearized	55.0	56.7
VGG-11 wide, linearized	61.0	61.7
Prior features (Oyallon et al.)	-	82.3
Random features (Recht et al.)	-	84.2
VGG-11 wide, standard	99.9	89.7
ResNet wide, standard	99.4	91.0

• Bach. Breaking the Curse of Dimensionality with Convex Neural Networks). 2017.