# §2. Escaping Saddle Points.

1. **Def.** $f: \mathbb{R}^d \to \mathbb{R}$. $x$ is said to be a **local minimizer** of $f$ if $\exists$ open $U \subseteq \mathbb{R}^d$ s.t. $x \in U$ and $f(x') \geq f(x) \ \forall x' \in U$.

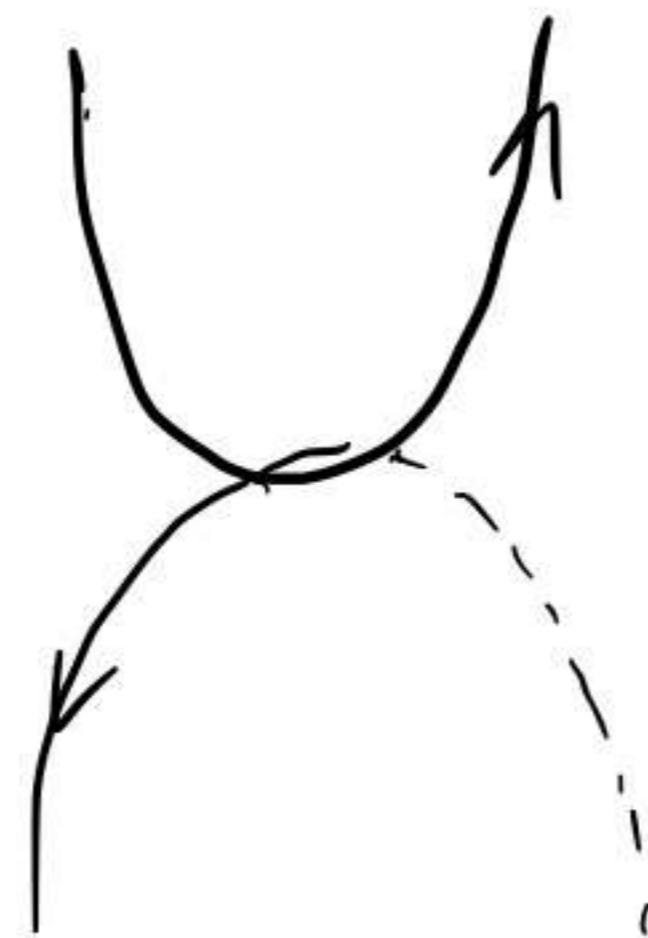**2nd order condition.** $\nabla f(x) = 0$, $\nabla^2 f(x) \succ 0$

$\Rightarrow x$ is a **strict** local minimizer

$*$ $\nabla f(x) = 0$, $\nabla^2 f(x) \succeq 0 \not\Rightarrow$ local minimizer.

2nd order stationary point

e.g. $f(x) = x^3$. $x = 0$

$f(x_1, x_2) = x_1^2 - x_2^4$.



2. **Def.** $f: \mathbb{R}^d \to \mathbb{R}$. $\rho$-Hessian Lipschitz. We say $x$ is a **$\varepsilon$-second order stationary point** of $f$ if.

$$\|\nabla f(x)\| \leq \varepsilon \quad \text{and} \quad \lambda_{min}(\nabla^2 f(x)) \geq -\sqrt{\rho \varepsilon}.$$

· error terms: Consider $f(x) = \Theta(\rho x^3)$

3. **Goal:** Assume $f: \mathbb{R}^d \to \mathbb{R}$ is $\ell$-smooth and $\rho$-Hessian Lipschitz. Noisy GD/Gr can efficiently find a $\varepsilon$-second order stationary point. ($\Leftrightarrow$ they can escape saddle points with $\lambda_{min}(\nabla^2 f(x)) < -\sqrt{\rho \varepsilon}$.)

## 4. The quadratic case

Consider $f(x) = x^T A x$ where $A = \text{diag}(-1, 1, \ldots, 1)$.

clear that $\hat{x} = 0$ is a saddle point.

Claim. For any reasonable $x_0$, GF will not get stuck at $\hat{x}$.

pf. $\frac{d}{dt} x_t = -\nabla f(x_t) = -2A x_t$

$\Rightarrow \quad \frac{d}{dt} x_{k,t} = -2 A_{kk} x_{k,t}$

$\Rightarrow \quad x_{k,t} = x_{k,0} \exp(-2 A_{kk} t)$

$\Rightarrow \quad x_{k,t} \to 0 \quad (k \neq 1), \quad |x_{1,t}| \to \infty$

As long as $|x_{1,0}| \geq 1/\text{poly}(d)$, $|x_{1,t}|$ will

become $\Omega(1)$ in $O(\log d)$ time.

Lemma. (Anti-concentration of ball volume).

$x \sim \text{Unif}(B(r)) \quad \forall \delta \in (0, 1)$

with proba. $\geq 1 - \delta$, we have

$$|x_1| \geq r\delta / (2\sqrt{d})$$

Thus, small ball perturbation

$\Rightarrow$ escaping the saddle point.

# 5. General loss function.

- Algorithm: Run GD/GF and occasionally add a small perturbation.

- High-level idea:

Suppose $x_0$ is near a saddle point with at least one descent direction. Assume

$\nabla f(x_t)$ is small $\forall t \in [0, T] \Rightarrow x_t \approx x_0$

$\Rightarrow \nabla^2 f(x_t) \approx \nabla^2 f(x_0)$.

$\frac{d}{dt} \| \nabla f(x_t) \|^2 = 2 \langle \nabla f(x_t), \frac{d}{dt} \nabla f(x_t) \rangle$

$= 2 \langle \nabla f(x_t), \nabla^2 f(x_t) \dot{x}_t \rangle$

$= -2 \langle \nabla f(x_t), \nabla^2 f(x_t) \nabla f(x_t) \rangle$

$\approx -2 \langle \nabla f(x_t), \nabla^2 f(x_0) \nabla f(x_t) \rangle \Rightarrow \nabla f(x_t)$ will blow up along the descent dir. Contradiction.

**Thm** (Thm 2 of Jin et al. (2017))

$f: \mathbb{R}^d \to \mathbb{R}$ $\ell$-smooth. $\rho$-Hessian Lipschitz.

$\forall \epsilon < \ell^2/\rho$, $\delta \in (0, 1)$, with proba. $\geq 1 - \delta$.

noisy GD can output a $\epsilon$-second order stationary pt. with

$$O\left( \frac{\ell \, (f(x_0) - f_t)}{\epsilon^2} \log^4 \left( \frac{d\ell \, (f(x_0) - f_*)}{\epsilon^2 \delta} \right) \right)$$

iterations.

* poly log (d).

# §3. Neural Tangent Kernel (NTK)

1. Themes of DL Theory.

- "Traditional" supervised learning:
  - Why GD + NN can (over)fit the training set?
    (Global convergence)
  - What solutions can GD find? Why they can generalize?
    (algorithmic regularization.
  - ~ · · · ·

- Not so supervised. (Pretraining + finetuning).
  - Different pre-training tasks. (contrasive learning.
                      reconstruction · · · · )

  - What representations do they learn?
  - Why they can be used in downstream tasks.
  - ~ · · · ·

- In-context learning (Prompting)
  - We don't know what questions
    to ask yet · · ·

- Generic template.
  - Why X works?

  - Why X is better than Y?

# 2. Background of NTK

- (Zhang et al., 2016). (experimental)

  GD + NN can often _globally_ minimize the loss, even when the labels are random.

- Over-parameterized networks.

  - Original meaning: # params $> n$
    
    $\llcorner$ # training samples

  - What people mean in DL theory:
    
    # neurons (per layer) $= \text{poly}(r) \quad \curvearrowleft$ latent dim.
    
    or $\text{poly}(d) \quad \curvearrowleft$ input dim.
    
    or $\text{poly}(d, n)$
    
    or $\exp(d)$ or $\infty$

# 3. Setup

- Training set. $\{(x_i, y_i)\}_{i=1}^{n} \subseteq \mathbb{R}^d \times \mathbb{R}$.

- Learner network.

  $$f(x; W, a) = \frac{1}{\sqrt{m}} a^T \sigma(Wx) = \frac{1}{\sqrt{m}} \sum_{k=1}^{m} a_k \sigma(w_k \cdot x)$$

  where $x \in \mathbb{R}^d$ input.
  
  $\quad m \in \mathbb{N}$. # neurons.
  
  $\quad W \in \mathbb{R}^{m \times d}$ first layer weights
  
  $\quad \sigma : \mathbb{R} \to \mathbb{R}$. ReLU.
  
  $\quad a \in \mathbb{R}^m$ output weights

- Initialization. $a_k \sim \text{Unif}(\{\pm 1\})$
  
  $\quad w_k \sim N(0, I_d)$.

- Training algorithm: Fix $a$. Train $W$ with
  
  GD + MSE. $L(W) = \frac{1}{2} \sum_{i=1}^{n} (y_i - f(x_i; W))^2$

# 4. Derivatives.

$$\nabla_{W_k} L(W) = \sum_{i=1}^{n} (y_i - f(x_i; W)) \nabla_{W_k} f(x_i; W)$$

$$= -\sum_{i=1}^{n} (y_i - f(x_i; W)) \frac{a_k}{\sqrt{m}} \nabla_{W_k} \sigma(W_k \cdot x)$$

$$= -\frac{a_k}{\sqrt{m}} \sum_{i=1}^{n} (y_i - f(x_i; W)) \sigma'(W_k \cdot x_i) x_i$$

Recall GF: $\frac{d}{dt} W_k = -\nabla_{W_k} L(W)$.

$$\frac{d}{dt} f(x_j) = \frac{1}{\sqrt{m}} \sum_{k=1}^{m} a_k \frac{d}{dt} \sigma(W_k \cdot x_j)$$

$$= \frac{1}{\sqrt{m}} \sum_{k=1}^{m} a_k \left\langle \sigma'(W_k \cdot x_j) x_j, \frac{d}{dt} W_k \right\rangle$$

$$= \frac{1}{\sqrt{m}} \sum_{k=1}^{m} a_k \left\langle \sigma'(W_k \cdot x_j) x_j, \frac{a_k}{\sqrt{m}} \sum_{i=1}^{n} (y_i - f(x_i; W)) \sigma'(W_k \cdot x_i) x_i \right\rangle$$

$$= \frac{1}{m} \sum_{i=1}^{n} (y_i - f(x_i; W)) \sum_{k=1}^{m} a_k^2 \, ^{1} \sigma'(W_k \cdot x_i) \sigma'(W_k \cdot x_j) \langle x_i, x_j \rangle$$

Define the **NTK**,

$$H(x, x') = \frac{1}{m} \sum_{k=1}^{m} \sigma'(W_k \cdot x) \sigma'(W_k \cdot x_j) \langle x_i, x_j \rangle$$

and $H_{ij} = H(x_i, x_j)$.   $i, j \in [n]$.

Define $y = (y_i)_{i=1}^{n}$,   $f_t = \left( f(x_i; W_t) \right)_{i=1}^{n}$

Then

$$\frac{d}{dt} f_t = H_t (y - f_t).$$

Remarks.

a) $H_t$ depends on $t$.

b) NTK and these formula themselves are quite generic, while the NTK technique is not.

c) If $H_t \succeq \lambda_0 I_d$ for some $\lambda_0 > 0$, $\forall t$, then $f_t \to y$ linearly.

5. The NTK technique.

Choose $m$ and the initialization scale (im)properly to ensure $H_t \approx H_0$ throughout training, and reduce the problem to a convex one.

6. Define $H^\infty \in \mathbb{R}^{n \times n}$ by

$$H_{i,j}^\infty = \lim_{m \to \infty} H_{i,j}(0)$$
$$= \lim_{m \to \infty} \frac{1}{m} \sum_{k=1}^{m} \sigma'(w_k(0) \cdot x_i) \, \sigma'(w_k(0) \cdot x_j) \langle x_i, x_j \rangle$$
$$= \mathop{\mathbb{E}}_{W \sim N(0, J_d)} \left[ \sigma'(W \cdot x_i) \, \sigma'(W \cdot x_j) \langle x_i, x_j \rangle \right].$$

Define $\lambda_0 := \lambda_{min}(H^\infty)$.

<u>Fact</u>. If $x_i \not\parallel x_j$, $\forall i \neq j$, then $\lambda_0 > 0$.

<u>Assume</u> $\lambda_0 > 0$.

Du et al., 2019. Gradient descent provably optimizes over-parametrized neural networks

7. **Lemma.** Choose $m = \Omega\left(\frac{n^2}{\lambda_0^2} \log\left(\frac{n}{\delta}\right)\right)$.

With proba. $\geq 1-\delta$, we have

$\|H_0 - H^{(\infty)}\|_2 \leq \lambda_0/4$, whence $\lambda_{min}(H_0) \geq \frac{3}{4}\lambda_0$.

**Lemma.** Initialization $w_1, \ldots, w_m$.

Any $\hat{w}_1, \ldots, \hat{w}_m \in \mathbb{R}^d$ with $\|w_k - \hat{w}_k\| \leq \mathcal{O}\left(\frac{\delta \lambda_0}{n^2}\right) =: R$.

Let $H$ and $\hat{H}$ be the corresponding NTKs. With

proba. $\geq 1-\delta$ over the random init., we have.

$\|\hat{H} - H\|_2 \leq \lambda_0/4$, whence $\lambda_{min}(\hat{H}) \geq \lambda_0/2$.

Deterministic.

The infinite-width limit.

$H^\infty$ ← Depends only on the initialization scheme

Depends only on the actual random initialization. ⤳ $H_0$

$H_t$ ← Also depends on the training procedure.

**8. Thm.** Assume $\|x_i\|, |y_i| \leq C$. Choose $m = \Omega\left(\frac{n^6}{\delta^3 \lambda_0^4}\right)$.

With proba. $\geq 1 - O(\delta)$, GF converges to a point with

$\|y - f_t\| \leq \varepsilon$ within $O\left(\frac{1}{\lambda_0} \log\left(\frac{\|y - f_0\|}{\varepsilon}\right)\right)$ amount of time.

In words. GF will fit the training set before $W_k$ moves too far away from the initialisation.

**Pf.** Define $T_* := \min\{T_1, T_2\}$ where

$$T_1 := \inf\{t \geq 0 : \|f_t - y\| \leq \varepsilon\}$$

$$T_2 := \inf\{t \geq 0 : \exists k, \|W_k(t) - W_k(0)\| \geq R\}$$

By the previous lemmas, $\forall t \leq T_* \leq T_2$, we have

$$\frac{d}{dt}\|y - f_t\|^2 = -2\langle y - f_t, H_t(y - f_t)\rangle$$

$$\leq -\lambda_0 \|y - f_t\|^2$$

$$\Rightarrow \|y - f_t\|^2 \leq \|y - f_0\|^2 \exp(-\lambda_0 t).$$

$$\left(\Rightarrow \|y - f_t\| \leq \|y - f_0\| \exp(-\lambda_0 t/2)\right)$$

Meanwhile, we have

$$\|\nabla_{W_k} \mathcal{L}(w)\| = \left\|\frac{a_k}{\sqrt{m}} \sum_{i=1}^{n} (y_i - f(x_i)) \sigma'(W_k \cdot x) x_i\right\|$$

$$\leq \frac{C}{\sqrt{m}} \sum_{i=1}^{n} |y_i - f(x_i)|$$

$$\leq C \sqrt{\frac{n}{m}} \|y - f_t\|$$

$$\Rightarrow \|W_k(t) - W_k(0)\| \leq C\sqrt{\frac{n}{m}} \int_0^t \|y - f_s\| \, ds$$

$$\leq C\sqrt{\frac{n}{m}} \|y - f_0\| \int_0^t \exp(-\lambda_0 s/2) \, ds$$

$$\leq 2C\sqrt{\frac{n}{m}} \underbrace{\|y - f_0\|}_{\leq R}$$

$$\Rightarrow T_* \text{ cannot be attained by } T_2.$$

$$\Rightarrow T_* = T_2 \leq O\left(\frac{1}{\lambda_0} \log\left(\frac{\|y - f_0\|}{\varepsilon}\right)\right).$$

**Question.** Where did we use $a_k \sim \text{Unif}\{\pm 1\}$?

**Answer.**
$$\mathbb{E}_{a,W} \|f_0\|^2 = \mathbb{E}_{a,W} \left\| \frac{1}{\sqrt{m}} \sum_{k=1}^{m} a_k \, \sigma(w_k \cdot x) \right\|^2$$

$$= \frac{1}{m} \sum_{i,j=1}^{m} \mathbb{E}_{a,W} \left\{ a_i a_j \, \sigma(w_i \cdot x) \, \sigma(w_j \cdot x) \right\}$$

If $a_k \sim \text{Unif}\{\pm 1\}$, then
$$= \frac{1}{m} \sum_{k=1}^{m} \mathbb{E}_{W} \, \sigma^2(w_k \cdot x) = \mathbb{E}_{W} \, \sigma^2(w_k \cdot x)$$

2) $a_k = 1$, then
$$= \frac{1}{m} \sum_{i,j} \mathbb{E}_{W} \, \sigma(w_i \cdot x) \, \sigma(w_j \cdot x)$$

$$\geq \Omega \left( m \left( \mathbb{E}_{w_k} \sigma(w_k \cdot x) \right)^2 \right).$$

As a result -
$$\| w_k(t) - w_k(s) \| \leq 2c \sqrt{\frac{n}{m}} \, \| y - f_0 \|$$

$$\not\to 0 \quad \text{as } m \to \infty.$$

9. NTK $\longleftrightarrow$ random feature.

If the movement is small.

$$f(x; W(t)) = \frac{1}{\sqrt{m}} \sum_{k=1}^{m} a_k \sigma(W_k \cdot x)$$

$$\approx f(x; W(0)) + \frac{1}{\sqrt{m}} \sum_{k=1}^{m} a_k \langle W_k(t) - W_k(0), x \rangle \sigma'(W_k \cdot x)$$

$\hookrightarrow \approx 0$ when $m$ is large.

Define $\alpha_k = (W_k(t) - W_k(0))/\sqrt{m}$. Then

$$f(x; W(t)) \approx \sum_{k=1}^{m} \langle \alpha_k, a_k \sigma'(W_k(0) \cdot x) x \rangle.$$

$\hookrightarrow$ a linear model over the

fixed random feature mapping

$$x \mapsto \left( a_k \sigma'(W_k(0) \cdot x) x \right)_{k=1}^{m}$$

10. Random linear features cannot
learn a single linear model

- Input distr. $x \sim N(0, I_d)$
- Target function. $f_*(x) = \langle W_*, x \rangle$.

    for some fixed unit vector $W_*$.

- Learner: $f(x; u) = \sum_{k=1}^{m} u_k \langle W_k, x \rangle$

    $$= \langle \sum_{k=1}^{m} u_k W_k, x \rangle$$

    $$\underbrace{\phantom{\sum_{k=1}^{m} u_k W_k}}_{=: W_u}$$

    where $W_k \sim N(0, I_d/d)$

Claim. If $m \leq d/2$. then w.h.p.

$\min_{u} MSE = \min_{u} \mathbb{E}_X \left\{ (f(x; u) - f_*(x))^2 \right\} \geq \frac{1}{4}$

Pf. $MSE = \mathbb{E}_X \langle W_u - W_*, x \rangle^2 = \|W_u - W_*\|^2$

$\Rightarrow \min_{u} MSE = $ dist. from $W_*$ to $\text{span}(W_1, \dots, W_m)$

By symmetry. we may assume w.l.o.g. that
$(W_1, \dots, W_m)$ are fixed and $W_*$ is a random
unit vector. Moreover assume $W_k = e_k$.

$\Rightarrow L(W_*) := \min_{u} MSE = 1 - \sum_{k=1}^{m} [W_*]_k^2$

$$\geq 1 - \sum_{k=1}^{d/2} [W_*]_k^2$$

$\hookrightarrow$ $O(1)$-Lipschitz with median $= 1/2$

$\overset{\text{Lévy}}{\Longrightarrow}$ with proba. $\geq 1 - \exp(-\Theta(d))$.

$L(W_*) \geq 1/4$.

**Lemma.** (Lévy's inequality). Let $V \sim \text{Unif}(S^{d-1})$,

$g: \mathbb{R}^d \to \mathbb{R}$. $L$-Lipschitz. $\exists C, c > 0$ s.t. $\forall \varepsilon$.

$$\mathbb{P}\left[ \, |g(V) - \text{median}(g)| \geq \varepsilon \, \right] \leq C \exp\left(-\frac{c \varepsilon^2 d}{L^2}\right).$$

- Informal version of Thm. 1.2 of Thm 4.2 of Yehudai and Shamir (2019).

$$\mathcal{F} = \left\{ x \mapsto f(Wx) : W = \begin{bmatrix} W_1 \\ \vdots \\ W_p \end{bmatrix}, \, W_k \sim \text{Unif}(S^{d-1}) \right\} \qquad N(x) = \sum_{k=1}^{m} u_i f(W_i x)$$

$\forall w^* \in \mathbb{R}^d$ with $\|w^*\| = d^2$. $\exists$ bias $b_* \in \mathbb{R}$ with $|b| = O(d^3)$ s.t. w.h.p.

$$\mathbb{E}\left\{ (N(x) - \sigma(w_* \cdot x + b_*))^2 \right\} \leq 1/50.$$

$$\Rightarrow \quad m \cdot \max_k |u_k| \geq \exp(\Omega(d)).$$

11. Continuity argument aka. induction/minimal counterexample in continuous time

- Example $f \in C^0(\mathbb{R})$. $\dot{x}_t = f(x_t)$. $f(0) > 0$. $x_0 > 0$ $\Rightarrow$ $x_t > 0$ $\forall t$.

- Example. Comparative Gronwall.

  "Progress" $X_t$. "error" $Y_t$. Suppose when $Y_t \le \delta = \frac{1}{poly(d)}$. we have

  $$\begin{cases} \dot{X}_t \le -A_t X_t \\ \dot{Y}_t \le K A_t Y_t. \end{cases} \quad \text{or} \quad \begin{cases} \dot{X}_t \le -A_t X_t \\ \dot{Y}_t \le K A_t Y_t + B. \end{cases}$$

- Example. Suppose when $Y_t \le \delta \le 1/poly(d)$, we have $\dot{X}_t \le -A X_t$. $\dot{Y}_t \le X_t Y_t$.

  $*$. Present this before the NTK results.