# TRUSTWORTHY MACHINE LEARNING THROUGH THE LENS OF HIGH-DIMENSIONAL PROBABILITY

Marco Mondelli, Institute of Science and Technology Austria (ISTA)

DARMSTADT SPRING SCHOOL IN PROBABILITY
"ARTIFICIAL INTELLIGENCE : PROBABILISTIC CONCEPTS"


## OVERVIEW

$\longrightarrow$ Role of overparameterization

*) How many parameters to interpolate the labels ?

*) How many parameters to interpolate robustly ?

*) How many parameters to reconstruct training data ?


$\longrightarrow$ Dynamics of differentially private algorithms

*) Privacy for free (when enough samples)

*) Optimal rate via aggressive clipping

# BASIC SETTING AND TERMINOLOGY

## Supervised learning

Training data : $\{ (x_i, y_i) \}_{i \le n} \overset{iid}{\sim} \mathbb{P}( \mathbb{R}^d \times \mathbb{R})$

response / label $\in \mathbb{R}$

vector of covariates in $\mathbb{R}^d$

unknown

Goal : Find $f : \mathbb{R}^d \longrightarrow \mathbb{R}$ to predict $y_{test}$ given $x_{test}$

with $(x_{test}, y_{test}) \sim \mathbb{P}$

test distribution = training distribution

(no distribution shift)

Measure __performance__ via test error / population error / population risk :

$$R(f) = \mathbb{E} \left[ \ell ( y_{test}, f( x_{test} )) \right]$$

loss function $\ell : \mathbb{R} \times \mathbb{R} \longrightarrow \mathbb{R}$

expectation over the distribution of $(x_{test}, y_{test})$

Running example : Square loss

$$R(f) = \mathbb{E} \left[ | y_{test} - f( x_{test} ) |^2 \right]$$

Parametric models (neural networks) : $f(x) = f(x ; \theta) \qquad \theta \in \mathbb{R}^p$

vector of parameters
(weights of the neural network)

Empirical risk minimization (ERM), standard approach to learn $f(\cdot\,; \theta)$

$$\hat{R}_n(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f(x_i; \theta))$$

empirical ↗

↙ # samples

$$\text{Solve} \quad \min_{\theta} \hat{R}_n(\theta)$$

We will solve the optimization problem via (several variants of) gradient descent:

*) Gradient descent (GD)   (Full batch)

$$\theta^{k+1} = \theta^k - \eta_k \nabla_\theta \hat{R}_n(\theta^k) = \theta^k - \eta_k \frac{1}{n} \sum_{i=1}^{n} \nabla_\theta \ell(y_i, f(x_i; \theta^k))$$

step size ↗ , can choose different step sizes for different groups of parameters

*) Gradient flow (GF)

$$\dot{\theta}_t = -\nabla_\theta \hat{R}_n(\theta_t)$$

When $\eta_k$ is small, GD $\approx$ GF

*) Stochastic gradient descent (SGD)   (Online, 1-pass)

$$\theta^{k+1} = \theta^k - \eta_k \nabla_\theta \ell(y_i, f(x_i; \theta^k)) \qquad i \sim \text{Unif}(\{1, \dots, n\})$$

Later in the mini-course also:

*) Differentially-private gradient descent / stochastic gradient descent

# HOW MANY PARAMETERS TO INTERPOLATE THE LABELS? ( MEMORIZATION #1 )

$$\exists \theta \quad \text{s.t.} \quad y_i = f(x_i, \theta) \quad \forall i \in \{1, \dots, n\}$$

Classical problem dating back to [Cover, 1965]: For a neural network with a single neuron and sign activation, label interpolation for data in generic position is possible if and only if $p/n > 1/2$.

*) Single neuron and sign activation: $f(x, \theta) = \text{sign}(\langle x, w \rangle)$

   ( spherical perceptron )

*) Data in generic position: Every subset of $x_1, \dots, x_n$ containing $p$ or fewer vectors is linearly independent ( satisfied by i.i.d. data with high probability).

More complex architectures? $\theta$ obtained by training via gradient descent?

<u>IDEA</u> : Under certain conditions, $f(x, \theta)$ is close to its linearization

$$f_{\text{lin}}(x, \theta) = f(x, \theta_0) + \langle \theta - \theta_0, \nabla_\theta f(x, \theta_0) \rangle$$

initialization

throughout the training (GD/GF) dynamics.

$f_{\text{lin}}(x, \theta)$ is <u>linear</u> in $\theta$ so its GD training dynamics $\theta_{\text{lin}}^k$ can be solved explicitly.

Some more notation :

*) $y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$

$$\ast) \quad f_{ein}(\theta) = \begin{bmatrix} f_{ein}(x_1, \theta) \\ \vdots \\ f_{ein}(x_n, \theta) \end{bmatrix} \in \mathbb{R}^n$$

We now compute:

$$\nabla_\theta f_{ein}(\theta) = \begin{bmatrix} \nabla_\theta f_{ein}(x_1, \theta) \\ \vdots \\ \nabla_\theta f_{ein}(x_1, \theta) \end{bmatrix} \overset{f_{ein}(x_i,\theta)\ \text{is linear in}\ \theta}{=} \begin{bmatrix} \nabla_\theta f(x_1, \theta_0) \\ \vdots \\ \nabla_\theta f(x_n, \theta_0) \end{bmatrix} = \Phi \in \mathbb{R}^{n \times p}$$

does not depend on $\theta$ !

Then, we have

$$\theta_{lin}^{k+1} = \theta_{lin}^k - \eta_k \nabla_\theta \hat{R}_n(\theta_{lin}^k)$$

square loss, constant step size $\eta_k \equiv \eta$

$$= \theta_{lin}^k - \eta \nabla_\theta \frac{1}{2} \| y - f_{ein}(\theta_{lin}^k) \|^2$$

$$= \theta^k - \eta \Phi^T (f_{ein}(\theta_{lin}^k) - y)$$

$$= \theta^k - \eta \Phi^T (\Phi \theta_{lin}^k - y)$$

pick $\theta_0 = 0$ and $f$ such that $f(x_i, \theta_0) = 0$ for all $x_i$.

Define $r^k = \Phi \theta_{lin}^k - y$. Then,

$$r^{k+1} = (I - \eta \Phi \Phi^T) r^k \implies \| r^{k+1} \|_2 \leq \| I - \eta \Phi \Phi^T \| \, \| r^k \|_2$$

$$\leq \left(1 - \frac{\lambda_{min}(\Phi \Phi^T)}{\|\Phi\|^2}\right) \| r^k \|_2$$

pick $\eta \leq \frac{1}{\|\Phi\|^2}$

Loss converges geometrically to $0$ as long as $\lambda_{min}(\widehat{\mathcal{I}}\widehat{\mathcal{I}}^T) > 0$.

Neural tangent kernel (NTK) $\in \mathbb{R}^{n \times n}$

NTK introduced by [Jacot, Gabriel, Hongler, 2018] and then lots of follow-up work [Du, Zhai, Poctor, Singh, 2019 ; Chizat, Oyallon, Bach, 2019;...]

# ONE FORMAL RESULT

Let $f(\theta) = \begin{bmatrix} f(x_1, \theta) \\ \vdots \\ f(x_n, \theta) \end{bmatrix} \in \mathbb{R}^n$, $J(\theta) \in \mathbb{R}^{n \times p}$ the corresponding Jacobian,

generic parametric model

$(J(\theta))_{ij} = \dfrac{\partial f(x_i, \theta)}{\partial \theta_j}$

and $\ell(y_i, f(x_i ; \theta)) = \frac{1}{2}(y_i - f(x_i, \theta))^2$ the square loss.

Consider a set $\mathcal{D} \subset \mathbb{R}^p$ containing the initialization $\theta_0$ such that

(A1) $\qquad \alpha \leq \sigma_{min}(J(\theta)) \leq \|J(\theta)\| \leq \beta \qquad \forall \; \theta \in \mathcal{D}$

Bounded spectrum of the Jacobian

(A2) $\qquad \|J(\theta_1) - J(\theta_2)\| \leq \dfrac{\alpha^2}{2\beta} \qquad \forall \; \theta_1, \theta_2 \in \mathcal{D}$

Bounded deviations of the Jacobian

THEOREM [Oymak, Soltanolkotabi, 2019] Assume that (A1) and (A2) hold when $\mathcal{D}$ is a ball centered at $\theta_0$ with radius $R = 4\|f(\theta_0) - y\|_2 / \alpha$. Run GD with constant step size $\eta \leq \dfrac{1}{2\beta^2}$ and initialization $\theta_0$. Then,

$$\|f(\theta^k) - y\|_2^2 \leq \left(1 - \dfrac{\eta \alpha^2}{2}\right)^k \|f(\theta_0) - y\|_2^2 .$$

INTERPRETATION Loss converges geometrically to $0$ as long as

$$\lambda_{min}(\mathcal{J} \mathcal{J}^T) = \alpha^2 > 0 .$$

Much more is true:

*) $\| \theta^{\kappa} - \theta^{\kappa}_{lin} \|_2 \ll \| \theta^{\kappa} - \theta_0 \|_2$ for all $\kappa$ until convergence

GD dynamics close to GD dynamics of the linearization

*) $\| f(\cdot, \theta^{\kappa}) - f_{lin}(\cdot, \theta^{\kappa}_{lin}) \|_{L^2} \ll \| f(\cdot, \theta^{\kappa}) \|_{L^2}$

with $\| g \|_{L^2} = (\mathbb{E}\, g^2(x))^{1/2}$ where $\mathbb{E}$ is taken on a test point.

Model learned by GD close to model learned by the linearized flow on test points

*) Total gradient path is bounded:

$$\sum_{\kappa=0}^{\infty} \| \theta^{\kappa+1} - \theta^{\kappa} \|_2 \leq \frac{4 \| f(\theta_0) - y \|_2}{\alpha}$$

Gradient descent iterates remain close to initialization (they never leave a neighborhood of radius $\frac{4}{\alpha} \| f(\theta_0) - y \|_2$ around initialization)

*) Gradient descent follows a short path:

$$\| \theta^{\kappa} - \theta_0 \| \leq 4 \frac{\beta}{\alpha} \| \theta^* - \theta_0 \|$$

$$\sum_{\kappa=0}^{\infty} \| \theta^{\kappa-1} - \theta^{\kappa} \| \leq 4 \frac{\beta}{\alpha} \| \theta^* - \theta_0 \| \quad ,$$

with $\theta^*$ a global optimum of the loss CLOSEST (in $\ell_2$) to initialization.

GD follows an almost direct route from initialization to a global optimum (the length of the GD path is within a factor of the distance between initialization and CLOSEST global optimum)

$\longrightarrow$ Also known as "lazy" training [Chizat, Oyallon, Bach, 2019]

$\longrightarrow$ This "linear regime" is discussed in detail in the review [Bartlett, Montanari, Rakhlin, 2021], see Section 5 therein.

# BOUNDING THE SMALLEST EIGENVALUE OF THE NTK

Key quantity in this analysis is $\alpha^2 = \lambda_{min}(\Phi\Phi^T)$

Here, $\Phi$ is the Jacobian at initialization which is a random quantity (weights at initialization are random, data is also random).

If $p < n$, then $\lambda_{min}(\Phi\Phi^T) = 0$ $\therefore$

How large does $p$ have to be so that $\lambda_{min}(\Phi\Phi^T) > 0$ ?

$\longrightarrow$ Two-layer networks

$$f(x,\vartheta) = a^T \phi(W\underset{\in \mathbb{R}^d}{x}) = \sum_{i=1}^{N} a_i \phi(<w_i, x>)$$

$\underset{\in \mathbb{R}^N}{a} \quad \underset{\text{activation function applied component-wise}}{\phi} \quad \underset{\in \mathbb{R}^{N\times d}}{W}$

$d$ = input dimension, $N$ = # neurons, $p = Nd$

$p \gg n \implies \lambda_{min}(\Phi\Phi^T) > 0$  [Montanari, Zhong, 2022]

(improvement upon earlier work by [Soltanolkotabi, Javanmard, Lee, 2018] )

This is optimal (up to poly-log factors) !

$\longrightarrow$ Deep networks

$$f(x,\vartheta) := f_L(x,\vartheta) \quad \text{with} \quad f_\ell(x,\vartheta) = \begin{cases} x & \ell = 0 \\ \phi(W_\ell^T f_{\ell-1}) & \ell \in \{1, \cdots, L-1\} \\ W_\ell^T f_{L-1} & \ell = L \end{cases}$$

$W_\ell \in \mathbb{R}^{N_{\ell-1} \times N_\ell}$, $N_0 = d$, $N_L = 1$, $N_\ell$ = # neurons at layer $\ell$

$(\ell \in \{1, \cdots, L-1\})$

A direct calculation gives

$$\Xi \, \Xi^T = \sum_{\ell=0}^{L-1} F_\ell \, F_\ell^T \circ B_{\ell+1} \, B_{\ell+1}^T$$

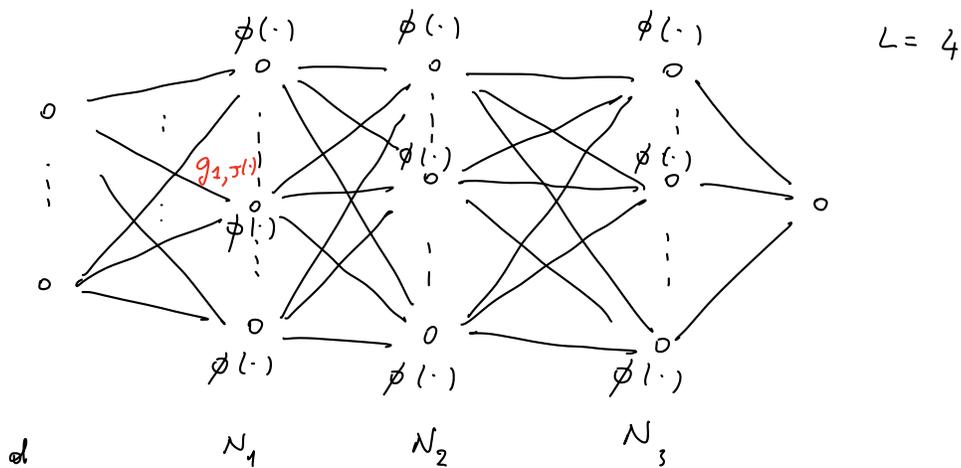<span style="color:blue">Hadamard (component-wise) product</span>

*) $\quad F_\ell = \begin{bmatrix} f_\ell(x_1) \\ \vdots \\ f_\ell(x_n) \end{bmatrix} \in \mathbb{R}^{n \times N_\ell} \quad$ feature matrix of layer $\ell$

*) $\quad B_\ell \in \mathbb{R}^{n \times N_\ell} \quad$ backpropagation matrix

$$(B_\ell)_{i:} = \begin{cases} \Sigma_\ell(x_i) \left( \prod_{\ell'=\ell+1}^{L-1} W_{\ell'} \, \Sigma_{\ell'}(x_i) \right) W_L \, , & \ell \in \{1, \cdots, L-2\} \\[2mm] \Sigma_{L-1}(x_i) \, W_L \, , & \ell = L-1 \\[2mm] 1 \, , & \ell = L \end{cases}$$

$$\Sigma_\ell(x) = \text{diag}\left( \left[ \phi'(g_{\ell, J}(x)) \right]_{J=1}^{N_\ell} \right)$$

<span style="color:blue">pre-activation neuron</span>

Weight assumption : $(W_\ell)_{ij} \overset{iid}{\sim} N(0, 1/n_{\ell-1})$ , $(W_L)_i \overset{iid}{\sim} N(0,1)$

standard in practice ( He, LeCun initialization )

Data assumptions :

(a) $\int \|x\|_2 \, dP_x(x) = \Theta(\sqrt{d})$

(b) $\int \|x\|_2^2 \, dP_x(x) = \Theta(d)$

(c) $\int \| x - \int x' \, dP_x(x') \|_2^2 \, dP_x(x) = \Omega(d)$ \} scaling

(d) $P\left( | \varphi(x) - \int \varphi(x') \, dP_x(x') | > t \right) \leq 2e^{- t^2/2c \cdot \text{Lip}(\varphi)^2}$, $\forall \varphi$ Lipschitz

dimension - independent constant

↓

(C-)Lipschitz concentration holds in a variety of settings :

*) Standard Gaussian

*) Uniform distribution on the sphere and the hypercube

*) Data produced by a fully connected neural network with well-conditioned weights and Lipschitz activations

*) Distributions satisfying log-Sobolev (with dimension-independent constant)

⌈ EXTRA :

A probability measure $\mu$ satisfies the log-Sobolev inequality with constant $C > 0$ if for any smooth function $f$

$$\text{Ent}_\mu(f^2) \leq C \int |\nabla f(x)|^2 \, d\mu(x) ,$$

with $\text{Ent}_\mu(f) = \int f \ln f \, d\mu - \int f \, \ln\left(\int f \, d\mu\right) d\mu$
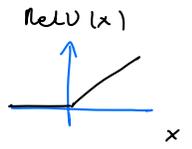
We now give a lower bound on $\lambda_{min}(\Phi\Phi^T)$ in two settings :

① Single wide layer $\qquad N_K \gg n \qquad$ for some $K \in \{1, \ldots, L-1\}$

② Minimum over parameterization $\qquad$ (smallest overall number of neurons)

$$N_{L-2} \; N_{L-1} \gg n$$

① Single wide layer $N_k \gg n$ for some $k \in \{1, \ldots, L-1\}$
+ no exponential bottleneck before the wide layer

$$\left( \prod_{\ell=1}^{k-2} \log N_\ell \ll \min_{\ell \in \{0, \ldots, k-1\}} N_\ell \right)$$

ReLU $(x)$

**THEOREM** (Simplification of [Nguyen, M., Montufar, 2021] ). Let $\phi(\cdot) = \mathrm{ReLU}(\cdot)$.

Assume $N_k = \tilde{\Omega}(n)$ and $N_\ell = \widehat{\Theta}(1)$ for $\ell \neq k$, where $\sim$ omits polylog

factors. Then, with high probability,

$$\lambda_{min}\left( \Phi \Phi^T \right) = \Theta(N_k).$$

**INTERPRETATION** One layer whose width is linear in the number of samples

suffices to guarantee well-behavedness of the NTK.


**PROOF SKETCH**.

**STEP 1**: From NTK to feature matrices.

For PSD matrices $P, Q$, it holds

$$\lambda_{min}(P \circ Q) \geq \lambda_{min}(P) \min_i Q_{ii}$$

$$\Downarrow$$

$$\lambda_{min}\left( \Phi \Phi^T \right) \geq \sum_{\ell=0}^{L-1} \lambda_{min}\left( F_\ell F_\ell^T \right) \min_{i \in \{1, \ldots, n\}} \| (B_{\ell+1})_{i:} \|_2^2$$

$$\geq \lambda_{min}\left( F_k F_k^T \right) \min_{i \in \{1, \ldots, n\}} \| (B_{k+1})_{i:} \|_2^2$$

$$\| (B_{k+1})_{i:} \|_2^2 = \| \underbrace{\Sigma_{k+1}(x_i) \left( \prod_{\ell=k+2}^{L-1} W_\ell \Sigma_\ell (x_i) \right)}_{B} W_L \|_2^2$$

$$\underset{W_L}{\overset{\approx}{}} \mathbb{E} \| B W_L \|_2^2 = \| B \|_F^2 = \Theta(1)$$

Hansan-Wright $(A = B^2)$          induction + Bernstein

## STEP 2: Concentration of smallest eigenvalue

$$\lambda_{min}\left( F_k \, F_k^T \right) \geqslant \frac{1}{4} \, \lambda_{min}\left( \underset{W_k}{\mathbb{E}} \, F_k \, F_k^T \right)$$

$$= \frac{N_k}{4} \, \lambda_{min}\left( \underset{w \sim N(0, \frac{1}{N_{k-1}} I)}{\mathbb{E}} \, \phi\left( F_{k-1} \, w \right) \, \phi\left( F_{k-1} \, w \right)^T \right)$$

$\phi(\cdot)$ is homogeneous

$$= \frac{N_k}{4 N_{k-1}} \, \lambda_{min}\left( D \, \underset{\hat{w} \sim N(0, I)}{\mathbb{E}}\left[ \phi(\hat{F}_{k-1} \, w) \, \phi(\hat{F}_{k-1} \, w)^T \right] D \right)$$

$\phi_r = r$-th Hermite coefficient of $\phi$     $D = \text{diag}\left( \|(F_k)_{1:}\|_2, \cdots, \|(F_k)_{n:}\|_2 \right)$

$$= \frac{N_k}{4 N_{k-1}} \, \lambda_{min}\left( D \left( \phi_0^2 \, 1_n \, 1_n^T + \sum_{s=1}^{\infty} \phi_s^2 \left( \hat{F}_{k-1} \, \hat{F}_{k-1}^T \right)^{\circ s} \right) D \right)$$

$$\geqslant \frac{N_k}{4 N_{k-1}} \, \phi_r^2 \, \lambda_{min}\left( D \left( \hat{F}_{k-1} \, \hat{F}_{k-1}^T \right)^{\circ r} D \right)$$

$$= \frac{N_k}{4 N_{k-1}} \, \phi_r^2 \, \lambda_{min}\left( D^{-(r-1)} \left( F_{k-1} \, F_{k-1}^T \right)^{\circ r} D^{-(r-1)} \right)$$

$$\geqslant \underbrace{\frac{N_k}{4 N_{k-1}} \, \phi_r^2}_{\substack{/ \\ \text{polylog}(n)}} \frac{\lambda_{min}\left( \left( F_{k-1} \, F_{k-1}^T \right)^{\circ r} \right)}{\underbrace{\max_{i \in \{1, \cdots, n\}} \|(F_{k-1})_{i:}\|_2^{2(r-1)}}_{\text{polylog}(n)}}$$

$$\overset{\sim}{=} \tilde{\Theta}\left( N_k \, \lambda_{min}\left( \left( F_{k-1} \, F_{k-1}^T \right)^{\circ r} \right) \right)$$

Let $\tilde{F}_{k-1} = F_{k-1} - \underset{(x_1,\cdots,x_n)}{\mathbb{E}} F_{k-1}$, $\mu = \underset{x}{\mathbb{E}} f_{k-1}(x) \in \mathbb{R}^{N_{k-1}}$ and $\Lambda = \text{diag}\left(F_{k-1}\,\mu - \|\mu\|^2 1_n\right)$.

Then,

$$F_{k-1} F_{k-1}^T = \tilde{F}_{k-1} \tilde{F}_{k-1}^T + \|\mu\|^2 1_n 1_n + \Lambda 1_n 1_n^T + 1_n 1_n^T \Lambda$$

$$= \tilde{F}_{k-1} \tilde{F}_{k-1}^T + \underbrace{\left(\|\mu\| 1_n + \frac{\Lambda 1_n}{\|\mu\|}\right)\left(\|\mu\| 1_n + \frac{\Lambda 1_n}{\|\mu\|}\right)^T}_{\succeq 0} - \frac{\Lambda 1_n 1_n^T \Lambda}{\|\mu\|^2}$$

$$\succeq \tilde{F}_{k-1} \tilde{F}_{k-1}^T - \frac{\Lambda 1_n 1_n^T \Lambda}{\|\mu\|^2}$$

$$\lambda_{min}\left(\left(F_{k-1} F_{k-1}^T\right)^{\circ r}\right) \geq \lambda_{min}\left(\left(\tilde{F}_{k-1} \tilde{F}_{k-1}^T - \frac{\Lambda 1_n 1_n^T \Lambda}{\|\mu\|^2}\right)^{\circ r}\right)$$

$$= \lambda_{min}\left(\underbrace{\left(\tilde{F}_{k-1} \tilde{F}_{k-1}^T\right)^{\circ r}}_{\approx I}\right) \left(1 + o(1)\right) = \Theta(1).$$

by picking large enough $r$

$$\left(\left(\tilde{F}_{k-1} \tilde{F}_{k-1}^T\right)^{\circ r}\right)_{ij} = \langle \underbrace{\tilde{f}_{k-1}(x_i)}, \tilde{f}_{k-1}(x_j)\rangle^r \approx \delta_{ij}$$

$$\overset{''}{f}_{k-1}(x_i) - \underset{x_i}{\mathbb{E}} f_{k-1}(x_i) \quad \text{centered feature vector at layer } k-1$$
computed with input $x_i$

$f_{k-1}(\cdot)$ is Lipschitz, so by Assumption (d) on the data,

$$\langle \tilde{f}_{k-1}(x_i), \tilde{f}_{k-1}(x_j)\rangle \ll \|\tilde{f}_{k-1}(x_i)\| \, \|\tilde{f}_{k-1}(x_j)\| \quad \forall i \neq j$$

This gives the lower bound $\lambda_{min}(\mathbb{F}\,\mathbb{F}^r) = \Omega(N_k)$.

The upper bound is easy:

$$\lambda_{min}(\underline{\Phi}\,\underline{\Phi}^T) \leq (\underline{\Phi}\,\underline{\Phi}^T)_{1,1} = \sum_{\ell=0}^{L-1} \|(F_\ell)_{1:}\|^2 \overbrace{\|(B_{\ell+1})_{1:}\|^2}^{\Theta(1)}$$

$$= \Theta(N_K).$$

🔁

⌐ EXTRA:

Hanson-Wright inequality.

THEOREM    Let $x \in \mathbb{R}^n$ be a random vector with independent, mean zero, sub-Gaussian coordinates (with sub-Gaussian norm of constant order).

Let $A \in \mathbb{R}^{n \times n}$. Then,

$$\mathbb{P}\left(\,|\,x^T A x - \mathbb{E}\,x^T A x\,| > t\,\right) \leq 2\, e^{-c\,min\left(\frac{t^2}{\|A\|_F^2}\,,\,\frac{t}{\|A\|}\right)}$$

Matrix Chernoff.

THEOREM. Let $A \in \mathbb{R}^{n \times m}$ and assume $\|(A_{:J})(A_{:J})^T\| \leq t^2$. Then,

$$\mathbb{P}\left(\,\lambda_{min}(A A^T) \leq (1-\epsilon)\,\lambda_{min}(\mathbb{E}\,A A^T)\,\right) \leq n\left[\frac{e^{-\epsilon}}{(1-\epsilon)^{1-\epsilon}}\right]^{\frac{\lambda_{min}(\mathbb{E}\,A A^T)}{t^2}}$$

Hermite expansion.

THEOREM    Let $\phi : \mathbb{R} \to \mathbb{R}$ with Hermite coefficients $\{\phi_r\}_{r\in\mathbb{N}}$, let $u, v \in \mathbb{R}^d$ such that $\|u\|_2 = \|v\|_2 = 1$. Then,

$$\mathbb{E}_{w \sim N(0, Id)}\, \phi(\langle w, u\rangle)\,\phi(\langle w, v\rangle) = \sum_{r\in\mathbb{N}} \phi_r^2\,\langle u, v\rangle^r.$$
⌐

② Minimum over parameterization (smallest number of neurons)    $N_{L-2} N_{L-1} \gg n$

THEOREM [Bombari, Amani, M., 2022]    Let $\phi$ be non-linear, Lipschitz,
with Lipschitz derivative. Assume $N_\ell = O(N_{\ell-1})$ for $\ell \in \{1, \cdots, L-1\}$ and
$n \, \text{poly} \log(n) = o(N_{L-2} N_{L-1})$. Then, with high probability,

$$\lambda_{min}(\Phi \Phi^T) = \Omega(N_{L-2} N_{L-1}) \qquad \text{and} \qquad \lambda_{min}(\Phi \Phi^T) = O(d N_{L-1})$$

INTERPRETATION    $\tilde{\Omega}(\sqrt{n})$ neurons and thus $\tilde{\Omega}(n)$ parameters are enough
to interpolate $n$ data points. This is in agreement with back-of-the-envelope
calculations on CIFAR-10 and Image Net:

*) CIFAR-10 :    $n = 50000$ images and $10^6$ parameters suffice to fit
    random labels

*) Image Net :    $n = 1.2 \cdot 10^6$ images and $2.4 \, 10^7$ parameters suffice to fit
    random labels

PROOF SKETCH.    Upper bound as before. Let's focus on the lower bound.

If $N_\ell < n$, then $\lambda_{min}(F_\ell F_\ell^T) = 0$ and step 1 of the previous
approach fails.

$$\Phi \Phi^T = \sum_{\ell=0}^{L-1} F_\ell F_\ell^T \circ B_{\ell+1} B_{\ell+1}^T \succeq F_{L-2} F_{L-2}^T \circ B_{L-1} B_{L-1}^T := \Phi_{L-2} \Phi_{L-2}^T$$

<span style="color:blue">( need to consider the Hadamard product together</span>

Let $A, B \in \mathbb{R}^{n \times \sqrt{n}}$. Then, $A A^T, B B^T \in \mathbb{R}^{n \times n}$ have rank $\sqrt{n}$,

but $A A^T \circ B B^T \in \mathbb{R}^{n \times n}$ can have rank $n$.

$$\left( \overline{\Phi}_{L-2} \right)_{i:} = f_{L-2} (x_i) \otimes \text{diag} (W_L) \sigma' \left( W_{L-1}^T f_{L-2} (x_i) \right)$$

Kronecker (tensor) product

rows are iid w.r.t. the randomness in $(x_1, \dots, x_n)$.

$\underline{\text{THEOREM}}$ (Simplification of $[$Adamczak, Litvak, Pajor, Tomczak-Jaegermann, 2011$]$).

Let $\Phi$ be a matrix with i.i.d. rows having $\ell_2$ norm $\ell$ and

sub-exponential norm $\psi$. Then, with high probability,

$$\lambda_{\min} \left( \Phi \Phi^T \right) \geq \ell^2 - \tilde{O} \left( \psi^2 \ell \sqrt{n} \right)$$

STEP 1 : Centering

*) Necessary otherwise the sub-exponential norm is too large (dominated by

the mean).

*) Need to center $F_{L-2}$ and $B_{L-1}$ together, otherwise cross-terms

become too large.

$$\lambda_{\min} \left( F_{L-2} F_{L-2}^T \circ B_{L-1} B_{L-1}^T \right) \geq \lambda_{\min} \left( \overbrace{\tilde{F}_{L-2} \tilde{F}_{L-2}^T}^{\tilde{\Phi}_{L-2}^{(1)} \; \tilde{\Phi}_{L-2}^{(1)T}} \circ \tilde{B}_{L-1} \tilde{B}_{L-1}^T \right) - o \left( N_{L-2} N_{L-1} \right)$$

$$\tilde{F}_{L-2} = F_{L-2} - \underset{x_1, \dots, x_n}{\mathbb{E}} F_{L-2}, \quad \tilde{B}_{L-1} = B_{L-1} - \underset{x_1, \dots, x_n}{\mathbb{E}} B_{L-1}$$

$$\left( \widetilde{\Phi}_{L-2}^{(1)} \right)_{i:} = \widetilde{f}_{L-2}(x_i) \otimes \operatorname{diag}(W_L) \, \widetilde{\sigma}'\left( W_{L-1}^T f_{L-2}(x_i) \right)$$

$$\underbrace{\phantom{f_{L-2}(x_i)}}_{} \qquad \underbrace{\phantom{\widetilde{\sigma}'}}_{}$$

$$\parallel \qquad\qquad\qquad\qquad \parallel$$

$$f_{L-2}(x_i) - \mathbb{E}_{x_i} f_{L-2}(x_i) \qquad \sigma'\left( W_{L-1}^T f_{L-2}(x_i) \right) - \mathbb{E}_{x_i} \sigma'\left( W_{L-1}^T f_{L-2}(x_i) \right)$$

this is still not 0-mean w.r.t. $x_i$ because of the quadratic terms

$$x \otimes x = \left[ x_1^2 \quad x_1 x_2 \quad x_1 x_3 \quad \cdots \quad x_1 x_n \quad x_2 x_1 \quad x_2^2 \quad \cdots \quad x_n^2 \right]$$

$$\uparrow \qquad\qquad\qquad\qquad\qquad\qquad \uparrow \qquad\qquad \uparrow$$

$$\text{not} \quad \text{0-mean}$$

$$\left( \widetilde{\Phi}_{L-2} \right)_{i:} = \left( \widetilde{\Phi}_{L-2}^{(1)} \right)_{i:} - \mathbb{E}_{x_i} \left( \widetilde{\Phi}_{L-2}^{(1)} \right)_{i:}$$

$$\widetilde{F}_{L-2} \widetilde{F}_{L-2}^T \circ \widetilde{B}_{L-1} \widetilde{B}_{L-1}^T = \widetilde{\Phi}_{L-2} \widetilde{\Phi}_{L-2}^T + \underbrace{\text{rank-}1}_{} + \text{PSD}$$

operator norm bounded via (generalized) Hanson-Wright

$$\Downarrow$$

$$\lambda_{\min}\left( \widetilde{F}_{L-2} \widetilde{F}_{L-2}^T \circ \widetilde{B}_{L-1} \widetilde{B}_{L-1}^T \right) \geqslant \lambda_{\min}\left( \widetilde{\Phi}_{L-2} \widetilde{\Phi}_{L-2}^T \right) + o\left( N_{L-1} N_{L-2} \right)$$

STEP 2 : Control norms

$$\ell = \left\| \left( \widetilde{\Phi}_{L-2} \right)_{i:} \right\| = \Theta\left( \sqrt{N_{L-1} N_{L-2}} \right)$$

$$\psi = \left\| \left( \widetilde{\Phi}_{L-2} \right)_{i:} \right\|_{\psi_1} = \widetilde{O}(1)$$

$$\underbrace{\phantom{\psi_1}}_{\text{sub-exponential norm}}$$

Applying the result by Adamczak et al. gives

$$\lambda_{\min}\left( \widetilde{\Phi}_{L-2} \widetilde{\Phi}_{L-2}^T \right) \geqslant N_{L-1} N_{L-2} - \widetilde{O}\left( \sqrt{n N_{L-1} N_{L-2}} \right) = \Omega\left( N_{L-1} N_{L-2} \right)$$

# How Many Parameters To Interpolate Robustly?

Machine learning models are vulnerable to adversarial perturbations, and this has been known for over 10 years now [Szegedy, Baremba, Sutskever, Bruna, Erhan, Goodfellow, Fergus, 2014].

*) Prototypical examples: given the image of e.g. a cat, one can construct a perturbation which is imperceptible to the human eye but makes the neural network classify the image as e.g. a gibbon.

Immense literature on the topic, including "adversarial training" methods aimed at reducing this effect [Madry, Makelov, Schmidt, Tsipras, Vladu, 2018].

Our focus is on the ROLE of OVERPARAMETERIZATION :

*) How many parameters NECESSARY to have robustness?

*) How many parameters SUFFICIENT to have robustness?

# A NECESSARY CONDITION: UNIVERSAL LAW OF ROBUSTNESS

**THEOREM** ( [Bubeck, Sellke, 2021] ). Let $\mathcal{F}$ be a class of functions from $\mathbb{R}^d \to \mathbb{R}$

Let $\{(x_i, y_i)\}_{i=1}^n$ be input-output pairs in $\mathbb{R}^d \times [-1,1]$ and fix $\epsilon, \delta \in (0,1)$.

Assume that

1. $\mathcal{F} = \{ f_w, w \in W, Lip(f_w) \leq L \}$ with $W \subseteq \mathbb{R}^p$, $diam(W) \leq W$

   and for any $w_1, w_2 \in W$

   $$\| f_{w_1} - f_{w_2} \| \leq J \| w_1 - w_2 \|$$

2. The distribution $\mu$ of the covariates $x_i$ can be written as $\mu = \sum_{\ell=1}^{k} \alpha_\ell \mu_\ell$,

   where $\alpha_\ell \geq 0$, $\sum_{\ell=1}^{n} \alpha_\ell = 1$, $k \log(\beta k/\delta) \leq c \cdot n \epsilon^2$ and each $\mu_\ell$ satisfies

   Lipschitz concentration:

   $$\mathbb{P}\left[ | f(x) - \mathbb{E}\{f\} | \geq t \right] \leq 2 e^{-\frac{c d t^2}{Lip(f)^2}}$$
   $$x \sim \mu_\ell$$

3. The expected conditional variance of the output is strictly positive, i.e.

   $$\sigma^2 \equiv \mathbb{E}_\mu Var[y_i | x_i] > 0$$

4. The dimension $d$ is large compared to $\epsilon$, i.e., $d \geq C_1 \left( \frac{c L^2 \sigma^2}{\epsilon^2} \right)$.

Then, with probability at least $1-\delta$ with respect to the sampling of the data,

one has simultaneously for all $f \in \mathcal{F}$:

$$\frac{1}{n} \sum_{i=1}^{n} (f(x_i) - y_i)^2 \leq \sigma^2 - \epsilon$$

$$\Downarrow$$

$$Lip(f) \geq \frac{\epsilon}{\sigma} \sqrt{c} \sqrt{\frac{nd}{p \log(1 + W J \epsilon^{-1} c^{-1}) + \log(4/\delta)}}.$$

## INTERPRETATION.

$$x_{adv} = x + \Delta_{adv},$$

adversarial perturbation which is "imperceptible"

$$\| \Delta_{adv} \| \leq \delta \cdot \|x\|$$

adversarial example close to $x$

$$| f(x_{adv}, \vartheta) - f(x, \vartheta) | \simeq | \nabla_x f(x, \vartheta)^T \Delta_{adv} |$$

$$\leq \| \Delta_{adv} \| \, \| \nabla_x f(x, \vartheta) \|$$

$$\leq \delta \cdot \|x\| \cdot \| \nabla_x f(x, \vartheta) \|$$

$$\leq c \, \delta \, \text{Lip}(f)$$

*) under Assumption 2, $\|x\| = \Theta(1)$

*) $\sup_x \| \nabla_x f(x, \vartheta) \| = \text{Lip}(f)$

Let us take $\Delta_{adv}$ aligned with the direction of the gradient (w.r.t. $x$) having the largest norm. Then, if $\text{Lip}(f) \gg 1$, $f(x_{adv})$ is far away from $f(x)$.


## RESTATING THE RESULT.

Consider a family of functions $\mathcal{F}$ that admits a Lipschitz parameterization by $p$ parameters each of size poly$(n, d)$ (Assumption 1 with $J = \text{poly}(n, d)$). Assume that the distribution of the data features is a mixture of (not too many) Lipschitz-concentrated distributions (Assumption 2). Then, if the network fits below the noise level $\sigma^2 > 0$ (e.g., it interpolates the training data), then, with high probability, $p = \Omega(nd)$ is NECESSARY to have a robust predictor, namely, such that $f(x_{adv})$ has to be close to $f(x)$ for any choice of $\Delta_{adv}$.

# IDEA OF THE PROOF.

Consider an $f$ that interpolates random $\pm 1$ labels

*) Lipschitz concentration implies that either $0$-level set or $1$-level set of $f$

have probability $\leq e^{-\frac{d}{\text{Lip}(f)^2}}$

*) Probability of fitting all $n$ points is $\leq e^{-\frac{nd}{\text{Lip}(f)^2}}$

*) Union bound over a function class of size $N$ gives probability of

$$N e^{-\frac{nd}{\text{Lip}(f)^2}} = e^{\log N - \frac{nd}{\text{Lip}(f)^2}}$$

*) Discretization argument gives that $\log N = O(p)$ for smoothly parametrized

family of functions with $p$ parameters

Need $\qquad p - \frac{nd}{\text{Lip}(f)^2} = \Omega(1)$ (otherwise probability of finding a fitting function

is very small )

$$\Downarrow$$

$$\text{Lip}(f) = \Omega\left(\sqrt{\frac{nd}{p}}\right).$$

PROOF    For simplicity pick $k = 1$ and assume $\sigma, \epsilon$ are constants $> 0$ (independent of $n, d, p$) so that we don't need to track their dependency in the calculations. We use $c$ to denote a constant $> 0$ independent of $n, d, p$ whose value may change from passage to passage.

STEP 1 : If $f$ fits below noise level, then the predictions are correlated with noise.

Let $\underbrace{g(x) = \mathbb{E}[y|x]}$ and $\underbrace{z_i = y_i - g(x_i)}$. Then,

$\quad\quad\quad\quad$ target function $\quad\quad\quad\quad$ noise part of the observed label $y_i$

$$\mathbb{P}\left( \exists f \in \mathcal{F} : \frac{1}{n}\sum_{i=1}^{n}(y_i - f(x_i))^2 \le \sigma^2 - \epsilon \right)$$

$$\le 2e^{-cn\epsilon^2} + \mathbb{P}\left( \exists f \in \mathcal{F} : \frac{1}{n}\sum_{i=1}^{n} f(x_i)\, z_i \ge \epsilon/4 \right)$$

(1)

We now prove (1).

*) $(z_i^2)$ are iid. with mean $\sigma^2$ $(\sigma^2 = Var[y|x])$ and $|z_i|^2 \le 4$ $(y_i \in [-1,1])$. Thus, by Hoeffding's inequality,

$$\mathbb{P}\left( \frac{1}{n}\sum_{i=1}^{n} z_i^2 \le \sigma^2 - \frac{\epsilon}{6} \right) \le e^{-cn\epsilon^2}$$

*) $(z_i\, g(x_i))$ are iid with mean $0$ $(\mathbb{E}[z_i | x_i] = 0)$ and $|z_i\, g(x_i)| \le 2$ $(y_i \in [-1,1])$. Thus, by Hoeffding's inequality,

$$\mathbb{P}\left( \frac{1}{n}\sum_{i=1}^{n} z_i\, g(x_i) \le -\frac{\epsilon}{6} \right) \le e^{-cn\epsilon^2}$$

*) Define $z = \frac{1}{\sqrt{n}}(z_1, \cdots, z_n)$, $G = \frac{1}{\sqrt{n}}(g(x_1), \cdots, g(x_n))$,

$F = \frac{1}{\sqrt{n}}(f(x_1), \cdots, f(x_n))$.

We can rewrite (1) as

$$\mathbb{P}\left(\exists f \in \mathcal{F} : \|G + z - F\|^2 \leq \sigma^2 - \epsilon\right) \leq 2e^{-cn\epsilon^2} + \mathbb{P}\left(\exists f \in \mathcal{F} : \langle F, z \rangle \geq \frac{\epsilon}{4}\right)$$

Note that we have just proved that, with probability at least $1 - 2e^{-cn\epsilon^2}$,

$$\|z\|^2 \geq \sigma^2 - \frac{\epsilon}{6} \qquad \text{and} \qquad \langle z, G \rangle \geq -\frac{\epsilon}{6} .$$

Thus,

$$\sigma^2 - \epsilon \geq \|G + z - F\|^2 = \|z\|^2 + 2\langle z, G - F\rangle + \|G - F\|^2$$

$$= \underbrace{\|z\|^2}_{\geq \sigma^2 - \epsilon/6} + \underbrace{2\langle z, G\rangle}_{\geq -\epsilon/6} - 2\langle z, F\rangle + \underbrace{\|G - F\|^2}_{\geq 0} \geq \sigma^2 - \frac{\epsilon}{2} - 2\langle z, F\rangle$$

$$\Downarrow$$

$$\langle F, z \rangle \geq \epsilon/4 \qquad, \text{ which gives (1).}$$

STEP 2: Chance that $f$ fits below noise level is $e^{\log|\mathcal{F}| - c\underbrace{\frac{nd}{L^2}}_{}}$

All functions in $\mathcal{F}$ are $L$-Lipschitz

$$\mathbb{P}\left(\exists f \in \mathcal{F} : \frac{1}{n}\sum_{i=1}^{n}(y_i - f(x_i))^2 \leq \sigma^2 - \epsilon\right) \leq 4e^{-cn\epsilon^2} + 2e^{\log|\mathcal{F}| - c\frac{\epsilon^2 nd}{L^2}}$$

(2)

We now prove (2).

*) $f$ $L$-Lipschitz + distribution of $x_i$ satisfying $\mathbb{P}(|f(x_i) - \mathbb{E}[f]| > t) \leq 2e^{-\frac{cdt^2}{L^2}}$

$$\Downarrow$$

$$\frac{f(x_i) - \mathbb{E}[f]}{L}\sqrt{d} \qquad \text{is} \qquad c\text{-subGaussian}$$

$$\Downarrow$$

$$\frac{f(x_i) - \mathbb{E}[f]}{L}\sqrt{d} z_i \qquad \text{is} \qquad c\text{-subGaussian} \qquad (|z_i| \leq 2)$$

+ it is 0 mean $\quad$ ( $\mathbb{E}[(f(x_i) - \mathbb{E}[f])z_i] = 0 \quad$ as $\quad \mathbb{E}[z_i \mid x_i] = 0$ )

$$\Downarrow$$

$$\frac{1}{\sqrt{n}} \sum_{i=1}^{n} \frac{f(x_i) - \mathbb{E}[f]}{L} \sqrt{d}\, z_i \quad \text{is} \quad C-\text{sub Gaussian}$$

$$\Downarrow$$

$$\mathbb{P}\left( \sqrt{\frac{d}{nL^2}} \sum_{i=1}^{n} (f(x_i) - \mathbb{E}[f])\, z_i \geq t \right) \leq 2 e^{-ct^2}$$

$$\Downarrow$$

$$\mathbb{P}\left( \frac{1}{n} \sum_{i=1}^{n} (f(x_i) - \mathbb{E}[f])\, z_i \geq \frac{\epsilon}{8} \right) \leq 2 e^{-c \frac{\epsilon^2 n d}{L^2}}$$

*1 $\quad \mathbb{P}\left( \exists f \in \mathcal{F} : \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}[f] z_i \geq \frac{\epsilon}{8} \right) \leq \mathbb{P}\left( \left| \frac{1}{n} \sum_{i=1}^{n} z_i \right| \geq \frac{\epsilon}{8} \right) \quad$ ( $\mathbb{E}[f] \in [-1,1]$ )

$$\leq 2 e^{-cn\epsilon^2} \quad \text{( Hoeffding } \quad \text{since } |z_i| \leq 2 \text{)}$$

Putting everything together, we have

$$\mathbb{P}\left( \exists f \in \mathcal{F} : \frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2 \leq \sigma^2 - \epsilon \right)$$

$$\underset{\underbrace{\quad}_{(1)}}{\leq} 2 e^{-cn\epsilon^2} + \mathbb{P}\left( \exists f \in \mathcal{F} : \frac{1}{n} \sum_{i=1}^{n} f(x_i) z_i \geq \epsilon/4 \right)$$

$$\leq 2 e^{-cn\epsilon^2} + |\mathcal{F}| \underbrace{\mathbb{P}\left( \frac{1}{n} \sum_{i=1}^{n} (f(x_i) - \mathbb{E}[f]) z_i \geq \frac{\epsilon}{8} \right)}_{\text{union bound}} + \mathbb{P}\left( \left| \frac{1}{n} \sum_{i=1}^{n} z_i \right| \geq \frac{\epsilon}{8} \right)$$

$$\underbrace{\leq}_{\text{calculations above}} 2|\mathcal{F}|\, e^{-c \frac{\epsilon^2 n d}{L^2}} + 4 e^{-cn\epsilon^2} \quad , \quad \text{which gives} \quad (2).$$

**STEP 3 : Covering argument (ε-net)**

Let $W_\epsilon$ be are $\frac{\epsilon}{16 J}$ - net of $W$. Then,

$$|W_\epsilon| \leq \left( 1 + \frac{JW}{c\epsilon} \right)^p .$$

Apply (2) to $\mathcal{F}_\epsilon = \{ f_w, w \in W_\epsilon \}$ :

$$\mathbb{P}\left( \exists f \in \mathcal{F}_\epsilon : \frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2 \leq \sigma^2 - \epsilon/2 \right) \leq 4\, e^{-cn\epsilon^2} + 2\, e^{p \log(1 + JW\epsilon^{-1}c^{-1}) - c\frac{\epsilon^2 nd}{L^2}}$$

If $\| f - h \|_\infty \leq \frac{\epsilon}{16} \leq 1$, $\| f \|_\infty \leq 1$, $\| h \|_\infty \leq 1$ and $|y_i| \leq 1$ for all $i$, then

$$\frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2 \leq \frac{1}{n} \sum (y_i - h(x_i))^2 + \frac{\epsilon}{2}$$

$$(y_i - f(x_i))^2 = (y_i - h(x_i) + (h(x_i) - f(x_i)))^2$$

$$= (y_i - h(x_i))^2 + (h(x_i) - f(x_i))^2 + 2(y_i - h(x_i))(h(x_i) - f(x_i))$$

$$\leq (y_i - h(x_i))^2 + \left(\frac{\epsilon}{16}\right)^2 + 2 \cdot 2 \frac{\epsilon}{16} \leq (y_i - h(x_i))^2 + \frac{\epsilon}{2}$$

For any $w \in W$, $\exists w' \in W_\epsilon$ s.t. $\| w - w' \| \leq \frac{\epsilon}{16 J}$

$$\Downarrow$$

For any $f_w \in \mathcal{F}$, $\exists f_{w'} \in \mathcal{F}_\epsilon$ s.t. $\| f_w - f_{w'} \|_\infty \leq \frac{\epsilon}{16}$

$$\Downarrow$$

$$\frac{1}{n} \sum_{i=1}^{n} (y_i - f_w(x_i))^2 \leq \frac{1}{n} \sum (y_i - f_{w'}(x_i))^2 + \frac{\epsilon}{2}$$

$$\Downarrow$$

$$\mathbb{P}\left( \exists f \in \mathcal{F} : \frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2 \leq \sigma^2 - \epsilon \right) \leq 4\, e^{-cn\epsilon^2} + 2\, e^{p \log(1 + JW\epsilon^{-1}c^{-1}) - c\frac{\epsilon^2 nd}{L^2}}$$

$$4 e^{-cn\epsilon^2} \leq \frac{\delta}{2} \qquad (\log(8/\delta) \leq c \cdot n\epsilon^2)$$

$$2 e^{p\log(1+JW\epsilon^{-1}c^{-1}) - c\frac{\epsilon^2 nd}{L^2}} = \frac{\delta}{2} \implies L \geq \epsilon\sqrt{c} \sqrt{\frac{nd}{p\log(1+JW\epsilon^{-1}c^{-1}) + \log(4/\delta)}}$$

(In the statement we had $\frac{\epsilon}{\sigma}$ in place of $\epsilon$, but we have assumed $\sigma$ bounded away from 0 for simplicity. The paper [Bubeck, Sellke, 2021] has a slightly more refined analysis that captures the dependence on $\frac{\epsilon}{\sigma}$ ).

EXTRA :

Hoeffding's inequality for bounded random variables.

THEOREM   Let $X_1, \dots, X_N$ be independent random variables. Assume that $X_i \in [m_i, M_i]$ for all $i$. Then, for every $t > 0$, we have

$$\mathbb{P}\left( \sum_{i=1}^{N} (X_i - \mathbb{E}X_i) \geq t \right) \leq e^{-2t^2 / \sum_{i=1}^{N}(M_i - m_i)^2}$$

$\epsilon$-net :   $N \subseteq K$ is an $\epsilon$-net of $K$ if every point in $K$ is within distance $\epsilon$ of some point of $N$

Formally,   $\forall x \in K \quad \exists x_0 \in N : \quad \underbrace{\|x - x_0\|}_{} \leq \epsilon$

valid for general metric spaces
$(d(x, x_0) \leq \epsilon)$

THEOREM   Let $N$ be an $\epsilon$-net of the Euclidean ball in $n$ dimensions with unit radius. Then,

$$\left(\frac{1}{\epsilon}\right)^n \leq |N| \leq \left(\frac{2}{\epsilon} + 1\right)^n$$

# A SUFFICIENT CONDITION FOR NTK FEATURES

Let's go back to :

$$x_{adv} = x + \Delta_{adv}, \quad \|\Delta_{adv}\| \leq \delta \cdot \|x\|$$

$$|f(x_{adv}, \vartheta) - f(x, \vartheta)| \simeq |\nabla_x f(x, \vartheta)^T \Delta_{adv}| \leq \delta \cdot \underbrace{\|x\| \cdot \|\nabla_x f(x, \vartheta)\|}_{=:}$$

$$S_f(x)$$

SENSITIVITY of the model evaluated in $x$

size of the labels

The model $f$ is ROBUST if $S_f(x) = O(1)$ for most test samples $x$. This means that the output change $|f(x_{adv}, \vartheta) - f(x, \vartheta)|$ is of the same order as $|f(x, \vartheta)|$. In contrast, if $S_f(x) \gg 1$, we expect the model to be vulnerable to adversarial perturbations.

## NTK FEATURES.

Consider the two-layer network

$$f_{NN}(x, w) = \sum_{i=1}^{k} \phi(\langle w_i^{(1)}, x \rangle) - \sum_{i=1}^{k} \phi(\langle w_i^{(2)}, x \rangle)$$

*) $2k$ neurons

*) Activation function $\phi$

*) # parameters $p = 2kd$, $W^{(1)} = [w_1^{(1)}, \cdots, w_k^{(1)}]$, $W^{(2)} = [w_1^{(2)}, \cdots, w_k^{(2)}]$
$W = [W^{(1)}, W^{(2)}]$

Initialization $W_0 = [W_0^{(1)}, W_0^{(2)}]$ with $[W_0^{(1)}]_{is} \overset{iid}{\sim} N(0, 1/d)$ and $W_0^{(2)} = W_0^{(1)}$.

Under certain conditions, the GD dynamics of $f_{NN}(x, w)$ is close to the GD dynamics of its linearization around $w_0$
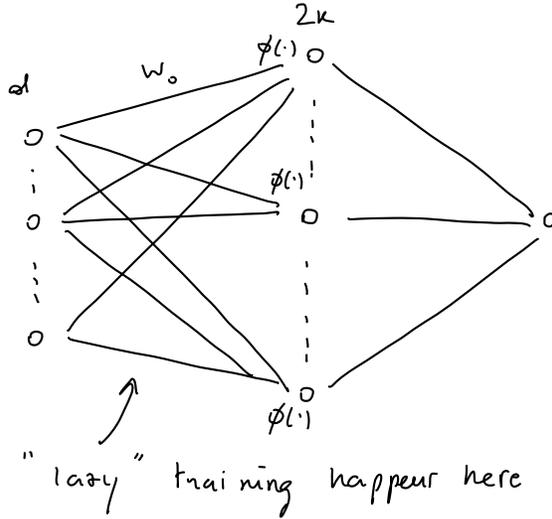
$$\underbrace{f_{NN}(x, w_0)}_{= 0 \text{ by symmetry}} + \langle \nabla_w f_{NN}(x, w_0), w - w_0 \rangle$$

Define
$$f_{NTK}(x, \vartheta) = \Phi_{NTK}(x)^T \vartheta, \qquad \Phi_{NTK}(x) = \nabla_w f_{NN}(x, w_0)$$

*) $\vartheta \in \mathbb{R}^p$, $p = 2kd$, vector of trainable parameters

$$\vartheta^* = \vartheta_0 + \Phi_{NTK}^T \left( \Phi_{NTK} \Phi_{NTK}^T \right)^{-1} y, \quad \text{with} \quad \Phi_{NTK} = \begin{bmatrix} \Phi_{NTK}(x_1) \\ \vdots \\ \Phi_{NTK}(x_n) \end{bmatrix} \in \mathbb{R}^{n \times p}$$

obtained by running GD from initialization $\vartheta_0$



"lazy" training happens here

$$S_{f_{NTK}}(x) = \|x\| \cdot \left\| \nabla_x \Phi_{NTK}^T(x) \, \Phi_{NTK}^T \left( \Phi_{NTK} \Phi_{NTK}^T \right)^{-1} y + \underbrace{\nabla_x \Phi_{NTK}^T(x) \vartheta_0}_{= 0 \text{ by symmetry}} \right\|$$

Data assumptions:

*) $\|x_i\| = \sqrt{d}$ for all $i$  —  data on the sphere of radius $\sqrt{d}$

*) $\mathbb{E} x_i = 0$  —  data centered

*) The data distribution $P_x$ satisfies $c$-Lipschitz concentration   (dimension-independent constant)

$$\mathbb{P}\left( \left| \varphi(x) - \int \varphi(x') \, dP_{x|x'} \right| > t \right) \le 2e^{-t^2 / 2c \cdot Lip(\varphi)^2}$$

__THEOREM__ [Bombari, Kiyani, M., 2023]. Let $\phi$ be non-linear, even, with Lipschitz derivative. Assume $n \, polylog(n) = o(kd)$, $k = O(d)$, $n = O(k)$. Then, with high probability,

$$S_{f_{NTK}}(x) = O\left( \log n \sqrt{\frac{nd}{p}} \right).$$

# INTERPRETATION

*) Same condition on overparameterization as in the lower bound by

[Bubeck, Sellke, 2021]

*) Lower bound is for Lipschitz constant $\left( \sup_x \| \nabla_x f(x) \| \right)$, while the upper bound above holds with high probability over a test point.

# PROOF SKETCH

$$S_{f_{NTK}}(x) \leq \| x \|_2 \; \| \nabla_x \Phi_{NTK}^T(x) \; \Phi_{NTK}^T \|_{op} \; \| ( \Phi_{NTK} \Phi_{NTK}^T )^{-1} \|_{op} \; \| y \|$$

*) $\| x \|_2 = \sqrt{d}$      data normalization

*) $\| y \| = \sqrt{n}$      iid data with $O(1)$ $y_i$

*) $\| ( \Phi_{NTK} \Phi_{NTK}^T )^{-1} \|_{op} = \dfrac{1}{\lambda_{min}( \Phi_{NTK} \Phi_{NTK}^T )} = O\left( \dfrac{1}{d\kappa} \right)$

previous analysis

*) $\| \nabla_x \Phi_{NTK}^T(x) \; \Phi_{NTK}^T \|_{op} = O\left( \log \kappa \, ( \sqrt{\kappa} + \sqrt{n} ) \sqrt{d} \right)$

direct calculation via chain rule

Thus, $S_{f_{NTK}}(x) = O\left( \sqrt{nd} \; \dfrac{1}{d\kappa} \; \log \kappa \; ( \overbrace{\sqrt{\kappa}}^{O(\sqrt{\kappa})} + \sqrt{n} ) \sqrt{d} \right)$

$$= O\left( \log \kappa \; \sqrt{\dfrac{nd}{d\kappa}} \right).$$

$\rho$

☑

Both model ( NTK features ) and activation function ( even ) matter for robustness ...

RANDOM FEATURES (RF) [Rahimi, Recht, 2007]

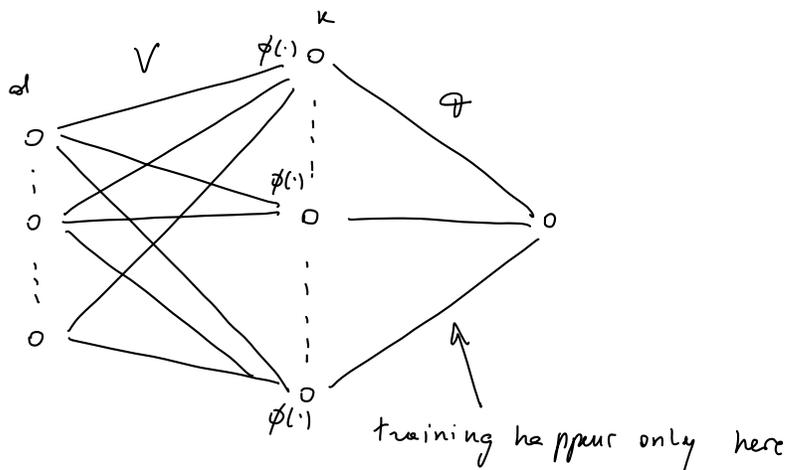$$f_{RF}(x, \vartheta) = \Phi_{RF}^T(x) \vartheta \quad , \quad \Phi_{RF}(x) = \phi(Vx)$$

*) $\quad V \in \mathbb{R}^{K \times d} \qquad V_{ij} \overset{iid}{\sim} N(0, 1/d)$

*) $\quad \vartheta \in \mathbb{R}^p \quad , \quad p = K \quad ,$ vector of trainable parameters

$$\vartheta^* = \Phi_{RF}^T \left( \Phi_{RF} \Phi_{RF}^T \right)^{-1} y \quad , \quad \text{with} \quad \Phi_{RF} = \begin{bmatrix} \Phi_{RF}(x_1) \\ \vdots \\ \Phi_{RF}(x_n) \end{bmatrix} = \phi(X V^T) \in \mathbb{R}^{n \times K}$$

$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n \times d}$

obtained by running GD from initialization $\vartheta_0 = 0$



training happens only here

*) RF model can be regarded as a two-layer neural network with first-layer weights frozen at initialization.

*) Simplest model allowing a degree of freedom in the number of parameters ( in linear regression , # parameters = input dimension )

*) Lots of theoretical work on RF model , and we'll see more later in this course...

$$S_{f_{RF}}(x) = \|x\| \cdot \| \nabla_x \Phi_{RF}^T(x) \Phi_{RF}^T \left( \Phi_{RF} \Phi_{RF}^T \right)^{-1} y \|$$

## THEOREM [Bombari, Kiyani, M., 2023].

Assume the data assumptions hold, and let $y_i = g(x_i) + \epsilon_i$ with $\epsilon_i$ zero-mean and sub-Gaussian. Let $\phi$ be non-linear, Lipschitz, with two Lipschitz derivatives and such that $\underset{G \sim N(0,1)}{\mathbb{E}} \phi'(G) \neq 0$. Assume $n \, \text{polylog}(n) = o(k)$, $d \, \text{polylog}(d) = o(k)$, and $n \, \text{polylog}(n) = o(d^{3/2})$. Then, with high probability,

$$S_{f_{RF}}(x) = \Omega\left(n^{1/6}\right) \gg 1.$$

## INTERPRETATION

*) RF model NOT ROBUST for any degree of overparameterization

*) $\underset{G \sim N(0,1)}{\mathbb{E}} \phi'(G) \neq 0$ may be fundamental (confirmed by simulations)

## PROOF SKETCH

STEP 1 : Fitting the noise lower bounds the sensitivity

Let $\mathbb{E}[\epsilon^2] = \varepsilon^2$ and $A(x) = \nabla_x \Phi_{RF}^T(x) \, \Phi_{RF}^T \left(\Phi_{RF} \Phi_{RF}^T\right)^{-1} \in \mathbb{R}^{d \times n}$. Then, with high probability,

$$S_{f_{RF}}(x) \geq \frac{\varepsilon}{2} \|x\| \, \|A(x)\|_F$$

Idea : Decompose $y_i$ into $g(x_i)$ and $\epsilon_i$ + use Hanson-Wright

$$S_{f_{RF}}(x) = \|x\| \cdot \|A(x) y\| = \|x\| \cdot \|A(x) \underbrace{(g(X) + \epsilon)}\|$$

$$\begin{bmatrix} g(x_1) + \epsilon_1 \\ \vdots \\ g(x_n) + \epsilon_n \end{bmatrix}$$

remove dependence on $x$ from $A(x)$ for simplicity of notation

$$\|A y\|^2 = \epsilon^T A^T A \epsilon + g(X)^T A^T A \, g(X) + 2 g(X)^T A^T A \epsilon$$

$$\geq \epsilon^T A^T A \epsilon + g(X)^T A^T A \, g(X) - \varepsilon \|A g(X)\| \|A\|_F$$

holds with high probability upon concentrating $2 g(X)^T A^T A \epsilon$

$$= \epsilon^T A^T A \epsilon + M^2 - \epsilon M \|A\|_F$$

$$\left( \quad M := \|A g(x)\| \right.$$

$$\geq \epsilon^T A^T A \epsilon - \frac{\epsilon^2 \|A\|_F^2}{4} \geq \frac{\epsilon^2}{2} \|A\|_F^2 - \frac{\epsilon^2}{4} \|A\|_F^2 = \frac{\epsilon^2}{4} \|A\|_F^2$$

$$\left( \text{minimize over } M \right. \qquad \left( \text{holds with high probability upon concentrating} \right.$$

$$\epsilon^T A^T(x) A(x) \epsilon \quad \text{via Hanson-Wright}$$

## STEP 2: Split interaction and kernel components

Let $\quad \underbrace{\widetilde{\mathcal{I}}_{nf}(x)}_{\text{interaction matrix}} = \nabla_x \Phi_{nf}^T(x) \, \widetilde{\Phi}_{nf}^T \in \mathbb{R}^{n \times n}$ with $\widetilde{\Phi}_{nf} = \Phi_{nf} - \underset{x_1,\cdots,x_n}{\mathbb{E}} \Phi_{nf}$.

Then, with high probability,

$$\|A(x)\|_F \geq \frac{\|\mathcal{I}_{nf}(x)\|_F}{\lambda_{max}\left(\widetilde{\Phi}_{nf} \widetilde{\Phi}_{nf}^T\right)} - O\left(\frac{\sqrt{n+d}}{d}\right) = \Omega\left(\frac{1}{\kappa} \frac{d}{n+d} \|\mathcal{I}_{nf}(x)\|_F - \frac{\sqrt{n+d}}{d}\right)$$

$$\|A(x)\|_F \leq \frac{\|\mathcal{I}_{nf}(x)\|_F}{\lambda_{min}\left(\widetilde{\Phi}_{nf} \widetilde{\Phi}_{nf}^T\right)} + O\left(\frac{\sqrt{n+d}}{d}\right) = O\left(\frac{1}{\kappa} \|\mathcal{I}_{nf}(x)\|_F + \frac{\sqrt{n+d}}{d}\right)$$

Ideas:

*) Center $A(x)$

$\overbrace{\nabla_x \Phi_{nf}^T(x) \, \Phi_{nf}^T \left(\Phi_{nf} \Phi_{nf}^T\right)^{-1}}$

$$\left| \|A(x)\|_F - \|\nabla_x \Phi_{nf}^T(x) \, \widetilde{\Phi}_{nf}^T \left(\widetilde{\Phi}_{nf} \widetilde{\Phi}_{nf}^T\right)^{-1}\|_F \right| = O\left(\frac{\sqrt{n+d}}{d}\right)$$

(extract two rank-1 terms via Sherman-Morrison + various estimates)

*) Bounds on the spectrum of $\widetilde{\Phi}_{nf} \widetilde{\Phi}_{nf}^T$

$\overbrace{\nabla_x \mathcal{I}_{nf}^T(x) \, \widetilde{\Phi}_{nf}^T}$

$$\underbrace{\frac{\|\mathcal{I}_{nf}(x)\|_F}{\lambda_{max}\left(\widetilde{\Phi}_{nf} \widetilde{\Phi}_{nf}^T\right)}}_{O(\kappa(n+d)/d)} \leq \|\nabla_x \Phi_{nf}^T(x) \, \widetilde{\Phi}_{nf}^T \left(\widetilde{\Phi}_{nf} \widetilde{\Phi}_{nf}^T\right)^{-1}\|_F \leq \underbrace{\frac{\|\mathcal{I}_{nf}(x)\|_F}{\lambda_{min}\left(\widetilde{\Phi}_{nf} \widetilde{\Phi}_{nf}^T\right)}}_{\Omega(\kappa)}$$

$$\left| \| \mathcal{I}_{RF}(x) \|_F - \left( \mathop{\mathbb{E}}_{G \sim N(q_1)} \phi'(G) \right)^2 \frac{\kappa \sqrt{n}}{\sqrt{d}} \right| = o\left( \frac{\kappa \sqrt{n}}{\sqrt{d}} \right)$$

If $\mathop{\mathbb{E}}_{G \sim N(0,1)} \phi'(G) = 0$, then $\| \mathcal{I}_{RF}(x) \|_F$ is of lower order and the

sensitivity can also be much smaller.

Idea :

$\overset{\text{direct calculation}}{\searrow}$

$$\| \mathcal{I}_{RF}(x) \|_F^2 = \sum_{i=1}^{n} \| V^T \text{diag}(\phi'(Vx)) \, \tilde{\phi}(Vx_i) \|^2$$

$$\tilde{\phi}(Vx_i) = \phi(Vx_i) - \underset{\kappa_i}{\mathbb{E} \, \phi(Vx_i)}$$

*) $\text{diag}(\phi'(Vx))$ and $\tilde{\phi}(Vx_i)$ depend on $V$ only via a single projection

*) We split $\text{diag}(\phi'(Vx))$ and $\tilde{\phi}(Vx_i)$ via Taylor expansion into a component correlated with $V$ and an independent one

*) Correlated components are computed exactly and independent ones are negligible

$$\| \mathcal{I}_{RF}(x) \|_F^2 \approx n \underbrace{\| V^T \text{diag}(\phi'(Vx)) \, \tilde{\phi}(Vx_1) \|^2}_{\sum_{J=1}^{d} (V^T \text{diag}(\phi'(Vx)) \, \tilde{\phi}(Vx_1))_J^2}$$

$$= n \sum_{J=1}^{d} \left[ |(x_1)_J \left( \mathop{\mathbb{E}}_{G \sim N(q_1)} \phi'(G) \right)^2 \frac{\kappa}{d} | + o\left( \frac{\kappa}{d} \right) \right]^2$$

$$= n \cancel{d} \left( \mathop{\mathbb{E}}_{G \sim N(0,1)} \phi'(G) \right)^4 \frac{\kappa^2}{d^2} (1 + o(1))$$

Putting everything together :

$$S_{\mathcal{F}_{MF}}(x) = \|x\| \cdot \|A(x) y\|$$

$$= \Omega\left(\sqrt{d}\,\|A(x)\|_F\right)$$

$$= \Omega\left(\sqrt{d}\left[\frac{1}{k}\,\frac{d}{n+d}\,\|\mathcal{I}_{MF}(x)\|_F - \frac{\sqrt{n+d}}{d}\right]\right)$$

$$= \Omega\left(\sqrt{d}\left[\frac{1}{k}\,\frac{d}{n+d}\,\frac{k\sqrt{n}}{\sqrt{d}} - \frac{\sqrt{n+d}}{d}\right]\right)$$

$$= \Omega\left(n^{1/6}\right)$$

$$\frac{1}{k}\,\frac{d}{n+d}\,\frac{k\sqrt{n}}{\sqrt{d}} = \frac{\sqrt{nd}}{n+d} \gg \frac{\sqrt{n+d}}{d} \iff d^{3/2}\,n^{1/2} \gg (n+d)^{3/2} = \Theta\left(\max\left(n^{3/2}, d^{3/2}\right)\right)$$

Distinguish two cases :

① $n \leq d$

$$(n+d)^{3/2} = \Theta(d^{3/2}) \ll d^{3/2}\,n^{1/2}$$

$$\sqrt{d}\left[\frac{1}{k}\,\frac{d}{n+d}\,\frac{k\sqrt{n}}{\sqrt{d}} - \frac{\sqrt{n+d}}{d}\right] = \Omega\left(\sqrt{d}\,\frac{\sqrt{nd}}{d}\right) = \Omega(\sqrt{n})$$

② $n > d$

$$(n+d)^{3/2} = \Theta(n^{3/2}) \ll n^{1/2}\,d^{3/2} \iff n \ll d^{3/2} \quad (\text{by hypothesis})$$

$$\sqrt{d}\left[\frac{1}{k}\,\frac{d}{n+d}\,\frac{k\sqrt{n}}{\sqrt{d}} - \frac{\sqrt{n+d}}{d}\right] = \Omega\left(\sqrt{d}\,\frac{\sqrt{nd}}{n}\right) = \Omega\left(\frac{d}{\sqrt{n}}\right) = \Omega(n^{1/6})$$

$$n \ll d^{3/2}$$

# HOW MANY PARAMETERS TO RECONSTRUCT TRAINING DATA? ( MEMORIZATION #2)

[ start with slides ]

# RANDOM FEATURES (RF)

$$f_{RF}(x, \vartheta) = \Phi_{RF}^T(x) \vartheta \quad , \quad \Phi_{RF}(x) = \phi(Vx)$$

\* ) $\quad V \in \mathbb{R}^{p \times d} \quad \quad V_{ij} \overset{iid}{\sim} N(0, 1/d)$

# ASSUMPTIONS :

(B1) Let the training samples $(x_i)_{i=1}^n$ be iid, sub-Gaussian (with sub-Gaussian norm having constant order) and $\ell_2$ norm equal to $\sqrt{d}$.

(data distribution)

(B2) The activation function $\phi : \mathbb{R} \longrightarrow \mathbb{R}$ is non-linear, Lipschitz, with Lipschitz derivative. Letting $\phi_\ell$ denote its $\ell$-th Hermite coefficient, we assume $\phi_0 = \phi_2 = 0$, $\phi_1 \neq 0$ and that there exist two non-zero Hermite coefficients of order $\geqslant 3$ with different parity.

(activation function)

resolves sign ambiguity in the reconstruction (more on this later)

(B3) $\quad n = O(d)$ and $p = \tilde{\Omega}(nd)$

(overparameterisation)

THEOREM [Iurada, Bombari, Tommasi, M., 2026] Let (B1)-(B2)-(B3) hold.

Let $\hat{X} = \begin{bmatrix} \hat{x}_1 \\ \vdots \\ \hat{x}_n \end{bmatrix} \in \mathbb{R}^{n \times d}$ be such that $\|\hat{x}_i\| = \sqrt{d}$ and, for all $i \in \{1, \dots, n\}$,

$\phi(Vx_i) \in \text{span} \{ \phi(V\hat{x}_1), \dots, \phi(V\hat{x}_n) \}$. Then, with high probability,

for any $\hat{i} \in \{1, \dots, n\}$, there exists $i \in \{1, \dots, n\}$ such that

$$\| \hat{x}_{\hat{i}} - x_i \| = o(\sqrt{d})$$

Each row of $\hat{X}$ is close to a training sample

<u>INTERPRETATION</u>  If random features of training samples are spanned by random features of $X$, the rows of $\hat{X}$ must be close to original training samples, as long as $p \gg dn$.

\*) Existence of two non-zero Hermite coefficients with different parity is a necessary condition to recover the sign of the training samples.

$\longrightarrow$ For example, if $\phi$ is either even or odd, the problem is under-determined in terms of the signs of the $\hat{x}_i$'s, as span $\{\phi(V\hat{x}_1), \cdots, \phi(V\hat{x}_n)\}$ does not depend on them.

$\longrightarrow$ Evident in numerical simulations (for $\phi = ReLU$, where $\phi_{2\ell+1} = 0$ for $\ell > 1$, negatives of training samples may be reconstructed).

$\longrightarrow$ Can drop the condition if only interested in overlap :

$$\frac{|<\hat{x}_\ell , x_i >|}{\| \hat{x}_\ell \| \, \| x_i \|} = 1 - o(1)$$

# PROOF SKETCH.

## STEP 1 : Decompose features of $\hat{X}$ into features of training samples

For any $\hat{x}$ row of $\hat{X}$, we have

$$\phi(V\hat{x}) = \sum_{i=1}^{n} a_i \phi(Vx_i) \quad \text{for some } (a_i)_{i=1}^{n} \qquad (3)$$

*) Using ideas similar to those discussed when bounding the smallest eigenvalue of the NTK, we have that $\lambda_{min}(\Phi\Phi^T) > 0$, with $\Phi = \begin{bmatrix} \phi(Vx_1) \\ \vdots \\ \phi(Vx_n) \end{bmatrix}$.

*) By hypothesis, $\phi(Vx_i) \in \text{span} \{ \phi(V\hat{x}_1), \cdots, \phi(V\hat{x}_n) \}$ for all $i$, which gives

$$\text{span} \{ \phi(Vx_1), \cdots, \phi(Vx_n) \} \subseteq \text{span} \{ \phi(V\hat{x}_1), \cdots, \phi(V\hat{x}_n) \}$$

As $\dim(\text{span} \{ \phi(Vx_1), \cdots, \phi(Vx_n) \}) = n$ $(\lambda_{min}(\Phi\Phi^T) > 0)$, the two spans are in fact equal, which gives

$$\phi(V\hat{x}) \in \text{span} \{ \phi(Vx_1), \cdots, \phi(Vx_n) \} \quad \text{for any } \hat{x} \text{ row of } \hat{X}$$

↑ this is equivalent to (3)

## STEP 2 : Looking at the decomposition in one direction

Let $\tilde{\phi}(t) = \phi(t) - \phi_1 t$ be obtained by removing the linear part in $\phi(\cdot)$.

Linear part is problematic and, in fact, if $\phi$ was linear, then there is a simple counterexample to the statement:

Pick $\hat{x}_i = \sqrt{d} \dfrac{x_i + x_{i+1}}{\| x_i + x_{i+1} \|}$ $i \in \{1, \cdots, n-1\}$, $\hat{x}_n = \sqrt{d} \dfrac{x_n - x_1}{\| x_n - x_1 \|}$.

Then $Vx_i \in \text{span} \{ V\hat{x}_1, \cdots, V\hat{x}_n \}$ for all $i \in \{1, \cdots, n\}$, but

$(\hat{x}_i)_{i=1}^{n}$ are not close to $(x_i)_{i=1}^{n}$.

$$\tilde{\phi}(V\hat{x})^T \phi(V\hat{x}) \overset{(3)}{=} \sum_{i=1}^{n} a_i \, \tilde{\phi}(V\hat{x})^T \phi(Vx_i) = \sum_{i=1}^{n} a_i \sum_{J=1}^{p} \tilde{\phi}(<v_J, \hat{x}>)\phi(<v_J, x_i>)$$

$$V = \begin{bmatrix} v_1 \\ \vdots \\ v_p \end{bmatrix} \in \mathbb{R}^{p \times d}$$

$$= \mathbb{E}_V \left[ \sum_{i=1}^{n} a_i \sum_{J=1}^{p} \tilde{\phi}(<v_J, \hat{x}>)\phi(<v_J, x_i>) \right] \left( 1 + \tilde{O}\left( \|a\| \sqrt{\frac{dn}{p}} \right) \right)$$

Bernstein inequality as $\tilde{\phi}(<v_J, \hat{x}>)\phi(<v_J, x_i>)$ is sub-exponential

$$= p \sum_{i=1}^{n} a_i \sum_{\ell=3}^{\infty} \phi_\ell^2 \, \frac{<\hat{x}, x_i>^\ell}{d^\ell} \left( 1 + \tilde{O}\left( \|a\| \sqrt{\frac{dn}{p}} \right) \right)$$

Hermite expansion:

$$\mathbb{E}_{v_J} \tilde{\phi}(<v_J, \hat{x}>)\phi(<v_J, x_i>) = \mathbb{E}_{\bar{v} \sim N(0, I_d)} \tilde{\phi}\left(<\bar{v}, \frac{\hat{x}}{\sqrt{d}}>\right)\phi\left(<\bar{v}, \frac{x_i}{\sqrt{d}}>\right)$$

$$= \sum_{\ell \in \mathbb{N}} \phi_\ell \tilde{\phi}_\ell \, \frac{<x_i, \hat{x}>^\ell}{d^\ell} = \sum_{\ell \geqslant 3} \phi_\ell^2 \, \frac{<\hat{x}, x_i>^\ell}{d^\ell}$$

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \text{Hadamard product}$$

$$\leq p \|a\| \sum_{\ell=3}^{\infty} \phi_\ell^2 \left\| \frac{(X\hat{x})^{\circ\ell}}{d^\ell} \right\| \left( 1 + \tilde{O}\left( \|a\| \sqrt{\frac{dn}{p}} \right) \right)$$

$$\leq p \sum_{\ell=3}^{\infty} \phi_\ell^2 \left\| \frac{(X\hat{x})^{\circ\ell}}{d^\ell} \right\| \left( 1 + \tilde{O}\left( \sqrt{\frac{dn}{p}} \right) \right)$$

$$\left| \|a\|^2 - 1 \right| = \tilde{O}\left( \sqrt{\frac{dn}{p}} \right)$$

$$\leq p \sum_{\ell=3}^{\infty} \phi_\ell^2 \left( \max_i \left| \frac{x_i^T \hat{x}}{d} \right| \right)^{\ell-1} \left( 1 + \tilde{O}\left( \sqrt{\frac{dn}{p}} \right) \right)$$

$$\underbrace{\max_i \left| \frac{x_i^T \hat{x}}{d} \right|}_{:= C \in [0,1]}.$$

$\hat{x}$ can't be too correlated with too many $x_i$'s

$$\leq p \underbrace{\sum_{\ell=3}^{\infty} \phi_\ell^2}_{\bar{\phi}^2} C^2 \left( 1 + \tilde{O}\left( \sqrt{\frac{dn}{p}} \right) \right)$$

However, we also have

(similar argument with Bernstein + Hermite expansion)

$$\tilde{\phi}(V\hat{x})^T \phi(V\hat{x}) = \bar{\phi}^2 p (1 + o(1))$$

This implies $C = 1 + o(1)$, so $\left| \frac{x_i^T \hat{x}}{d} \right|$ is close to 1 for some $i$.

*) Bounds need to hold uniformly over $\hat{x}$, so we need (a few) $\epsilon$-net arguments.

*) To pass from overlap $\left( \left| \frac{x_i^T \hat{x}}{d} \right| = 1 + o(1) \right)$ to $\ell_2$ distance

$\left( \| \hat{x} - x_i \| = o(\sqrt{d}) \right)$, we use the existence of Hermite coefficients with different parities to resolve the ambiguity on the sign.

With a similar argument, we have that $\left| \bar{\phi}^2 - a_i \sum_{\ell=3}^{\infty} \phi_\ell^2 \frac{(x_i^T \hat{x})^\ell}{d^\ell} \right| = o(1)$.

This is possible only if $(x_i^T \hat{x})^\ell$ has the same sign for all $\ell \geq 3$ s.t. $\phi_\ell \neq 0$.

∎

⌐ EXTRA :

Bernstein inequality

THEOREM   Let $X_1, \dots, X_N$ be independent, mean zero, sub-exponential random variables. Then, for every $t \geq 0$, we have

$$\mathbb{P}\left( \left| \sum_{i=1}^{N} X_i \right| > t \right) \leq 2 \exp\left[ -c \min\left( \frac{t^2}{\sum_{i=1}^{N} \| X_i \|_{\psi_1}^2}, \frac{t}{\max_i \| X_i \|_{\psi_1}} \right) \right],$$

where $\| \cdot \|_{\psi_1}$ denotes the sub-exponential norm.

\*) Result above allows for duplicates (multiple rows of $\hat{X}$ corresponding to the same training sample). We can further show that <u>ALL</u> training samples are reconstructed (i.e., no duplicates) for $n = 2$.

<span style="color:blue">

↗

proof quite ad-hoc for $n = 2$. For $n = 3$, we would need to consider separately the case in which $\hat{x}_1, \hat{x}_2$ and $\hat{x}_3$ are close and the case in which just a pair of reconstructions is close. Thus, # cases increases combinatorially with $n$.
</span>

$\boxed{\text{OPEN PROBLEM}}$ Show all training samples are reconstructed for general $n$.

\*) Assumption that $\phi(Vx_i) \in \text{span} \{\phi(V\hat{x}_1), \cdots, \phi(V\hat{x}_n)\}$ is quite strong ...

$$\phi(Vx_i) \in \overbrace{\text{span} \{\phi(V\hat{x}_1), \cdots, \phi(V\hat{x}_n)\}}^{, S}$$

$$\Downarrow$$

$$\| P_{\hat{X}}^{\perp} \phi(Vx_i) \|_2 = 0 \quad \text{with} \quad P = \text{projector on } S$$

$$\theta^* = \Phi_{RF}^T (\Phi_{RF} \Phi_{RF}^T)^{-1} y, \quad \text{with} \quad \Phi_{RF} = \begin{bmatrix} \Phi_{RF}(x_1) \\ \vdots \\ \Phi_{RF}(x_n) \end{bmatrix} = \phi(XV^T) \in \mathbb{R}^{n \times p}$$

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n \times d}$$

<span style="color:blue">obtained by running GD from initialization $\theta_0 = 0$</span>

$\theta^*$ is a linear combination of $\phi(Vx_i)$ so instead of optimizing for

$\| P_{\hat{X}}^{\perp} \phi(Vx_i) \|_2 = 0$ for all $i$ (which we cannot do in practice since we don't know $\phi(Vx_i)$), we just optimize for $\| P_{\hat{X}}^{\perp} \theta^* \|_2 = 0$.

$$\hat{X}^* = \underset{\hat{X} : \|\hat{x}_i\| = \sqrt{d}}{\text{argmin}} \| P_{\hat{X}}^{\perp} \theta^* \|^2$$

*) Numerical evidence (shown before) that $\hat{x}^*$ reconstructs $X$ when $p \gg dn$.

*) When $p \gg nd$, $\phi(Vx_i)$ approximately lies in span $\{\phi(V\hat{x}_i), \cdots, \phi(V\hat{x}_n)\}$.

[OPEN PROBLEM] Show that optimizing $\| P_{\hat{X}}^{\perp} \theta^* \|$ leads to the reconstruction of the training samples.

[OPEN PROBLEM] Show similar results beyond the RF model (e.g. NTK model).

# RECAP

* $p \gg n$    gives    interpolation

*) $p \gg dn$    gives    smooth    interpolation    ( no adverrarial exaumpler )

but it allows an adverrary to recourruct the training datalet
from the paramerers of the trained network.

How to defend from this adverrary ?

One option is DIFFERENTIAL PRIVACY which will be the subject of the
rest of the course.

# DEEP LEARNING WITH DIFFERENTIAL PRIVACY

<u>What is differential privacy?</u>

Textbook : [Dwork, Roth, 2014] "The Algorithmic Foundations of Differential Privacy"

<u>DEFINITION</u>  A dataset $D'$ is <u>adjacent</u> to a dataset $D$ if they differ by only one sample.

<u>DEFINITION</u>  $((\varepsilon, \delta) - DP)$  [Dwork, McSherry, Nissim, Smith, 2006]. A randomized algorithm $A$ satisfies <u>$(\varepsilon, \delta)$- differential privacy</u> if, for any pair of adjacent datasets $D, D'$ and for any subset of the parameter space $S \subseteq \mathbb{R}^p$, we have

$$\mathbb{P}\left( A(D) \in S \right) \leq e^{\varepsilon} \mathbb{P}\left( A(D') \in S \right) + \delta$$

probability over the randomness induced by the algorithm

holds uniformly on all adjacent datasets $D, D'$

$A$ is private if one cannot distinguish $A(D)$ from $A(D')$. This means that given the output of the algorithm, it is (information - theoretically) impossible to recover a single training sample.

\*)  $(\varepsilon, \delta) = (0, 0) \implies$ perfect privacy $\left( \mathbb{P}(A(D) \in S) = \mathbb{P}(A(D') \in S) \right)$

\*)  $\varepsilon \gg 1$  or  $\delta = 1 \implies$ no privacy at all

\*)  In practice  $\varepsilon \sim 1$, $\delta \ll 1/n$

number of training samples

Many other notions in the literature. We will see later Renyi- DP and zero -concentrated DP.

# How to enforce it?

Essentially add noise ...

*) Different types of noise: Laplace mechanism / noise, Gaussian mechanism / noise

*) Different ways of adding it: inside the empirical risk minimization objective (objective perturbation: perturb the problem and then solve it), after calculating the ERM solution (output perturbation: perturb the solution) during the iterative algorithm (DP-GD / DP-SGD: perturb the iterations)

Here, we focus uniquely on adding noise to gradient descent - type algorithms.

## Differentially - private gradient descent (DP-GD)

$$\theta^{t+1} = \theta^t - \eta \frac{1}{n} \sum_{i=1}^{n} \nabla_\theta \ell(y_i, f(x_i; \theta^t)) / \max\left(1, \frac{\|\nabla_\theta \ell(y_i, f(x_i;\theta^t))\|_2}{C_{clip}}\right)$$

$$+ \sqrt{\eta} \frac{2 C_{clip}}{n} \sigma N(0, I_p)$$

→ Gradient is clipped (its norm is at most $C_{clip}$)

→ Additive Gaussian noise (Gaussian mechanism for privacy)

Informally, this avoids that the algorithm memorizes each individual data point.
More formally, we have the following:

**Proposition** For any $\delta \in (0,1)$, $\varepsilon \in (0, 8\log(1/\delta))$, if we set

$$\sigma \geq \sqrt{\eta T} \frac{\sqrt{8\log(1/\delta)}}{\varepsilon},$$ then the output $\theta^T$ of DP-GD is

$(\varepsilon, \delta)$ - differentially private.

*) Adaptation of the analysis from [Abadi, Chu, Goodfellow, 2016]

# What's the cost of differential privacy?

## Excess population risk :

$$R_E = \mathbb{E}\Big[\ell\big(y_{test}, f(x_{test}; \theta^T)\big)\Big] - \mathbb{E}\Big[\ell\big(y_{test}, f(x_{test}; \theta^*)\big)\Big]$$

solution of DP-GD with privacy guarantees    solution of GD with no privacy guarantees

## Excess empirical risk :

$$\hat{R}_E = \frac{1}{n}\sum_{i=1}^{n}\ell\big(y_i, f(x_i; \theta^T)\big) - \frac{1}{n}\sum_{j=1}^{n}\ell\big(y_i, f(x_i; \theta^*)\big)$$

More privacy $\Rightarrow$ $\varepsilon \downarrow$ , $\delta \downarrow$

$\Rightarrow$ $\sigma \uparrow$ (or equivalently $c_{clip} \downarrow$ )

proposition above

$\Rightarrow R_E, \hat{R}_E \uparrow$ ( worse performance )

Existing bounds on excess empirical / population risk tend to degrade as $p$ grows :

*) $\quad \hat{R}_E = \tilde{O}\left(\frac{\sqrt{p}}{n\varepsilon}\right)$ (objective perturbation, constrained, strongly convex

$(\varepsilon, \delta)$-DP optimization) [Kifer, Smith, Thakurta, 2012]

*) $\quad \hat{R}_E = \beta L\, \tilde{O}\left(\frac{\sqrt{p}}{n\varepsilon}\right)$ ( DP-GD, $\beta$ = diameter of optimization domain ,

$L$ = Lipschitz constant of the loss) , $R_E = \tilde{O}\left(\frac{p^{1/4}}{\sqrt{n\varepsilon}}\right)$

[Bassily, Smith, Thakurta, 2014]

*) $\mathbb{E}[R_E] = p L \; \widetilde{O}\left(\frac{1}{\sqrt{n}} + \frac{\sqrt{p}}{n\varepsilon}\right)$    [Bassily, Feldman, Talwar, Thakurta, 2019]

↑

expectation over randomness of algorithm and training data

*) Lots of work on unconstrained optimization too, but still dependence on dimension $p$

Noise introduced by DP-GD increases with dimension $p$ (each entry of $\theta \in \mathbb{R}^p$ perturbed with standard Gaussian noise $\Rightarrow$ $\ell_2$ norm of perturbation scales as $\sqrt{p}$).

*) DP algorithms acting over lower-dimensional subspaces

We will show that this bad dependence does not have to be there!

# DP TOOLS AND A PROOF OF THE PROPOSITION

**DEFINITION** Let $\mu : \mathcal{D} \longrightarrow \mathbb{R}^p$ be an arbitrary, deterministic, $p$-dimensional function, where $\mathcal{D}$ represents the space of datasets. Its $\ell_2$ sensitivity is defined as

$$\Delta_2 \mu = \sup_{D \text{ adjacent } D'} \| \mu(D) - \mu(D') \|_2$$

**DEFINITION** The Gaussian mechanism $M : \mathcal{D} \longrightarrow \mathbb{R}^p$ with parameter $\rho$ is the randomized mechanism that adds noise $N(0, \rho^2)$ to each of the $p$ components of the output of $\mu$, i.e.,

$$M(\cdot) := \mu(\cdot) + \rho \, N(0, I_p)$$

The privacy guarantees for the Gaussian mechanism are given by the following result.

**THEOREM** [THM A.1, Dwork, Roth, 2014] For every $\varepsilon \in (0,1)$ and $\delta > 0$, the Gaussian mechanism $M$ with parameter

$$\rho \geqslant \frac{\sqrt{2 \log(1.25/\delta)} \; \Delta_2 \mu}{\varepsilon}$$

is $(\varepsilon, \delta)$-differentially private.

In our setting,

$$\mu_{\theta^{t-1}, D} = \theta^{t-1} - \frac{\eta}{n} \sum_{(x_i, y_i \in D)} \nabla_\theta \ell\big(y_i, f(x_i; \theta^{t-1})\big) \Big/ \max\left(1, \frac{\| \nabla_\theta \ell(y_i, f(x_i; \theta^{t-1}) \|_2}{C_{\text{clip}}}\right)$$

Let us now compute the sensitivity.

$$\Delta_2 \mu = \sup_{\theta^{t-1} \in \mathbb{R}^p, \ D \text{ adjacent with } D'} \| \mu_{\theta^{t-1}, D} - \mu_{\theta^{t-1}, D'} \|_2$$

$$= \sup_{\substack{\theta^{t-1} \in \mathbb{R}^p \\ D \text{ adjacent } D'}} \frac{\eta}{n} \Big\| \sum_{(x_i, y_i \in D)} \nabla_\theta \ell(y_i, f(x_i; \theta^{t-1})) / \max\left(1, \frac{\|\nabla_\theta \ell(y_i, f(x_i; \theta^{t-1}))\|_2}{C_{clip}}\right)$$

$$- \sum_{(x_i, y_i \in D')} \nabla_\theta \ell(y_i, f(x_i; \theta^{t-1})) / \max\left(1, \frac{\|\nabla_\theta \ell(y_i, f(x_i; \theta^{t-1}))\|_2}{C_{clip}}\right) \Big\|_2$$

$$= \frac{\eta}{n} \sup_{\theta^{t-1} \in \mathbb{R}^p, (x,y), (x',y')} \Big\| \nabla_\theta \ell(y, f(x; \theta^{t-1})) / \max\left(1, \frac{\|\nabla_\theta \ell(y, f(x; \theta^{t-1}))\|_2}{C_{clip}}\right)$$

$$- \nabla_\theta \ell(y', f(x'; \theta^{t-1})) / \max\left(1, \frac{\|\nabla_\theta \ell(y', f(x'; \theta^{t-1}))\|_2}{C_{clip}}\right) \Big\|_2$$

$$\leq \frac{2\eta}{n} \sup_{\theta^{t-1} \in \mathbb{R}^p, (x,y)} \Big\| \nabla_\theta \ell(y, f(x; \theta^{t-1})) / \max\left(1, \frac{\|\nabla_\theta \ell(y, f(x; \theta^{t-1}))\|_2}{C_{clip}}\right) \Big\|_2$$

$$\leq \frac{2\eta}{n} C_{clip}.$$

For 1 step,  $\rho \gtrsim \frac{\sqrt{\log 1/\delta}}{\varepsilon} \quad \frac{C_{clip} \ \eta}{\cdot n}$

DP - GD  enforced  $\sqrt{\eta} \ \frac{2 C_{clip}}{n} \ \sigma \gtrsim \sqrt{\eta} \ \frac{C_{clip}}{n} \ \frac{\sqrt{\log 1/\delta} \ \sqrt{\eta T}}{\varepsilon}$

This matches for one step !

Now, we need to compose $T$ independent Gaussian mechanisms.

ATTEMPT 1 : Naive composition

THEOREM [THM 3.16, Dwork, Roth, 2014] For $t \in \{1, \ldots, T\}$, let $M_t$ be

$(\varepsilon', \delta')$-differentially private. Then, their composition is $(T\varepsilon', T\delta')$-differentially

private.

This would give $\qquad \sigma \gtrsim \sqrt{\eta} \dfrac{\sqrt{\log(T/\delta)}}{\varepsilon} T$

$\ddot{\frown}$ Too much! We want $\qquad \sigma \gtrsim \sqrt{\eta T} \dfrac{\sqrt{8\log(1/\delta)}}{\varepsilon}$

ATTEMPT 2: Advanced composition

THEOREM [THM 3.20, Dwork, Roth, 2014] For $t \in \{1, \ldots, T\}$, let $M_t$ be

$(\varepsilon', \delta')$-differentially private. Then, their composition is $(\varepsilon, T\delta' + \delta)$-differentially

private, with $\quad \varepsilon = \sqrt{2T \log(1/\delta)}\, \varepsilon' + T\varepsilon'(e^{\varepsilon'} - 1)$

This requires $\delta' < 1/T$, which would give at least

$$\sigma \gtrsim \sqrt{\eta T} \,\boxed{\sqrt{\log T}}$$

$\ddot{\frown}$ also too much since we will look at the

limit $T \longrightarrow +\infty$, $\eta \longrightarrow 0$ and $\eta T \longrightarrow t$

# ATTEMPT 3 : Moment accountant

<u>DEFINITION</u>   For adjacent datasets $D, D' \in \mathcal{D}$, a randomized mechanism $M_t : \mathbb{R}^P \times \mathcal{D} \longrightarrow \mathbb{R}^P$, auxiliary input $\theta^{t-1} \in \mathbb{R}^P$, the <u>privacy loss</u> at the output $\theta$ is defined as

$$\gamma(\theta ; M_t, \theta^{t-1}, D, D') = \log \frac{p(M_t(\theta^{t-1}, D) = \theta)}{p(M_t(\theta^{t-1}, D') = \theta)}$$

probability density function

We also define :

$$\alpha_{M_t}(\lambda ; \theta^{t-1}, D, D') = \log \mathbb{E}_{\theta \sim M_t(\theta^{t-1}, D)} \left[ \exp(\lambda \gamma(\theta ; M_t, \theta^{t-1}, D, D')) \right]$$

log of moment-generating function of the privacy loss evaluated at $\lambda$

$$\alpha_{M_t}(\lambda) = \sup_{\theta^{t-1} \in \mathbb{R}^P, \, D \text{ adjacent to } D'} \alpha_{M_t}(\lambda ; \theta^{t-1}, D, D')$$

supremum over all possible $\theta^{t-1}$ and adjacent datasets $D, D'$

<u>THEOREM</u>  [Abadi, Chu, Goodfellow, 2016]  Let $A$ consist of a sequence of independent mechanisms $M_1, \dots, M_T$. Then,

① <u>Composability</u>. For any $\lambda$,

$$\alpha_A(\lambda) \leq \sum_{t=1}^{T} \alpha_{M_t}(\lambda).$$

② <u>Tail bound</u>. For any $\varepsilon > 0$, $A$ is $(\varepsilon, \delta)$-differentially private for

$$\delta = \inf_{\lambda} \exp(\alpha_A(\lambda) - \lambda \varepsilon)$$

*) Moment accountant exploits the independence of the mechanisms while (advanced) composition allows for arbitrary correlations between them.

In our setting,

$$M_t \sim \mathcal{N}\left(\mu_{\theta^{t-1}, \Delta}, \rho^2\right) \qquad \text{with} \qquad \rho = \sqrt{\eta} \; \frac{2\, C_{\text{clip}}}{n} \, \sigma.$$

Thus,

$$\gamma(\theta; M_t, \theta^{t-1}, \Delta, \Delta') = \log \frac{p(M_t(\theta^{t-1}, \Delta) = \theta)}{p(M_t(\theta^{t-1}, \Delta') = \theta)}$$

$$= -\frac{1}{2\rho^2}\left( \|\theta - \mu_{\theta^{t-1}, \Delta}\|_2^2 - \|\theta - \mu_{\theta^{t-1}, \Delta'}\|_2^2 \right)$$

$$= -\frac{1}{2\rho^2}\left( 2\theta^T(\mu_{\theta^{t-1}, \Delta'} - \mu_{\theta^{t-1}, \Delta}) + \|\mu_{\theta^{t-1}, \Delta}\|_2^2 - \|\mu_{\theta^{t-1}, \Delta'}\|_2^2 \right)$$

$$= -\frac{1}{2\rho^2}\left( 2(\theta - \mu_{\theta^{t-1}, \Delta})^T(\mu_{\theta^{t-1}, \Delta'} - \mu_{\theta^{t-1}, \Delta}) - \|\mu_{\theta^{t-1}, \Delta}\|_2^2 \right.$$

$$\left. + 2\mu_{\theta^{t-1}, \Delta}^T \mu_{\theta^{t-1}, \Delta'} - \|\mu_{\theta^{t-1}, \Delta'}\|_2^2 \right)$$

$$= -\frac{1}{2\rho^2}\left( 2(\theta - \mu_{\theta^{t-1}, \Delta})^T \Delta_{\theta^{t-1}, \Delta, \Delta'} - \|\Delta_{\theta^{t-1}, \Delta, \Delta'}\|_2^2 \right)$$

$$\Delta_{\theta^{t-1}, \Delta, \Delta'} = \mu_{\theta^{t-1}, \Delta'} - \mu_{\theta^{t-1}, \Delta}$$

$$\alpha_{M_t}(\lambda; \theta^{t-1}, \Delta, \Delta') = \log \mathop{\mathbb{E}}_{\theta \sim M_t(\theta^{t-1}, \Delta)}\left[ \exp\left( \lambda \gamma(\theta; M_t, \theta^{t-1}, \Delta, \Delta') \right) \right]$$

$$= \frac{\|\Delta_{\theta^{t-1}, \Delta, \Delta'}\|_2^2}{2\rho^2}(\lambda + \lambda^2)$$

some calculations ...

$$\alpha_{M_t}(\lambda) = \sup_{\theta^{t-1} \in \mathbb{R}^p, \; \Delta \text{ adjacent to } \Delta'} \alpha_{M_t}(\lambda; \theta^{t-1}, \Delta, \Delta')$$

$$= \frac{(\Delta_2 \mu)^2}{2\rho^2}(\lambda + \lambda^2) \leq \frac{\eta}{2\sigma^2}(\lambda + \lambda^2)$$

$$\Delta_2 \mu \leq \frac{2\eta \, C_{clip}}{n}$$

$$\rho = \sqrt{\eta} \; \frac{2 C_{clip}}{n} \sigma$$

so $\dfrac{\Delta_2 \mu}{\rho} \leq \sqrt{\eta}/\sigma$

Let $A$ be the output of DP-GD after $T$ steps. Then, by composability,

$$\alpha_A(\lambda) \leq \frac{\eta T}{2\sigma^2}(\lambda + \lambda^2).$$

Thus,

$$\exp(\alpha_A(\lambda) - \lambda \varepsilon) \leq \exp\left(\frac{\eta T}{2\sigma^2}(\lambda + \lambda^2) - \lambda \varepsilon\right)$$

$$\leq \exp\left(\frac{\varepsilon^2}{16 \log(1/\delta)}(\lambda + \lambda^2) - \lambda \varepsilon\right)$$

$\sigma \geq \sqrt{\eta T} \; \dfrac{\sqrt{8 \log(1/\delta)}}{\varepsilon}$

$$= \exp\left(\frac{\varepsilon^2}{16 \log(1/\delta)}\lambda^2 - \left(1 - \frac{\varepsilon}{16 \log(1/\delta)}\right)\lambda \varepsilon\right)$$

$\Downarrow$

$\dfrac{\eta T}{2\sigma^2} \leq \dfrac{\varepsilon^2}{16 \log(1/\delta)}$

$$\leq \exp\left(\frac{\varepsilon^2}{16 \log(1/\delta)}\lambda^2 - \frac{\lambda \varepsilon}{2}\right)$$

$\varepsilon \in (0, 8 \log(1/\delta))$

$$\inf_{\lambda} \; \exp\left( \alpha_{\mathcal{A}}(\lambda) - \lambda \varepsilon \right) \leq \exp\left( \frac{\varepsilon^2}{16 \log(1/\delta)} \lambda_*^2 - \frac{\lambda_* \varepsilon}{2} \right)$$

$$= \exp\left( \log 1/\delta - 2\log(1/\delta) \right) = \delta$$

pick $\lambda_* = 4 \log(1/\delta)/\varepsilon$

By the tail bound, DP-GD is $(\varepsilon, \delta)$ − differentially private and the proof of the proposition is complete.

# PRIVACY FOR FREE IN THE OVERPARAMETERIZED REGIME

## RANDOM FEATURES (RF)

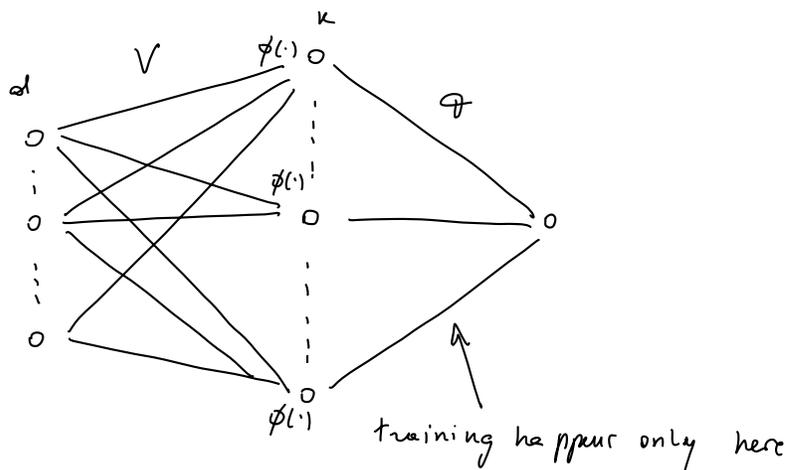$$f_{RF}(x, \vartheta) = \Phi_{RF}^T(x) \vartheta \quad , \quad \Phi_{RF}(x) = \phi(Vx)$$

*) $\quad V \in \mathbb{R}^{k \times d} \qquad V_{ij} \overset{iid}{\sim} N(0, 1/d)$

*) $\quad \vartheta \in \mathbb{R}^p$, $p = k$, vector of trainable parameters

$$\vartheta^* = \Phi_{RF,n}^T \left( \Phi_{RF,n} \Phi_{RF,n}^T \right)^{-1} y \quad , \quad \text{with} \quad \Phi_{RF,n} = \begin{bmatrix} \Phi_{RF}(x_1) \\ \vdots \\ \Phi_{RF}(x_n) \end{bmatrix} = \phi(XV^T) \in \mathbb{R}^{n \times k}$$

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n \times d}$$

obtained by running GD from initialization $\vartheta_0 = 0$



training happens only here

## Data assumptions:

(a) $\quad \int \|x\|_2 \, dP_x(x) = \sqrt{d} \quad \longleftarrow \quad$ scaling

(b) $\quad \|x\|_{\psi_2} = O(1) \quad \longleftarrow \quad$ sub Gaussian data (with subGaussian norm bdd by universal constant)

(c) $\quad \lambda_{min} \left( \underset{x \sim P_x}{\mathbb{E}}[xx^T] \right) = \Omega(1) \quad \longleftarrow \quad$ well conditioned covariance

(d) $\quad |y| \leq C \quad \longleftarrow \quad$ bounded labels

Square loss: $\quad \ell(y, \hat{y}) = (y - \hat{y})^2$

THEOREM [Bombari, M., 2015] Let $\phi(\cdot)$ be nonlinear, Lipschitz such that

$\phi_0 = \phi_2 = 0$, $\phi_1 \neq 0$. Let $n = O(\sqrt{p})$, $n = \hat{\Omega}(d)$, $n = \tilde{o}(d^{3/2})$ and

consider a privacy budget $\delta \in (0,1)$, $\varepsilon \in \left(0, 8\log(1/\delta)\right)$, $\dfrac{\varepsilon}{\sqrt{\log 1/\delta}} \gg \dfrac{d}{n}$ .

Then, by setting properly $C_{clip}$, $\sigma$, $T$, we have that $\theta^T$ is

$(\varepsilon, \delta)$ - differentially private and

$$R_E = \tilde{O}\left( \frac{d}{n\varepsilon} + \sqrt{\frac{d}{n}} + \sqrt{\frac{n}{d^{3/2}}} \right) .$$

<u>INTERPRETATION</u>   If $d \ll n \ll d^{3/2}$, then privacy comes for free (in the

sense that the excess population risk $R_E$ is $o(1)$) as long as $\varepsilon \gg \dfrac{d}{n}$, for

any degree of overparameterization.

$\longrightarrow$   $\varepsilon = \Theta(1)$ usual in practice. Here we can even guarantee the

strong privacy requirement $\varepsilon = o(1)$

$\longrightarrow$   dependence on $\delta$ only logarithmic (and hidden in $\tilde{O}(\cdot)$)

$\longrightarrow$   no dependence on $p$ as long as $p \gtrsim n^2$. We expect this

can be relaxed to $p \gtrsim n$ (which is the number of parameters

to interpolate as seen before)

$\longrightarrow$   $d \ll n \ll d^{3/2}$ corresponds to standard datasets, such as CIFAR-10

($n = 5 \cdot 10^4$, $d \approx 3 \cdot 10^3$) or ImageNet ($n \approx 1.3 \cdot 10^6$, $d \approx 9 \cdot 10^4$).

We also expect that it can be relaxed to $d \ll n \ll d^2$ and, in fact,

to $d^\ell \ll n \ll d^{\ell+1}$, as explained later.

$\longrightarrow$   $\phi_0 = \phi_2 = 0$ can again be potentially relaxed and role of $\phi_1 \neq 0$
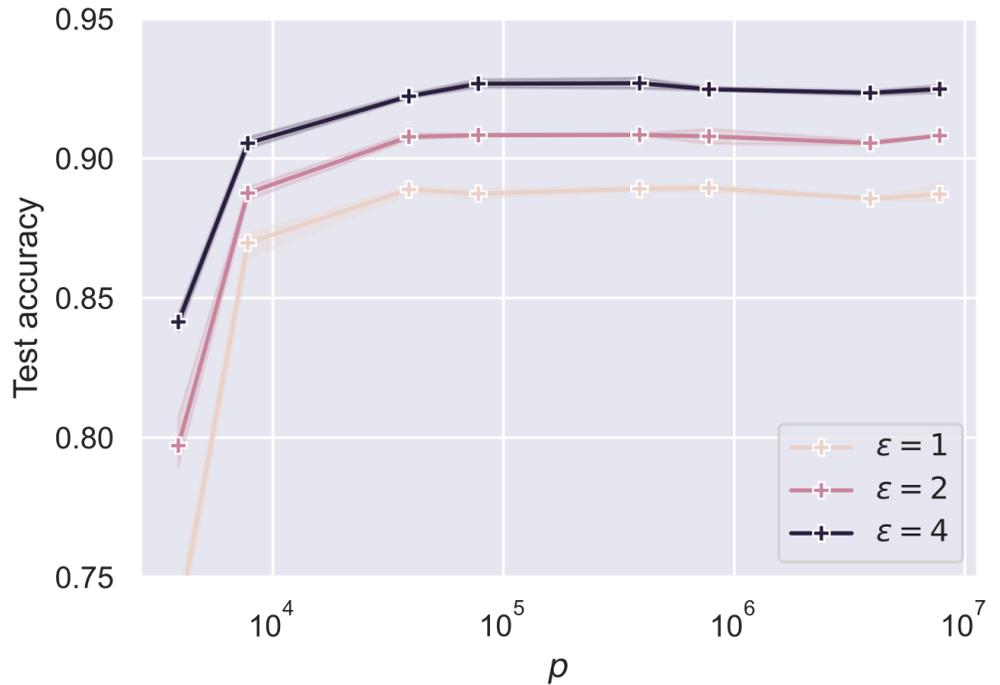
clarified later.

$\longrightarrow$   Proof describes correct scaling of hyperparameters $C_{clip}$, $\sigma$, $T$

# IDEA OF THE ARGUMENT IN TWO CARTOONS

CARTOON 1 : Test loss of DP-GD on MNIST for two-layer networks

(a)  n fixed

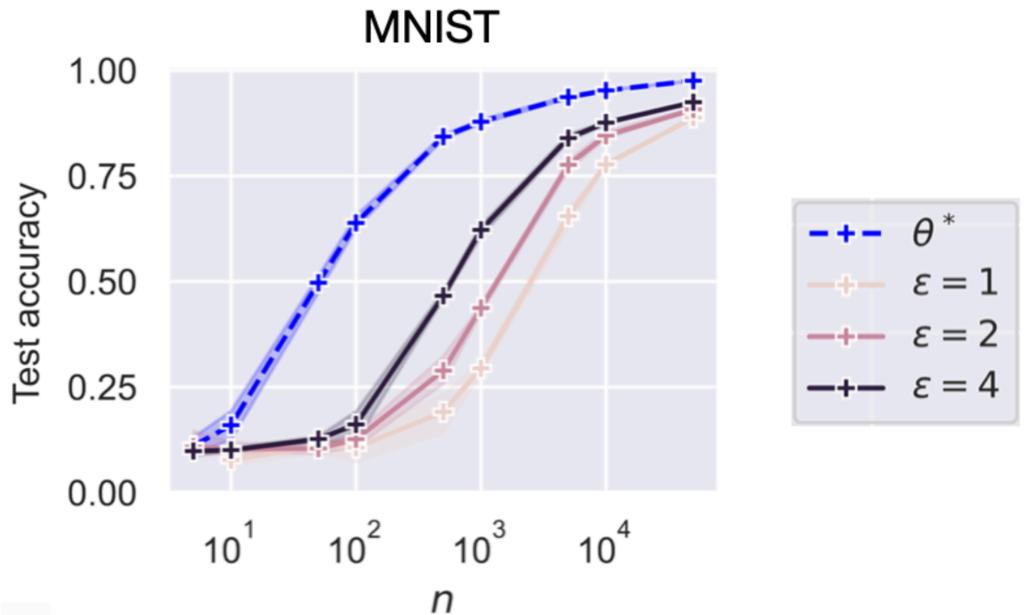Test accuracy
as a function of p



Test accuracy increases until the network is wide enough and then plateaus.

(b)  p fixed

Test accuracy
as a function of n



As # samples n grows, gap between DP-GD and GD shrinks

CARTOON 2 : Test loss of (non-private) GD on RF regression.

Before the actual cartoon, let us start with a result from the review

by [ Misiakiewicz, Montanari, 2024 ]

THEOREM Let $x_i \sim \text{Unif}\left(\mathbb{S}^{d-1}(\sqrt{d})\right)$, $y_i = f_*(x_i) + \varepsilon_i$, $f_* \in L^2$, $\varepsilon_i \perp x_i$,

$\mathbb{E}[\varepsilon_i] = 0$, $\mathbb{E}[\varepsilon_i^2] = \tau^2$. Assume $d^{\ell_1 + \delta} \leq n \leq d^{\ell_1 + 1 - \delta}$, $d^{\ell_2 + \delta} \leq p \leq d^{\ell_2 + 1 - \delta}$,

$\max\left(p/n, n/p\right) \geq d^\delta$ for some integers $\ell_1, \ell_2$ and constant $\delta > 0$. Denote

$\ell = \min(\ell_1, \ell_2)$. Furthermore assume that $\phi(\cdot)$ has all Hermite coefficients

bounded away from 0. Then,

$$\mathbb{E}_x \left| f_*(x) - f_{RF}(x, \theta^*) \right|^2 = \left\| P_{>\ell} f_* \right\|_{L^2}^2 + o(1)\left(\left\| f_* \right\|_{L^{2,1}}^2 + \tau^2\right)$$

$P_{>\ell} = I - P_{\leq \ell}$, with $P_{\leq \ell} : L^2 \to L^2$ the orthogonal projector onto the subspace

of polynomials of degree at most $\ell$

( project $f_*$ into the space of degree-$\ell$ polynomials and take the $L^2$ norm of the residual )

we expect the result on DP-GD to hold under these milder conditions

$p \gtrsim n$ and $d \ll n \ll d^2$ correspond to $\ell_2 \geq \ell_1 = 1$ which gives

$\ell = \min(\ell_1, \ell_2) = 1$. Then, the result above says that the test error

of the RF model is flat as a function of $n$ in the whole regime $d \ll n \ll d^2$

( when $p \gtrsim n$ ).

Loss plateaus for $d \ll n \ll d^2 \implies \Theta(d)$ samples used to achieve utility and

the surplus to achieve privacy.

More generally :

If $p \gg n$, then the test error of the RF model is well approximated
by the test error of the corresponding $p = \infty$ kernel method
[ Mei, Misiakiewicz, Montanari, 2022 ].

Formally, recall that
$$
\Phi_{RF,n} = \begin{bmatrix} \Phi_{RF}(x_1) \\ \vdots \\ \Phi_{RF}(x_n) \end{bmatrix} = \begin{bmatrix} \phi^T(Vx_1) \\ \vdots \\ \phi^T(Vx_n) \end{bmatrix} \in \mathbb{R}^{n \times N} , \text{ so that}
$$

$$
\hat{f}_{RF}(x, \theta^*) = (\phi(Vx))^T \Phi_{RF,n}^T \left( \Phi_{RF,n} \Phi_{RF,n}^T \right)^{-1} y
$$

$$
= \left[ \langle \phi(Vx), \phi(Vx_1) \rangle, \cdots, \langle \phi(Vx), \phi(Vx_n) \rangle \right] K_{RF,n}^{-1} y
$$

with $\quad (K_{RF,n})_{ij} = \langle \phi(Vx_i), \phi(Vx_j) \rangle$.

The corresponding $p = \infty$ kernel estimator is :

$$
\hat{f}_{RF, p=\infty}(x) = \left[ \underset{v}{\mathbb{E}}\left[ \phi(v^T x) \cdot \phi(v^T x_1) \right], \cdots, \underset{v}{\mathbb{E}}\left[ \phi(v^T x) \cdot \phi(v^T x_n) \right] \right] K_n^{-1} y
$$

with $\quad (K_n)_{ij} = \underset{v}{\mathbb{E}}\left[ \phi(v^T x_i) \cdot \phi(v^T x_j) \right]$

[ Mei, Misiakiewicz, Montanari, 2022 ] show that, when $p \gtrsim n$,

$$
R_{RF} := \underset{x}{\mathbb{E}} \left| f_*(x) - \hat{f}_{RF}(x, \theta^*) \right|^2 \approx \underset{x}{\mathbb{E}} \left| f_*(x) - \hat{f}_{RF, p=\infty}(x) \right|^2 =: R_{RF, p=\infty}
$$

this is the test loss of a kernel ridge(less) regression estimator

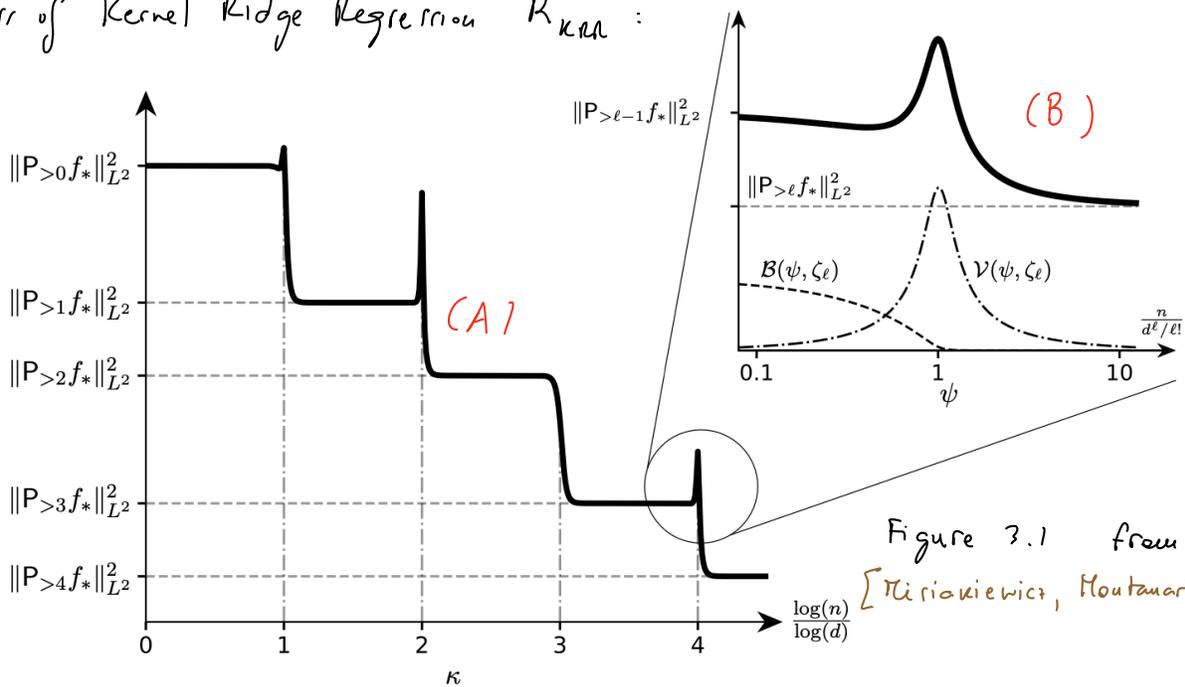Test loss of Kernel Ridge Regression $R_{KRR}$ :



Figure 3.1 from [Misiakiewicz, Montanari, 2024]

(A) If $d^{\ell+\delta} \leq n \leq d^{\ell+1-\delta}$ for some integer $\ell$ and constant $\delta > 0$, then

$$R_{KRR} = \| P_{>\ell} f_* \|_{L^2}^2 + o(1) \left( \| f_* \|_{L^2}^2 + \tau^2 \right)$$

we expect the result on DP-GD to hold under this milder condition

If $d^{\ell} \ll n \ll d^{\ell+1}$, then the model fits the best degree-$\ell$ approximation of the target function.

$\longrightarrow$ Proved first by [Ghorbani, Mei, Misiakiewicz, Montanari, 2021] and generalized to any RKHS under a spectral gap assumption in [Mei, Misiakiewicz, Montanari, 2022].

(b) If $\dfrac{n}{d^{\ell} \ell!} \longrightarrow \psi$ for some integer $\ell$ and constant $\psi > 0$, then

$$R_{KRR} = \| P_{\ell} f_* \|_{L^2}^2 \, \mathcal{B}(\psi) + \left( \| P_{>\ell} f_* \|_{L^2}^2 + \tau^2 \right) \mathcal{V}(\psi) + \| P_{>\ell} f_* \|_{L^2}^2$$
$$+ o(1) \left( \| f_* \|_{L^2}^2 + \tau^2 \right)$$

$P_{\ell} = P_{\leq \ell} \, P_{> \ell - 1}$

$\longrightarrow$ Proved by [Xiao, Hu, Misiakiewicz, Lu, Pennington, 2022]

Proof of this result based on diagonalization of inner-product kernels on the sphere and expansion in spherical harmonics.

To analyze DP-GD, need non-asymptotic control (in $n, d, p$) of the whole trajectory of the algorithm. We will next discuss the outline of the proof.
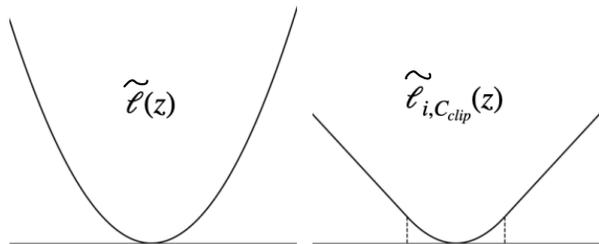
# PROOF OUTLINE

DP-GD :

$$\theta^{t+1} = \theta^t - \eta \; \frac{1}{n} \sum_{i=1}^{n} \nabla_\theta \, \ell(\, y_i \, , \, f(x_i \, ; \theta^t\, ))\Big/ \max\left(1 \, , \, \frac{\|\nabla_\theta \ell(y_i \, , f(x_i \, ; \theta^t))\|_2}{C_{clip}}\right)$$

$$+ \sqrt{\eta} \; \frac{2\, C_{clip}}{n} \; \sigma \, N(0, I_p)$$

$$= \theta^t - \eta \nabla_\theta \hat{R}_n^{(clip)}(\theta^t) + \sqrt{\eta} \; \frac{2\, C_{clip}}{n} \; \sigma \, N(0, I_p)$$

$$\hat{R}_n^{(clip)}(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell_{i, C_{clip}}(\, y_i \, , \, f(x_i \, ; \theta))$$  clipped / Huber loss

$$\ell(y_i , f(x_i, \theta)) = \tilde{\ell}(\, y_i - f(x_i, \theta))$$

$$\ell_{i, C_{clip}}(\, y_i \, f(x_i, \theta)) = \tilde{\ell}_{i, C_{clip}}(\, y_i - f(x_i, \theta))$$



$\tilde{\ell}(z)$   $\tilde{\ell}_{i, C_{clip}}(z)$

This is the Euler-Maruyama discretization scheme of the SDE :

$$d\,\Theta(t) = -\nabla \hat{R}_n^{(clip)}(\Theta(t))\, dt + \underbrace{\overbrace{\frac{2\, C_{clip}}{n}}^{:= \Sigma} \sigma \; dB(t)}$$

p-dimensional Wiener process

By the earlier proposition, if $\quad \Sigma \geq \frac{2\, C_{clip}}{n} \sqrt{T} \; \frac{\sqrt{8 \log(1/\delta)}}{\varepsilon}$ , then

$\Theta(T)$ is $(\varepsilon, \delta)$-differentially private. It remains to compute the test

error of $\Theta(T)$.

Consider now the SDE with the usual quadratic loss $\hat{R}_n(\cdot)$ :

$$d\,\hat{\Theta}(t) = - \nabla \hat{R}_n(\Theta(t))\,dt + \Sigma\,dB(t)$$

$$= -\frac{2\,\overline{\Phi}_{nf}^{\mathsf{T}}}{n}\left(\overline{\Phi}_{nf}\,\hat{\Theta}(t) - y\right)dt + \Sigma\,dB(t)$$

This is easier to control, as it is an Ornstein-Uhlenbeck (OU) process.

Let $\mathcal{C} := \{\,\Theta : \;\; \|\nabla_\Theta \ell(y_i,\,f(x_i;\Theta))\| < C_{clip} \quad \forall i\,\}$

( subset of parameter space where clipping does not happen , i.e., $\hat{R}_n(\Theta) = \hat{R}_n^{(clip)}(\Theta)$

If $\hat{\Theta}(t)$ stays in $\mathcal{C}$ , then $\Theta(t) = \hat{\Theta}(t)$ .

STEP 1 : $\hat{\theta}(t) \in \mathcal{C} \quad \forall t \in [0,T]$ with $C_{clip} = \sqrt{p} \log^2 n$ and $T = \frac{d}{p} \log^2 n$

*) Computing $\nabla_\theta \ell(\cdot, \cdot)$, we can express $\mathcal{C}$ as

$$\mathcal{C} = \left\{ \theta : \quad |\phi^T(Vx_i)\theta - y_i| < \frac{C_{clip}}{2\|\phi(Vx_i)\|} \quad \forall i \right\}$$

*) $\hat{\theta}(t) = \underbrace{\mathbb{E}_B\left[\hat{\theta}(t)\right]}_{:= \bar{\theta}(t) \text{ expectation over the OU process giving the usual gradient flow}} + \tilde{\theta}(t)$

We show that, with high probability (at least $1 - \exp(-\Omega(\log^2 n)) - \exp(-\Omega(p))$)

(1.A) $\quad \sup_{t \in [0,T]} |\phi^T(Vx_i)\bar{\theta}(t) - y_i| = O(\log n)$

(1.B) $\quad \sup_{t \in [0,T]} |\phi^T(Vx_i)\tilde{\theta}(t)| = O(\log n)$

(1.A) + (1.B) give the desired claim as $\frac{C_{clip}}{2\|\phi(Vx_i)\|} = \Omega(\log^2 n)$

PROOF OF (1.A) VIA LEAVE-ONE-OUT

*) $\quad \hat{\theta}(t) = \left(1 - e^{-2\,\Phi_{nF}^T \Phi_{nF} \frac{1}{n} t}\right) \underbrace{\Phi_{nF}^+}_{\text{Moore-Penrose inverse}} y \qquad (\hat{\theta}(0) = \dot{\hat{\theta}}(0) = 0)$

$$\Downarrow$$

$$y_i - \phi^T(Vx_i)\,\hat{\theta}(t) = \phi^T(Vx_i)\, e^{-2\,\Phi_{nF}^T \Phi_{nF} \frac{1}{n} t}\, \Phi_{nF}^+ y$$

Consider the leave-one-out quantities $\Phi_{nF,-i}$ and $y_{-i}$ obtained respectively from $\Phi_{nF}$ and $y$ after removing the $i$-th sample.

*) $\quad \left\| \Phi_{nF}^+ y - \Phi_{nF,-i}^+ y_{-i} \right\| = \tilde{O}\left(\frac{1}{\sqrt{p}}\right)$

GD sensitivity + lower bound on $\lambda_{min}\left(\Phi_{nF}\Phi_{nF}^T\right)$

*) Leave-one-out at the exponent:

$$\sup_{t \in [0,T]} \left| \phi^T(Vx_i)\, e^{-2\,\Phi_{nF}^T \Phi_{nF} \frac{1}{n} t}\, \Phi_{nF,-i}^+ y_{-i} \right|$$

$$\leq 2 \sup_{t \in [0,T]} \left| \phi^T(Vx_i)\, e^{-2\,\Phi_{nF,-i}^T \Phi_{nF,-i} \frac{1}{n} t}\, \Phi_{nF,-i}^+ y_{-i} \right|$$

By Lie's product formula

$$\left| \phi^T(Vx_i)\, e^{-2\,\Phi_{nF}^T \Phi_{nF} \frac{1}{n} t}\, \Phi_{nF,-i}^+ y_{-i} \right|$$

$$= \lim_{s \to +\infty} \left| \phi^T(Vx_i) \underbrace{\left(e^{-\frac{2\phi(Vx_i)\phi^T(Vx_i)}{ns} t}\, e^{-2\,\Phi_{nF,-i}^T \Phi_{nF,-i} \frac{1}{ns} t}\right)^s}_{(:= \Pi(s)}\, \Phi_{nF,-i}^+ y_{-i} \right|$$

Expand the $r$-th power in $\quad \Pi(r) = \left( \left( I + \alpha(r) \dfrac{\phi(Vx_i)\,\phi^{\top}(Vx_i)}{\|\phi(Vx_i)\|^2} \right) A(S) \right)^{s}$

$$A(r) = e^{-2\,\Phi_{\mathrm{aF},-i}^{\top}\,\Phi_{\mathrm{aF},-i}\,\frac{1}{ns}t}, \qquad \alpha(r) = -1 + e^{-\frac{2\|\phi(Vx_i)\|^2}{ns}t}$$

*) Bound $\quad \sup_{t \in [0,T]} \left| \underbrace{\phi(Vx_i)}\; \underbrace{e^{-2\,\Phi_{\mathrm{aF},-i}^{\top}\,\Phi_{\mathrm{aF},-i}\,\frac{1}{n}t}\; \Phi_{\mathrm{aF},-i}^{+}\, y_{-i}} \right|$

these two pieces are independent now!

by Dudley's chaining tail inequality

Upper bound the $\epsilon$-covering number $N(\epsilon, \widetilde{T})$ of the set $\widetilde{T} \subseteq \mathbb{R}^d$ described

by the curve $\gamma(t) = V^{\top} e^{-2\,\Phi_{\mathrm{aF},-i}^{\top}\,\Phi_{\mathrm{aF},-i}\,\frac{1}{n}t}\; \Phi_{\mathrm{aF},-i}^{+}\, y_{-i}$ :

$$\int_0^{\infty} \sqrt{\log N(\widetilde{T}, \epsilon)}\; d\epsilon = O(\log n), \qquad \operatorname{diam}(\widetilde{T}) = O(1)$$

# PROOF OF (1.B) BY COMPARISON INEQUALITIES

Consider the auxiliary process $dz_i(t) = \phi^T(Vx_i) \Sigma \, dB(t)$.

This is obtained by removing the attractive drift $-2 \widetilde{\Phi}_{RF}^T ( \widehat{\Phi}_{RF} \widehat{\Theta}(t) - y)/n$

to the SDE, and it is easier to analyze (a Wiener process).

$$\sup_{t \in [0,T]} |\phi^T(Vx_i) \widetilde{\Theta}(t)| \leq \mathbb{E}_B \left[ \sup_{t \in [0,T]} |\phi^T(Vx_i) \widetilde{\Theta}(t)| \right] + \log n$$

Borell-TIS + variance = $O(1)$

$$\leq \mathbb{E}_{z_i} \left[ \sup_{t \in [0,T]} |z_i(t)| \right] + \log n$$

Sudakov-Fernique

( removing the drift increases variance, so process less concentrated around the mean )

$$\leq O(1) + \log n$$

$$\mathbb{E}_{z_i} \left[ \sup_{t \in [0,T]} |z_i(t)| \right] = O(1) \quad \text{since the variance of the Wiener process}$$

is $\quad \Sigma^2 T \|\phi(Vx_i)\|^2 = O(1)$.

STEP 2 : Control noise and early stopping in $\hat{\Theta}(t)$

$$\hat{\Theta}(t) = \hat{\theta}(t) + \tilde{\Theta}(t) = \theta^* + (\underbrace{\hat{\theta}(t) - \theta^*}_{\text{early stopping}}) + \underbrace{\tilde{\Theta}(t)}_{\text{noise}}$$

(non-private) GD solution

(2.A) $\quad \mathbb{E}_X\left[\left(\underbrace{\phi^T(V_x)\,\tilde{\Theta}(T)}\right)^2\right] = \tilde{O}\left(\dfrac{d^2}{\varepsilon^2 n^2}\right)$

Gaussian with variance proportional to $\|\phi(V_x)\|^2$, $T$, $\Sigma^2$
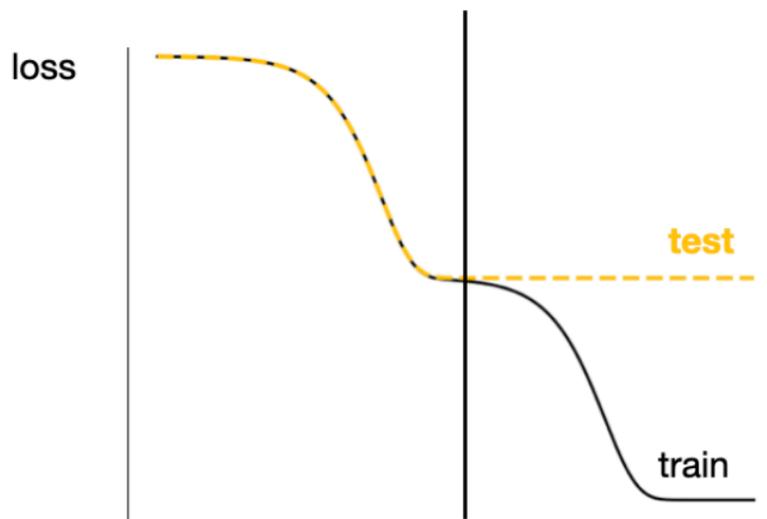so claim follows from choice of hyperparameters

(2.B) $\quad \mathbb{E}_X\left[\left(\phi^T(V_x)(\hat{\theta}(T) - \theta^*)\right)^2\right] = \tilde{O}\left(\dfrac{d}{n} + \dfrac{n}{d^{3/2}}\right)$

Gap between $d$-th and $(d+1)$-th eigenvalue of $\Phi_{nf}\Phi_{nf}^T$

$$\left(\lambda_d(\Phi_{nf}\Phi_{nf}^T) = \Omega\left(p\,\dfrac{n}{d}\right), \quad \lambda_{d+1}(\Phi_{nf}\Phi_{nf}^T) = O(p)\right)$$

$\Downarrow$

$T = \dfrac{d}{p}\log^2 n$ is enough time to reach the GD plateau and more time would not help anyway resulting only in overfitting (reduction in training error without improvement in test error)

Formally, we decompose:

$$\phi^T(V_x)\left(\hat{\theta}(T) - \theta^*\right) = \phi^T(V_x)\left(P_\Lambda + P_\Lambda^\perp\right)\left(\hat{\theta}(T) - \theta^*\right)$$

where $P_\Lambda$ is the projector on the space spanned by the eigenvectors associated to the $d$ largest eigenvalues of $\Phi_{RF}\Phi_{RF}^T$.

*1 $\left\| P_\Lambda\left(\hat{\theta}(T) - \theta^*\right)\right\|$ negligible

in that subspace $\hat{\theta}(T)$ already close to convergence despite early stopping

*1 $\mathbb{E}_x\left[\left(\phi^T(V_x)\ P_\Lambda^\perp\left(\hat{\theta}(T) - \theta^*\right)\right)^2\right]$

$$\leq 2\mathbb{E}_x\left[\left((V_x)^T P_\Lambda^\perp\left(\hat{\theta}(T) - \theta^*\right)\right)^2\right] + 2\mathbb{E}_x\left[\left(\tilde{\phi}^T(V_x)\ P_\Lambda^\perp\left(\hat{\theta}(T) - \theta^*\right)\right)^2\right]$$

$\phi(z) = z + \tilde{\phi}(z)$

$\tilde{O}\left(d/n + n/d^{3/2}\right)$ by bounding $\left\| V^T P_\Lambda^\perp \Phi_{RF}^+\right\|$

$\tilde{O}\left(d/n + n/d^{3/2}\right)$ by bounding $\left\|\mathbb{E}_x\left[\tilde{\phi}(V_x)\tilde{\phi}^T(V_x)\right]\right\|_{op}$

(2.A) + (2.B) give that $\left|\hat{R} - R^*\right| = \tilde{O}\left(\dfrac{d}{n\varepsilon} + \sqrt{\dfrac{d}{n}} + \sqrt{\dfrac{n}{d^{3/2}}}\right)$

test error of $\theta^*$

test error of $\hat{\Box}(T)$

Step 1 gives that $\hat{\Box}(T) = \Box(T)$ concluding the proof.

⌐ EXTRA :

Sudakov-Fernique's inequality

THEOREM Let $(X_t)_{t \in T}$ and $(Y_t)_{t \in T}$ be two mean zero Gaussian processes. Assume that for all $s, t \in T$, $\mathbb{E}(X_t - X_s)^2 \leq \mathbb{E}(Y_t - Y_s)^2$. Then,

$$\mathbb{E} \sup_{t \in T} X_t \leq \mathbb{E} \sup_{t \in T} Y_t$$

Tsirelson, Ibragimov, Sudakov

Borell -TIS inequality

THEOREM Let $(X_t)_{t \in T}$ be a mean zero Gaussian process which is almost surely finite. Let $\sigma_T^2 := \sup_{t \in T} \mathbb{E} X_t^2$. Then, $\sigma_T^2$ and $\mathbb{E} \sup_{t \in T} |X_t|$ are finite and, for all $u > 0$,

$$\mathbb{P}\left( \sup_{t \in T} |X_t| > \mathbb{E} \sup_{t \in T} |X_t| + u \right) \leq \exp\left( - \frac{u^2}{2\sigma_T^2} \right)$$

# DP-SGD with CLIPPING : SETTING

→ So far, clipping treated as a nuisance ( we set $C_{clip}$ so that clipping never happens with high probability )

→ Empirical evidence that aggressive clipping helps performance

We discuss a sharp analysis that captures the effect of clipping

\*) <u>Linear regression</u>

$$y_i = x_i^T \vartheta^* + z_i \quad , \quad x_i \sim N(0, \Sigma) \quad , \quad z_i \sim N(0, \gamma^2), \quad z_i \perp x_i$$

→ input dimension $d$ = # parameters $p$

→ $Tr(\Sigma) = d \quad , \quad \lambda_{max}(\Sigma) / \lambda_{min}(\Sigma) = \mathcal{O}(1) \quad$ ( well-conditioned covariance )

→ $\|\vartheta^*\| = \mathcal{O}(1)$

\*) <u>Proportional regime</u>

→ $n, d$ large with $d/n \longrightarrow \gamma \in (0, \infty)$

\*) <u>( One-pass ) Differentially-private stochastic gradient descent ( DP-SGD )</u>

$$\vartheta^{n+1} = \vartheta^n - \eta_n \nabla_\vartheta \ell(y_{n+1}, f(x_{n+1}; \vartheta^n)) \Big/ max \left( 1, \frac{\|\nabla_\vartheta \ell(y_{n+1}, f(x_{n+1}; \vartheta^n))\|_2}{C_{clip}} \right)$$
$$+ 2 C_{clip} \sigma_n N(0, I_d)$$

Square loss : $\ell(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$

Linear model : $f(x; \vartheta) = x^T \vartheta$

Test error : $R(\vartheta) := R(f(x; \vartheta)) = \frac{1}{2} \underset{(x,y)}{\mathbb{E}} |y - f(x; \vartheta)|^2 = \frac{1}{2} \underset{x}{\mathbb{E}} |x^T(\vartheta^* - \vartheta) + z|^2$
$$= \frac{1}{2}\left( \|\Sigma^{1/2}(\vartheta^* - \vartheta)\|^2 + \gamma^2 \right)$$

$\longrightarrow$ # samples $n$ = # iterations of the algorithm $T$

$\longrightarrow$ Each data point used only once

$\longrightarrow$ $C_{clip} = c\sqrt{d} \longleftarrow$ constant independent of $n, d$

Before, $C_{clip} = \sqrt{p} \log^2 n = \sqrt{d} \log^2 n$ so that clipping never happens w.h.p.

$( \quad p = d$ in linear regression

Here, $C_{clip} \sim \sqrt{d}$ and $\| \nabla_\theta \ell(y, f(x;\theta)) \| = \|x\| \ |x^T\theta - y| \sim \sqrt{d}$

so that clipping is frequent.

*) <u>zero-concentrated DP</u>   (z CDP)

<u>DEFINITION</u> ($\rho$- z CDP) [Bun, Steinke, 2016]  Given $\alpha \in (1, +\infty)$ and two random variables $X$ and $X'$ with laws $p_X$ and $p_{X'}$, their $\alpha$-Renyi divergence is

$$D_\alpha (X \| X') = \frac{1}{\alpha - 1} \ln \int \left( \frac{p_X(\theta)}{p_{X'}(\theta)} \right)^\alpha p_{X'}(\theta)\, d\theta .$$

Then, a randomized algorithm $A$ satisfies <u>$\rho$-zero concentrated differential</u> <u>privacy</u> if, for any pair of adjacent datasets $\Delta, \Delta'$ and any $\alpha \in (1, \infty)$, we have $D_\alpha (A(\Delta) \| A(\Delta')) \leq \alpha \rho$ .

$\rho^2/2$ - z CDP $\implies$ $(\rho^2/2 + \rho \sqrt{2\ln(1/\delta)}, \delta)$ - DP   [Bun, Steinke, 2016]

$\rho = \Theta(1) \not\Rightarrow \varepsilon = \Theta(1), \quad \delta \ll 1/n$ so $\rho$ - z CDP with $\rho = \Theta(1)$ is

a relaxed notion of privacy

<u>PROPOSITION</u> (following [Feldman, Koren, Talwar, 2020])  The output $\theta_n$

of DP-SGD is $(\rho^2/2)$ - z CDP, where

$$\rho = \max_{k \in \{1, \cdots, n\}} \frac{\eta_k}{\sqrt{\sum_{j=k}^{n} \sigma_j^2}}$$

$\longrightarrow$ Each sample $x_k$ protected by overall noise introduced in next updates $\sum_{j=k}^{n} \sigma_j^2$

$\longrightarrow$ Proof based on privacy amplification by iteration (rather than advanced composition)

$\longrightarrow$ Minimize noise introduced by the algorithm via the schedule :

$$\eta_n = \rho \sqrt{\sum_{j=k}^{n} \sigma_j^2} \iff \rho^2 \sigma_k^2 = \begin{cases} \eta_k^2 - \eta_{k+1}^2 & k \in \{1, \dots, n-1\} \\ \eta_k^2 & k = n \end{cases}$$

Idea : Track the test error of DP-SGD via an SDE in parameter space. This gives a __deterministic equivalent__ of the (stochastic) dynamics since the SDE is equivalent to a system of $d$ coupled ODEs.

__DEFINITION__ (Homogeneized DP-SGD). For any $t \in [0,1)$, we define the homogeneized DP-SGD (H-DP-SGD) as the solution of the SDE :

$$d\Theta_t = -\hat{\eta}(t) \mu_c(\Theta_t) \nabla R(\Theta_t) \, dt + \hat{\eta}(t) \sqrt{\frac{2 \nu_c(\Theta_t) R(\Theta_t) \mathcal{I}}{n}} \, dB_t^s + 2 \frac{\sqrt{d}}{n} c \hat{\sigma}(t) \, dB_t^p$$

$\longrightarrow$ $\Theta_0 = 0$

$\longrightarrow$ $B_t^s$, $B_t^p$ are two independent Wiener processes

$\longrightarrow$ $\hat{\eta}(t)$ is the continuous version of the discrete step size $\eta_k$ $\left(\eta_k = \frac{\hat{\eta}(k/n)}{n}\right)$

(both $\hat{\eta}^2(\cdot)$ and the absolute value of its first and second derivative are required to be uniformly bounded from above by a constant independent of $n, d$ )

$\longrightarrow$ $\hat{\sigma}(t)$ is the continuous version of the noise standard deviation $\sigma_k$ so that $\rho^2 \sigma_k^2 = \eta_k^2 - \eta_{k+1}^2$ becomes $\rho^2 \hat{\sigma}^2(t) = -\frac{d}{dt} \hat{\eta}^2(t)$

$\longrightarrow$ $R(\theta) = \frac{1}{2} \mathbb{E} |x^T \theta - y|^2$ is the risk

$$\longrightarrow \quad \mu_c(\theta) = \frac{\left\| \mathbb{E}_{x,y}\left[ r_c(\theta, x, y)\, x \right] \right\|}{\left\| \mathbb{E}_{x,y}\left[ r(\theta, x, y)\, x \right] \right\|} \qquad \text{is the descent reduction factor}$$

<span style="color:blue">↑ effect of clipping</span>

with $\quad r(\theta, x, y) = x^T \theta - y \quad$ the residual in $\theta$ and

$$r_c(\theta, x, y) = r(\theta, x, y) \,/\, \max\left( 1, \frac{|r(\theta, x, y)|}{c} \right) \quad \text{the clipped residual}$$

$$\longrightarrow \quad \nu_c(\theta) = \frac{\mathbb{E}\left[ r_c^2(\theta, x, y) \right]}{\mathbb{E}\left[ r^2(\theta, x, y) \right]} \qquad \text{is the variance reduction factor}$$

<span style="color:blue">↑ effect of clipping</span>

# DP-SGD WITH CLIPPING: SHARP ANALYSIS

__THEOREM__ [Bombari, Seroussi, M., 2025]  Assume $\rho = \Theta(1)$ and $\sup_{t \in [0,1]} \tilde{\eta}(t) < 2/\gamma$.

Let $\Theta_t$ and $\vartheta_u$ be independent realizations of H-DP-SGD and DP-SGD.

Then, with high probability,

$$\sup_{b \in [0,1)} | R(\Theta_t) - R(\vartheta^{\lfloor tn \rfloor}) | = O\left( \frac{\log^2 n}{\sqrt{n}} \right).$$

**\*)** $\sup_{t \in [0,1]} \tilde{\eta}(t) < 2/\gamma$ is a stability condition for SGD   (can be relaxed)

__INTERPRETATION__   The test error of DP-SGD is well approximated by

effect of clipping

$$d \Theta_t = \underbrace{- \tilde{\eta}(t) \mu_c(\Theta_t) \nabla R(\Theta_t) \, dt}_{(I)} + \underbrace{\tilde{\eta}(t) \sqrt{\frac{2 \nu_c(\Theta_t) R(\Theta_t) \mathcal{I}}{n}} \, dB_t^s}_{(II)} + \underbrace{2 \frac{\sqrt{d}}{n} c \tilde{\sigma}(t) dB_t^p}_{(III)}$$

$(I)$ = term corresponding to descent towards the minimizer of $R(\cdot)$

$(II)$ = term corresponding to noise inherent in SGD

$(III)$ = term corresponding to noise needed to enforce DP

## SDE $\longrightarrow$ System of ODES

SDE equivalent to the following deterministic system of $d$ coupled ODES:

$$dD_i = -2\lambda_i \hat{\eta}(t)\mu_c(R(t))D_i\,dt + \lambda_i \hat{\eta}^2(t)\nu_c(R(t)))(R(t)+\gamma^2/2)\gamma\,dt + 2c^2\hat{\sigma}^2(t)\gamma^2\,dt$$

$$(\bigstar)$$

where $R(t) = \dfrac{1}{d}\sum\limits_{i=1}^{d}\lambda_i D_i(t)$ and $\{\lambda_i\}_{i=1}^{d}$ are the eigenvalues of $\Sigma$,

so that

$$\sup_{t\in[0,1)} \left| R(t) - R^{exc}(\vartheta^{\lfloor tn\rfloor}) \right| = O\left(\frac{\log^2 n}{\sqrt{n}}\right).$$

$$R^{exc}(\vartheta) = R(\vartheta) - \gamma^2/2 \quad \rightsquigarrow \text{ excess error / noiseless test risk}$$

$$\text{such that } \quad R^{exc}(\vartheta^*) = 0.$$

It is easier (and already insightful!) to work with an upper and a lower bound:

$$\underline{R}(t) \;\leq\; R(t) \;\leq\; \overline{R}(t) \;,\quad \text{where } \overline{R}(t), \underline{R}(t): [0,1] \longrightarrow \mathbb{R} \text{ are the}$$

unique solutions to the following two (decoupled) ODEs:

$$d\overline{R}(t) = -2\lambda_{min}\hat{\eta}(t)\mu_c(\overline{R})\overline{R}\,dt + \lambda_{max}\hat{\eta}^2(t)\nu_c(\overline{R})(\overline{R}+\gamma^2/2)\gamma\,dt + 2c^2\hat{\sigma}^2(t)\gamma^2\,dt$$

$$d\underline{R}(t) = -2\lambda_{max}\hat{\eta}(t)\mu_c(\underline{R})\underline{R}\,dt + \hat{\eta}^2(t)\nu_c(\underline{R})(\underline{R}+\gamma^2/2)\gamma\,dt + 2c^2\hat{\sigma}^2(t)\gamma^2\,dt$$

where $\quad \overline{R}(0) = \underline{R}(0) = \|\Sigma^{1/2}\vartheta^*\|^2/2$ and $\quad \lambda_{max/min} := \lambda_{max/min}(\Sigma)$.

<u>REMARK</u> : $\overline{R}(t) = \underline{R}(t)$ when $\Sigma = I$ (all $d$ ODEs in $(\bigstar)$

decouple and coincide)

Idea of the argument : Doob's decomposition of the stochastic processes given by DP-SGD and H-DP-SGD allows to write any quadratic function of the parameters as sum of predictable component + (vanishing in $n$) martingale

Equivalence between high-dimensional dynamics of SGD and a "homogenized" SDE proved in a sequence of works :

$\longrightarrow$ Least squares (no clipping, no private noise)

[ Paquette, Paquette, Adlam, Pennington, 2022 ;
Paquette, Paquette, Adlam, Pennington, 2024 ]

$\longrightarrow$ GLMs and multi-index models (no clipping, no private noise)

[ Collins-Woodfin, Paquette, Paquette, Seroussi, 2024 ]

$\longrightarrow$ Clipped dynamics (no private noise)

[ Marshall, Xiao, Agarwala, Paquette, 2025 ]

To give some intuition about the terms (I) and (II) in the INTERPRETATION, we follow Example 9 in [ Paquette, 2023 ] and consider the simplest possible setting :

*) no DP noise $( \sigma_n = 0 )$

*) no clipping $( \mu_c \equiv 1, \quad \nu_c \equiv 1 )$

*) isotropic data $( \Sigma = I )$

In this case, the one-pass SGD iteration reads :

$$\theta^{k+1} = \theta^k - \eta_k \nabla_{\theta^k} \frac{1}{2} \left( x_{n+1}^T \theta^k - y_{n+1} \right)^2$$

$$= \theta^k - \eta_k \nabla_{\theta^k} \frac{1}{2} \left( x_{n+1}^T (\theta^k - \theta^*) - z_{n+1} \right)^2$$

$y_{n+1} = x_{n+1}^T \theta^* + z_{n+1}$

$$= \theta^k - \eta_n \left( x_{n+1}^T (\theta^k - \theta^*) - z_{n+1} \right) x_{n+1}$$

We now compute the first two moments of the gradient $\left(x_{n+1}^T (\vartheta^k - \vartheta^*) - z_{n+1}\right) x_{n+1}$ :

*) $\underset{(x_{n+1}, y_{n+1})}{\mathbb{E}}\left[\left(x_{n+1}^T (\vartheta^k - \vartheta^*) - z_{n+1}\right) x_{n+1}\right] = \mathbb{E}\left[x_{n+1} x_{n+1}^T (\vartheta^k - \vartheta^*)\right] = \vartheta^k - \vartheta^* = \nabla R(\vartheta^k)$

$z_{k+1} \perp x_{k+1}$ and $\mathbb{E}\left[z_{n+1}\right] = 0$ $\qquad \mathbb{E}\, x_{n+1} x_{n+1}^T = \Sigma = I$

Thus, the mean of the gradient gives $(\mathrm{I})$.

*) $\underset{(x_{n+1}, y_{n+1})}{\mathbb{E}}\left[\left(x_{n+1}^T (\vartheta^k - \vartheta^*) - z_{n+1}\right)^2 x_{n+1} x_{n+1}^T\right]$

$= \left(\underset{(x_{n+1}, y_{n+1})}{\mathbb{E}}\left[\left(x_{n+1}^T (\vartheta^k - \vartheta^*) - z_{n+1}\right)^2\right] + z^2\right) I + 2(\vartheta^k - \vartheta^*)(\vartheta^k - \vartheta^*)^T$

some manipulations
using Wick rule

$= 2 R(\vartheta^k) I + \underbrace{2(\vartheta^k - \vartheta^*)(\vartheta^k - \vartheta^*)^T}_{\text{lower order correction}} \approx 2R(\vartheta^k) I$

Thus, the covariance of the gradient gives something that looks like $(\mathrm{II})$.

Let us now compute

$\underset{(x_{n+1}, y_{n+1})}{\mathbb{E}}\left[R^{exc}(\vartheta^{k+1}) - R^{exc}(\vartheta^k)\right] = \underset{(x_{n+1}, y_{n+1})}{\mathbb{E}}\left[\frac{1}{2}\|\vartheta^{k+1} - \vartheta^*\|^2 - \frac{1}{2}\|\vartheta^k - \vartheta^*\|^2\right]$

$= \underset{(x_{n+1}, y_{n+1})}{\mathbb{E}}\left[\frac{1}{2}\left(\|\vartheta^{k+1} - \vartheta^k\|^2 + 2 \langle \vartheta^{k+1} - \vartheta^k, \vartheta^k - \vartheta^* \rangle\right)\right]$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \mathrm{Tr}\left[2 R(\vartheta^k) I\right.$

$= \frac{1}{2} \gamma_n^2 \underset{(x_{n+1}, y_{n+1})}{\mathbb{E}}\left[\left(x_{n+1}^T (\vartheta^k - \vartheta^*) - z_{n+1}\right)^2 \|x_{n+1}\|^2\right]$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \left. + 2(\vartheta^k - \vartheta^*)(\vartheta^k - \vartheta^*)^T\right]$

$+ \underset{(x_{n+1}, y_{n+1})}{\mathbb{E}}\left[- \gamma_n \left(x_{n+1}^T (\vartheta^k - \vartheta^*) - z_{n+1}\right) x_{n+1}^T\right](\vartheta^k - \vartheta^*)$

$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad}$
$(\vartheta^k - \vartheta^*)^T$

$$= -2\eta_k R^{exc}(\vartheta^u) + \eta_k^2 \left( d\, R^{exc}(\vartheta^u) + d\, \varsigma^2/2 + 2 R^{exc}(\vartheta^u) \right)$$

$$= -2\eta_k R^{exc}(\vartheta^u) + \eta_k^2 \gamma \left( n\, R^{exc}(\vartheta^u) + n\, \varsigma^2/2 + \underbrace{\frac{2 R^{exc}(\vartheta^u)}{\gamma}} \right)$$

*Can be neglected as the other terms in the parenthesis are of order $n$*

Recalling $\quad \eta_u = \dfrac{\widehat{\eta}(u/n)}{n} \quad$ and setting $\quad \rho(t) = \lim\limits_{n \to +\infty} \mathbb{E}\left[ R^{exc}(\vartheta^{\lfloor tn \rfloor}) \right],$

we have that the equation above is the Euler-Maruyama discretization of

$$d\rho = -2\,\widehat{\eta}(t)\,\rho\,dt + \widehat{\eta}^2(t)\left(\rho + \varsigma^2/2\right)\gamma\,dt$$

*this corresponds to the ODE for $R(t)$ when $\Sigma = I$, $\mu_c \equiv 1$, $\nu_c \equiv 1$, $\widetilde{\sigma} \equiv 0$*

Already when $\Sigma \neq I$ (still no clipping and no DP noise), the SGD dynamics depends on unboundedly many statistics — and not just the first two moments as above (see Example 13 in [Paquette, 2023]). Thus, we need the framework based on homogenized SGD (and Doob's decomposition).

# Noise in the last iteration

<u>REMARK</u> : sup is taken over $t \in [0, 1)$ so equivalence does not hold for the very last iterate $\theta^n$ which gives the private output of the algorithm.

For $k < n$,

$$\rho^2 \sigma_k^2 = -\frac{d}{dk} \eta_k^2 \approx -\frac{1}{n} \frac{d}{d(k/n)} \frac{\tilde{\eta}^2(k/n)}{n^2} = \frac{\rho^2}{n^3} \tilde{\sigma}^2(k/n)$$

For $k = n$,

$$\rho^2 \sigma_n^2 = \eta_n^2 = \frac{\tilde{\eta}^2(1)}{n^2}$$

<span style="color:blue">additional $n$ factor at the denominator on top</span>

This implies that, if $\hat{\eta}(1) > 0$, the noise added to the last iterate is much larger. After a few calculations, the final loss can be shown to be given by

$$\left| R(\theta^n) - R(\theta^{n-1}) - \frac{2c^2 \tilde{\eta}^2(1) \delta^2}{\rho^2} \right| = O\left( \frac{\log n}{\sqrt{n}} \right)$$

# OPTIMAL RATES & BENEFIT OF CLIPPING

To study DP-SGD, we can now analyze H-DP-SGD or, even earlier, the ODEr.

*) We stick to $\Sigma = I$ (see [Bombari, Seroussi, M., 2025] for tracking the dependence on $d_{max}/d_{min}$ in the bounds).

*) We work with $\gamma = \frac{d}{n} \to 0$ and use $o_\gamma(\cdot), O_\gamma(\cdot), \Omega_\gamma(\cdot)$.

*) We consider three schedules (corresponding to different choices of learning rate $\hat{\eta}(t)$ and noise standard deviation $\hat{\sigma}(t)$):

(a) **Output perturbation** $\hat{\eta}(t) = \hat{\eta}(0)$, $\hat{\sigma}(t) = 0$

Learning rate is fixed during all training and private noise added only at the end of the algorithm.

$\longrightarrow$ Corresponding DP-SGD output denoted by $\theta^n_{OUT}$

(b) **Constant noise** $\hat{\eta}(t) = \hat{\eta}(0)\sqrt{1-t}$, $\hat{\sigma}(t) = \frac{1}{\rho^2}\hat{\eta}^2(0)$

Linearly decaying $\hat{\eta}^2(t)$ corresponds to constant level of noise

$\longrightarrow$ Corresponding DP-SGD output denoted by $\theta^n_{CNST}$

(c) **Optimal schedule** $\hat{\eta}(t) = \dfrac{c_1}{t + \gamma \max(c_2, c_3/\rho)}$

(for some explicit constants $c_1, c_2, c_3$ independent of $\gamma, \rho$)

This is <u>OPTIMAL</u> in the sense that it gives the optimal decay rate of the excess error (as a function of $\gamma$).

$\longrightarrow$ Corresponding DP-SGD output denoted by $\theta^n_{OPT}$

THEOREM [Bombari, Seroussi, M., 2026] Let $\rho = \Omega_\gamma(\gamma^{1-h})$ for some $h > 0$.

① Pick $c = O_\gamma(1)$ and $\tilde{\eta}(0) c = C \log(1/\gamma)$ for a large enough constant $C$ (independent of $\gamma$). Then, with high probability,

$$R^{exc}(\theta^n_{OUT}) = O_\gamma\left( \gamma \log(1/\gamma) + \frac{\gamma^2 \log^2(1/\gamma)}{\rho^2} \right),$$

$$R^{exc}(\theta^n_{CNST}) = O_\gamma\left( \gamma \log^{2/3}(1/\gamma) + \frac{\gamma^2 \log^{4/3}(1/\gamma)}{\rho^2} \right).$$

Furthermore, a lower bound of the same order holds for any choice of hyperparameters $c, \tilde{\eta}(0)$.

② Pick $c$ to be a suitable constant (independent of $\gamma$). Then, with high probability,

$$R^{exc}(\theta^n_{OPT}) = O_\gamma\left( \gamma + \frac{\gamma^2}{\rho^2} \right).$$

Furthermore, a minimax lower bound of the same order holds: let $\theta^n$ be the solution found by a generic algorithm $M$ belonging to the set of all $\rho^2/2 - z$ CDP algorithms $\mathbb{M}$; then, we have

$$\inf_{M \in \mathbb{M}} \sup_{\|\theta^*\| < 1} \mathbb{E}\left[ R^{exc}(\theta^n) \right] = \Omega_\gamma\left( \gamma + \frac{\gamma^2}{\rho^2} \right).$$

# INTERPRETATION

*) By properly tuning the learning rate $\hat{\eta}(t)$ (and, consequently, the noise schedule $\tilde{\sigma}(t)$), DP-SGD achieves minimax optimal rates.

*) The sharp analysis of DP-SGD allows to identify optimal hyperparameters. In particular, aggressive clipping $(c = O_\gamma(1))$ leads to good performance: if $c = \omega_\gamma(1)$, then

$$R^{exc}(\theta^n_{OUT}) = \Omega_\gamma\left( \gamma \log(1/\gamma) + \underbrace{\max(1, c^2)}_{} \frac{\gamma^2 \log^2(1/\gamma)}{\rho^2} \right)$$

penalty paid for using large clipping constant

$$R^{exc}(\theta^n_{CNST}) = \Omega_\gamma\left( \gamma \log^{2/3}(1/\gamma) + \max(1, c^2) \frac{\gamma^2 \log^{4/3}(1/\gamma)}{\rho^2} \right)$$

# PROOF IDEAS

*) Bounds on $\theta^n_{OUT}$, $\theta^n_{CNST}$ proved by analyzing the corresponding ODEs for $R(t)$

*) Minimax lower bound obtained by extending the analysis of [Cai, Wang, Zhang, 2021] for $(\varepsilon, \delta)-DP$ to zCDP

*) Optimal schedule obtained by crafting $\hat{\eta}(t)$ in order to match the minimax lower bound

[OPEN PROBLEM] Sharp analysis of clipped dynamics for non-linear models

# References

[Abadi, Chu, Goodfellow, McMahan, Mironov, Talwar, Zhang, 2016] Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. "Deep Learning with Differential Privacy". In: *ACM SIGSAC Conference on Computer and Communications Security*. 2016.

[Adamczak, Litvak, Pajor, Tomczak-Jaegermann, 2011] Radoslaw Adamczak, Alexander E Litvak, Alain Pajor, and Nicole Tomczak-Jaegermann. "Restricted isometry property of matrices with independent columns and neighborly polytopes by random sampling". In: *Constructive Approximation*, 2011.

[Bartlett, Montanari, Rakhlin, 2021] Peter L Bartlett, Andrea Montanari, and Alexander Rakhlin. "Deep learning: a statistical viewpoint". In: *Acta numerica*, 2021.

[Bassily, Feldman, Talwar, Thakurta, 2019] Raef Bassily, Vitaly Feldman, Kunal Talwar, and Abhradeep Thakurta. "Private Stochastic Convex Optimization with Optimal Rates". In: *Advances in Neural Information Processing Systems*. 2019.

[Bassily, Smith, Thakurta, 2014] Raef Bassily, Adam Smith, and Abhradeep Thakurta. "Private Empirical Risk Minimization: Efficient Algorithms and Tight Error Bounds". In: *2014 IEEE 55th Annual Symposium on Foundations of Computer Science*. 2014.

[Bombari, Amani, Mondelli, 2022] Simone Bombari, Mohammad Hossein Amani, and Marco Mondelli. "Memorization and optimization in deep neural networks with minimum over-parameterization". In: *Advances in Neural Information Processing Systems*, 2022.

[Bombari, Kiyani, Mondelli, 2023] Simone Bombari, Shayan Kiyani, and Marco Mondelli. "Beyond the Universal Law of Robustness: Sharper Laws for Random Features and Neural Tangent Kernels". In: *International Conference on Machine Learning (ICML)*. 2023.

[Bombari, Mondelli, 2025] Simone Bombari and Marco Mondelli. "Privacy for free in the overparameterized regime". In: *Proceedings of the National Academy of Sciences*, 2025.

[Bombari, Seroussi, Mondelli, 2025] Simone Bombari, Inbar Seroussi, and Marco Mondelli. "Better rates for private linear regression in the proportional regime via aggressive clipping". In: *arXiv preprint arXiv:2505.16329*, 2025.

[Bubeck, Sellke, 2021] Sebastien Bubeck and Mark Sellke. "A Universal Law of Robustness via Isoperimetry". In: *Advances in Neural Information Processing Systems (NeurIPS)*. 2021.

[Bun, Steinke, 2016] Mark Bun and Thomas Steinke. "Concentrated differential privacy: Simplifications, extensions, and lower bounds". In: *Theory of cryptography conference*. Springer. 2016, pp. 635–658.

[Cai, Wang, Zhang, 2021] T. Tony Cai, Yichen Wang, and Linjun Zhang. "The cost of privacy: Optimal rates of convergence for parameter estimation with differential privacy". In: *The Annals of Statistics*, 2021.

[Chizat, Oyallon, Bach, 2019] Lenaic Chizat, Edouard Oyallon, and Francis Bach. "On lazy training in differentiable programming". In: *Neural Information Processing Systems (NeurIPS)*. Vol. 32. 2019.

[Collins-Woodfin, Paquette, Paquette, Seroussi, 2024] Elizabeth Collins-Woodfin, Courtney Paquette, Elliot Paquette, and Inbar Seroussi. "Hitting the High-dimensional notes: an ODE for SGD learning dynamics on GLMs and multi-index models". In: *Information and Inference: A Journal of the IMA*, 2024.

[Cover, 1965] Thomas M. Cover. "Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition". In: *IEEE Transactions on Electronic Computers*, 1965.

[Du, Zhai, Poczos, Singh, 2019] Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. "Gradient Descent Provably Optimizes Over-parameterized Neural Networks". In: *International Conference on Learning Representations (ICLR)*. 2019.

[Dwork, McSherry, Nissim, Smith, 2006] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. "Calibrating Noise to Sensitivity in Private Data Analysis". In: *Theory of Cryptography, Third Theory of Cryptography Conference, TCC 2006*. 2006.

[Dwork, Roth, 2014] Cynthia Dwork and Aaron Roth. "The algorithmic foundations of differential privacy". In: *Foundations and Trends in Theoretical Computer Science*, 2014.

[Feldman, Koren, Talwar, 2020] Vitaly Feldman, Tomer Koren, and Kunal Talwar. "Private stochastic convex optimization: optimal rates in linear time". In: *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*. STOC 2020. 2020. ISBN: 9781450369794.

[Ghorbani, Mei, Misiakiewicz, Montanari, 2021] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. "Linearized two-layers neural networks in high dimension". In: *The Annals of Statistics*, 2021.

[Iurada, Bombari, Tommasi, Mondelli, 2026] Leonardo Iurada, Simone Bombari, Tatiana Tommasi, and Marco Mondelli. "A Law of Data Reconstruction for Random Features (and Beyond)". In: 2026.

[Jacot, Gabriel, Hongler, 2018] Arthur Jacot, Franck Gabriel, and Clément Hongler. "Neural tangent kernel: Convergence and generalization in neural networks". In: *Neural Information Processing Systems (NeurIPS)*. 2018.

[Kifer, Smith, Thakurta, 2012] Daniel Kifer, Adam Smith, and Abhradeep Thakurta. "Private Convex Empirical Risk Minimization and High-dimensional Regression". In: *Conference on Learning Theory*. 2012.

[Madry, Makelov, Schmidt, Tsipras, Vladu, 2018] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. "Towards Deep Learning Models Resistant to Adversarial Attacks". In: *International Conference on Learning Representations (ICLR)*. 2018.

[Marshall, Xiao, Agarwala, Paquette, 2025] Noah Marshall, Ke Liang Xiao, Atish Agarwala, and Elliot Paquette. "To Clip or not to Clip: the Dynamics of SGD with Gradient Clipping in High-Dimensions". In: *The Thirteenth International Conference on Learning Representations*. 2025.

[Mei, Misiakiewicz, Montanari, 2022] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. "Generalization error of random feature and kernel methods: hypercontractivity and kernel matrix concentration". In: *Applied and Computational Harmonic Analysis*, 2022.

[Misiakiewicz, Montanari, 2024] Theodor Misiakiewicz and Andrea Montanari. "Six lectures on linearized neural networks". In: *Journal of Statistical Mechanics: Theory and Experiment*, 2024.

[Montanari, Zhong, 2022] Andrea Montanari and Yiqiao Zhong. "The interpolation phase transition in neural networks: Memorization and generalization under lazy training". In: *The Annals of Statistics*, 2022.

[Nguyen, Mondelli, Montufar, 2021] Quynh Nguyen, Marco Mondelli, and Guido Montufar. "Tight bounds on the smallest eigenvalue of the neural tangent kernel for deep ReLU networks". In: *International Conference on Machine Learning (ICML)*. 2021.

[Oymak, Soltanolkotabi, 2019] Samet Oymak and Mahdi Soltanolkotabi. "Overparameterized nonlinear learning: Gradient descent takes the shortest path?" In: *International Conference on Machine Learning (ICML)*. 2019.

[Paquette, Paquette, Adlam, Pennington, 2024] Courtney Paquette, Elliot Paquette, Ben Adlam, and Jeffrey Pennington. "Homogenization of SGD in high-dimensions: exact dynamics and generalization properties". In: *Mathematical Programming*, 2024.

[Paquette, Paquette, Adlam, Pennington, 2022] — . "Implicit Regularization or Implicit Conditioning? Exact Risk Trajectories of SGD in High Dimensions". In: *Advances in Neural Information Processing Systems*. 2022.

[Paquette, 2023] Elliot Paquette. *High-dimensional limits of stochastic gradient descent*. Lecture notes. 2023.

[Rahimi, Recht, 2007] Ali Rahimi and Benjamin Recht. "Random features for large-scale kernel machines". In: *Advances in Neural Information Processing Systems*, 2007.

[Soltanolkotabi, Javanmard, Lee, 2018] Mahdi Soltanolkotabi, Adel Javanmard, and Jason D Lee. "Theoretical insights into the optimization landscape of over-parameterized shallow neural networks". In: *IEEE Transactions on Information Theory*, 2018.

[Szegedy, Zaremba, Sutskever, Bruna, Erhan, Goodfellow, Fergus, 2014] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus. "Intriguing properties of neural networks". In: *International Conference on Learning Representations (ICLR)*. 2014.

[Xiao, Hu, Misiakiewicz, Lu, Pennington, 2022] Lechao Xiao, Hong Hu, Theodor Misiakiewicz, Yue Lu, and Jeffrey Pennington. "Precise learning curves and higher-order scalings for dot-product kernel regression". In: *Advances in Neural Information Processing Systems*, 2022.