

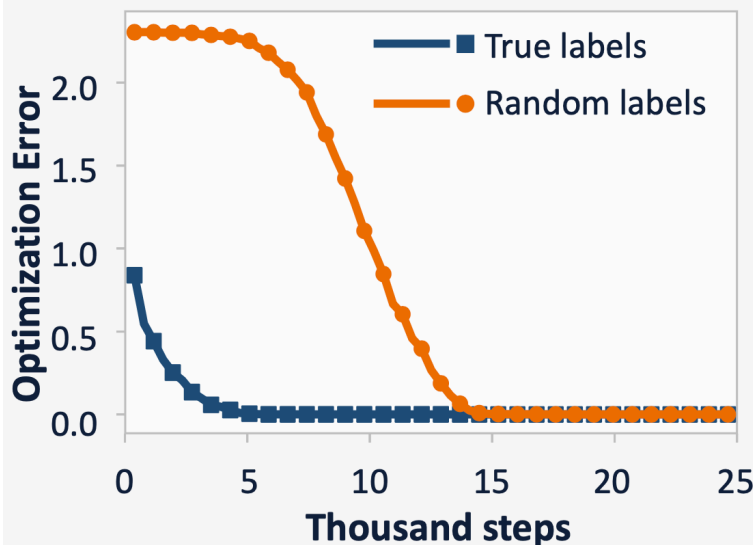
Non-convex Optimization Landscape



Gradient descent finds global minima

Practice: gradient descent

$$\theta(t+1) \leftarrow \theta(t) - \eta \frac{\partial L(\theta(t))}{\partial \theta(t)}$$



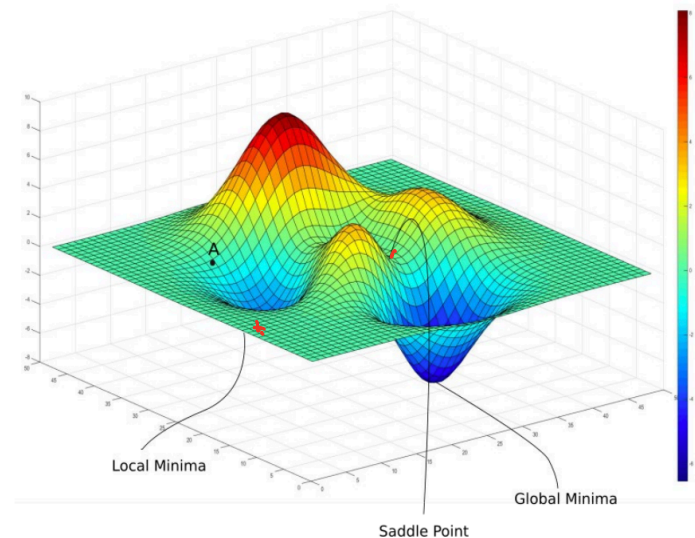
Optimization error $\rightarrow 0$ for both *true labels* and *random labels* !

Zhang Bengio Hardt Recht Vinyals 2017

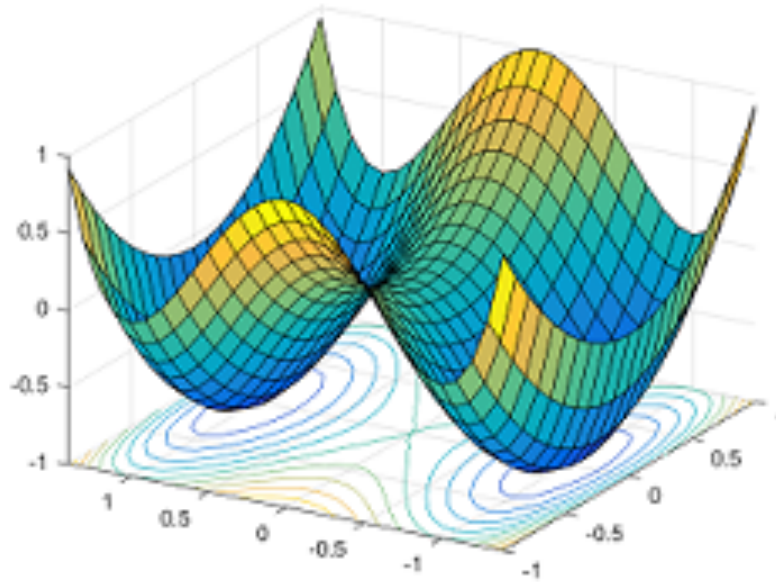
Understanding DL Requires Rethinking Generalization

Types of stationary points

- Stationary points: $x : \nabla f(x) = 0$
- Global minimum:
 $x : f(x) \leq f(x') \forall x' \in \mathbb{R}^d$
- Local minimum:
 $x : f(x) \leq f(x') \forall x' : \|x - x'\| \leq \epsilon$
- Local maximum:
 $x : f(x) \geq f(x') \forall x' : \|x - x'\| \leq \epsilon$
- Saddle points: stationary points that are not a local min/max

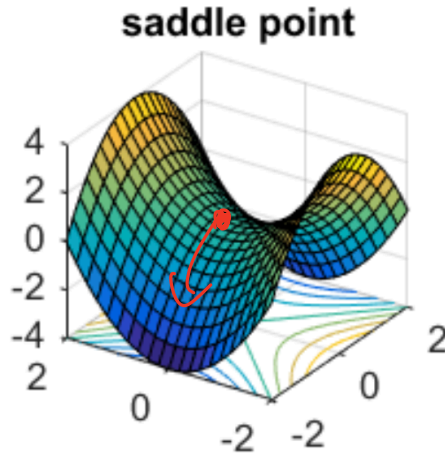


Landscape Analysis



- All local minima are global!
- Gradient descent can escape saddle points.

Strict Saddle Points (Ge et al. '15, Sun et al. '15)

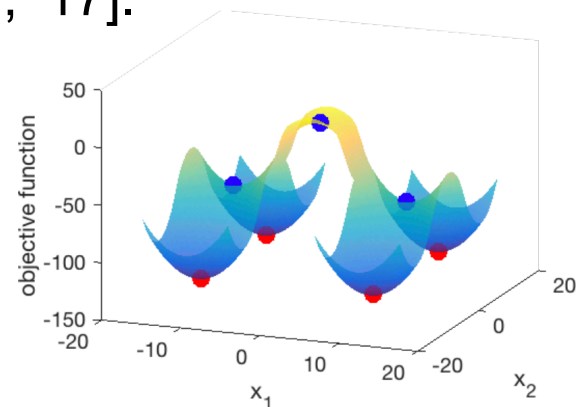


- Strict saddle point: a saddle point and $\lambda_{\min}(\nabla^2 f(x)) < 0$

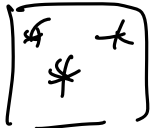

Escaping Strict Saddle Points

- **Noise-injected** gradient descent can escape strict saddle points in polynomial time [Ge et al., '15, Jin et al., '17].
- Randomly initialized gradient descent can escape all strict saddle points asymptotically [Lee et al., '15].
 - Stable manifold theorem.
- Randomly initialized gradient descent can take exponential time to escape strict saddle points [Du et al., '17].

If 1) all local minima are global, and 2) are saddle points are strict, then noise-injected (stochastic) gradient descent finds a global minimum in polynomial time



What problems satisfy these two conditions

- Matrix factorization $\min_{U, V} \|UV^T - \Sigma\|_F^2$
- Matrix sensing $i \in \{1, \dots, n\}, \langle A_i, \Sigma \rangle = y_i, \Sigma: \text{low-rank}$
- Matrix completion $\min_{U, V} \sum_{i=1}^n (\langle A_i, UV^T \rangle - y_i)^2$

- Tensor factorization $T = U \otimes U \otimes U$

- Two-layer neural network with quadratic activation

$$y_i = \sum_{j=1}^n \langle x_i, w_j \rangle^2$$

What about neural networks?

$$\sigma(z) = z$$

- Linear networks (neural networks with linear activations functions): **all local minima are global, but there exists saddle points that are not strict** [Kawaguchi '16].

$$\Rightarrow \text{saddle } x, \nabla_{\text{min}}(\nabla^2 f(x)) = 0 \quad \min_{w_0, \dots, w_n} (\sum_{i=1}^n (w_0 + w_1 x + \dots + w_n x^n - y)^2)$$

- Non-linear neural networks with:

- Virtually any non-linearity,

ReLU, Sigmoid

- Even with Gaussian inputs,

$$x \sim \mathcal{N}(0, I)$$

- Labels are generated by a neural network of the same architecture,

$$y = \mathcal{NN}(x)$$

There are many bad local minima [Safran-Shamir '18, Yun-Sra-Jadbaie '19].

not global minima

Global convergence of gradient descent



Global convergence of gradient descent

Theorem (Du et al. '18, Allen-Zhu et al. '18, Zou et al 19') If the width of each layer is $\text{poly}(n)$ where n is the number of data. Using random initialization with a particular scaling, gradient descent finds an approximate global minimum in polynomial time.

Neural Tangent Kernel

$$|f(\underline{w}) - f(\underline{w}^*)| \leq \underline{\varepsilon}$$

$$\text{poly}\left(\frac{1}{\varepsilon}, \text{width}, \text{depth}\right)$$

Gradient Flow: a Kernel Point of View

$$\bullet \mathcal{L}(\Theta) = \frac{1}{n} \sum_{i=1}^n \ell(f(\Theta, x_i), y_i)$$
$$\frac{\partial \mathcal{L}(\Theta)}{\partial \Theta} = \frac{1}{n} \sum_{i=1}^n \ell'(f(\Theta, x_i), y_i) \cdot \frac{\partial f(\Theta, x_i)}{\partial \Theta}$$

$$GF: \quad \frac{d\Theta(t)}{dt} = - \frac{\partial \mathcal{L}(\Theta)}{\partial \Theta}$$

if strongly convex, $\Rightarrow \Theta^*$, $\Theta(t) \rightarrow \Theta^*$

in NN, # of parameters $> n$

we know $t \rightarrow \infty$, $f(\Theta(t), x_i) \rightarrow y_i$

Gradient Flow: a Kernel Point of View

$$u_i(t) = f(\theta(t), x_i), \quad u(t) = \begin{pmatrix} u_1(t) \\ \vdots \\ u_n(t) \end{pmatrix}$$

$$\frac{d u_i(t)}{d t} = \left\langle \frac{\partial u_i(t)}{\partial \theta(t)}, \frac{d \theta(t)}{d t} \right\rangle$$

$$l'(u(t), y) \in \mathbb{R}^n \quad l' = \left\langle \frac{\partial u_i(t)}{\partial \theta(t)}, -\frac{1}{n} \sum_{j=1}^n l'(u_j(t), y_j) - \frac{\partial u_i(t)}{\partial \theta(t)} \right\rangle$$

$$[l'(u(t), y)]_i = l'(u_i(t), y_i)$$

$$= -\frac{1}{n} [l'(u_1(t), y_1), \dots, l'(u_n(t), y_n)]$$

$$H(t) \in \mathbb{R}^{n \times n}$$

$$\left(\left\langle \frac{\partial u_i(t)}{\partial \theta(t)}, \frac{\partial u_i(t)}{\partial \theta(t)} \right\rangle, \dots, \left\langle \frac{\partial u_i(t)}{\partial \theta(t)}, \frac{\partial u_i(t)}{\partial \theta(t)} \right\rangle \right)$$

$$= \left\langle \frac{\partial u_i(t)}{\partial \theta(t)}, \frac{\partial u_i(t)}{\partial \theta(t)} \right\rangle, \quad \frac{d u(t)}{d t} = -\frac{1}{n} H(t) \cdot l'(u(t), y)$$

Gradient Flow: a Kernel Point of View

If l is quadratic loss, $l(u(t), y) = \frac{1}{2} (u(t) - y)^2$

$$l'(u(t), y) = u(t) - y$$

$$(*) \quad \frac{du(t)}{dt} = -\frac{1}{n} \underbrace{H(t)}_{\text{does not depend on loss}} (u(t) - y)$$

If $\exists \lambda_0 > 0$, $\forall t, \lambda_{\min}(H(t)) \geq \lambda_0$

$$\frac{d \left(\underbrace{\frac{1}{2} \|u(t) - y\|_2^2}_{\text{loss}} \right)}{dt} = -\frac{1}{n} (u(t) - y)^T H(t) (u(t) - y) \leq -\frac{\lambda_0}{n} \|u(t) - y\|_2^2$$

Gradient Flow: a Kernel Point of View

$$\begin{aligned}
 & \text{Consider } \frac{d}{dt} \left(\exp\left(\frac{\lambda_0 t}{n}\right) \cdot \frac{1}{2} \|u(t) - y\|_2^2 \right) \\
 &= \frac{\lambda_0}{2n} \exp\left(\frac{\lambda_0 t}{n}\right) \|u(t) - y\|_2^2 + \frac{d\left(\frac{1}{2} \|u(t) - y\|_2^2\right)}{dt} \cdot \exp\left(\frac{\lambda_0 t}{n}\right) \\
 &\leq \exp\left(\frac{\lambda_0 t}{n}\right) \|u(t) - y\|_2^2 \left(\frac{\lambda_0}{2n} - \frac{\lambda_0 t}{n} \right) < 0
 \end{aligned}$$

$\Rightarrow \exp\left(\frac{\lambda_0 t}{n}\right) \cdot \frac{1}{2} \|u(t) - y\|_2^2$ is decreasing

$t=0, \quad \frac{1}{2} \|u(0) - y\|_2^2$ ()

$\forall t \quad \exp\left(\frac{\lambda_0 t}{n}\right) \cdot \frac{1}{2} \|u(t) - y\|_2^2 \leq \left(\frac{1}{2} \|u(0) - y\|_2^2 \right)$

$\log\left(\frac{1}{2} \text{ error}\right)$

$\frac{1}{2} \|u(t) - y\|_2 \leq \left(\frac{1}{2} \|u(0) - y\|_2 \right) \exp\left(-\frac{\lambda_0 t}{n}\right)$
 $t \rightarrow \infty, \text{ loss} \rightarrow 0 \|u(t) - y\|_2$

Gradient Flow: a Kernel Point of View

kernel: $[U(t)]_i = \phi(X_i)^T \Theta(t)$ ϕ : feature

$$[H(t)]_{ij} = \left\langle \frac{\partial u_i(t)}{\partial \Theta(t)}, \frac{\partial u_j(t)}{\partial \Theta(t)} \right\rangle$$

$$= \langle \phi(X_i), \phi(X_j) \rangle$$

$$= K(X_i, X_j) \quad \text{does not depend on } t$$

$$H(t) = \begin{pmatrix} K(X_i, X_j) \end{pmatrix} \quad \text{if Kernel is full rank}$$

$$\lambda_{\min}(H(t)) > 0 \quad \text{if kernel is universal}$$

(Gaussian) $\lambda_{\min}(H(t)) > 0$

Gradient Flow: a Kernel Point of View

- $f(\theta, x) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \phi(w_r^T x)$
- $x \in \mathbb{R}^d, w_r \in \mathbb{R}^d, a_r \in \mathbb{R}, \phi(\cdot): \text{ReLU}$
- Initialization: $a_r \sim \text{unif}\{1, -1\}$ *only for simplicity*
 $w_r \sim \mathcal{N}(0, I)$
- Training: only train w_1, \dots, w_m

$$\min_{w_1, \dots, w_m} \frac{1}{n} \sum_{i=1}^n f(x_i; a, w) - y_i)^2$$

$$\frac{dU(t)}{dt} = -\frac{1}{n} H(t) \cdot (U(t) - y)$$

Idea: show $H(t) \approx H^*$, $H^*_{ij} = \lim_{m \rightarrow \infty} \mathbb{E}_{\text{init}} \left[\frac{\partial f(x_i, w)}{\partial w_i} \frac{\partial f(x_j)}{\partial w_j} \right]$

Gradient Flow: a Kernel Point of View

$$H_{ij}(t) = \left\langle \frac{\partial U_i(t)}{\partial w(t)}, \frac{\partial U_j(t)}{\partial w(t)} \right\rangle, \quad w \in \mathbb{R}^{m \times d}$$

$$= \sum_{r=1}^m \left\langle \frac{\partial U_i(t)}{\partial w_r(t)}, \frac{\partial U_j(t)}{\partial w_r(t)} \right\rangle$$

$$\frac{\partial U_i(t)}{\partial w_r(t)} = \frac{1}{\sum_{k=1}^m} a_k \cdot x_i \cdot \mathbb{1}\{w_r^T x_i \geq 0\}$$

$$H_{ij}(t) = \sum_{r=1}^m \frac{1}{m} \left\langle a_k \cdot x_i \cdot \mathbb{1}\{w_r^T x_i \geq 0\}, a_k \cdot x_j \cdot \mathbb{1}\{w_r^T x_j \geq 0\} \right\rangle$$

$$= \frac{1}{m} x_i^T x_j \sum_{r=1}^m \mathbb{1}\{w_r^T x_i \geq 0, w_r^T x_j \geq 0\}$$

To show $H(t) \approx H^*$:

- (1) $H(0) \approx H^*$
- (2) $H(t) \approx H(0)$

Gradient Flow: a Kernel Point of View

Initialization:

Hoeffding Inequality:

Q.v. z_1, \dots, z_n i.i.d., $|z_i| \leq 1$
 if $n = \Omega\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$, $0 < \delta \leq 1$,
 $0 < \epsilon < 1$

w.p. $1 - \delta$, $\left|\frac{1}{n} \sum_{i=1}^n z_i - \mathbb{E}[z_i]\right| \leq \epsilon$

$$H_{ij}(0) = \underbrace{x_i^T x_j}_{\text{average}} \cdot \underbrace{\frac{1}{m} \sum_{k=1}^m \mathbb{1}\{w_k(0)^T x_i \geq 0, w_k(0)^T x_j \geq 0\}}_{\text{random variable}}$$

$$H_{ij}^* = \mathbb{E}_{w \sim \mathcal{N}(0, I)} x_i^T x_j \cdot \mathbb{1}\{w^T x_i \geq 0, w^T x_j \geq 0\}$$

when m is large enough, $|H_{ij}(0) - H_{ij}^*| \leq \epsilon$

Gradient Flow: a Kernel Point of View

$$\begin{aligned} & \|H^* - H(\theta)\|_F \\ & \leq \sum_{i,j} |H_{ij}^* - H_{ij}(\theta)| \\ & \leq h^2 \sum \quad \text{choose } \sum \text{ small} \end{aligned}$$

First : power of OUPC-parameterization
→ configuration

Gradient Flow: a Kernel Point of View

ε -training error

Want $H(t) \approx H(0)$

for simplicity:

- 1) just train till time t

- 2) $y_i = 0 (1)$

- 3) $\|x_i\|_2 = 1$

$$H_{ij}^* = x_i^T x_j \quad \cdot \quad \frac{\pi - \arccos(x_i^T x_j)}{2\pi}$$

Key idea: every weight vector only moves a little $(\frac{1}{\sqrt{n}})$, lazy training

$$\begin{aligned}
\|w(t) - w(0)\|_2 &= \left\| \int_0^t \frac{dw(\tau)}{d\tau} d\tau \right\|_2 \\
&\leq \int_0^t \left\| \frac{dw(\tau)}{d\tau} \right\|_2 d\tau \\
&= \int_0^t \left\| \frac{1}{\sqrt{m}} \frac{1}{n} \sum_{i=1}^n (u_i(\tau) - y_i) u_i x_i \cdot \mathbb{1} \{w_i^T x_i \geq 0\} \right\|_2 d\tau
\end{aligned}$$

$$\text{If } u_i(\tau) = 0 \quad (1)$$

$$\leq C \cdot \int_0^t \frac{1}{\sqrt{m}} 1 \cdot d\tau$$

$$\leq \frac{C \cdot t}{\sqrt{m}}$$