

Non-convex Optimization Landscape

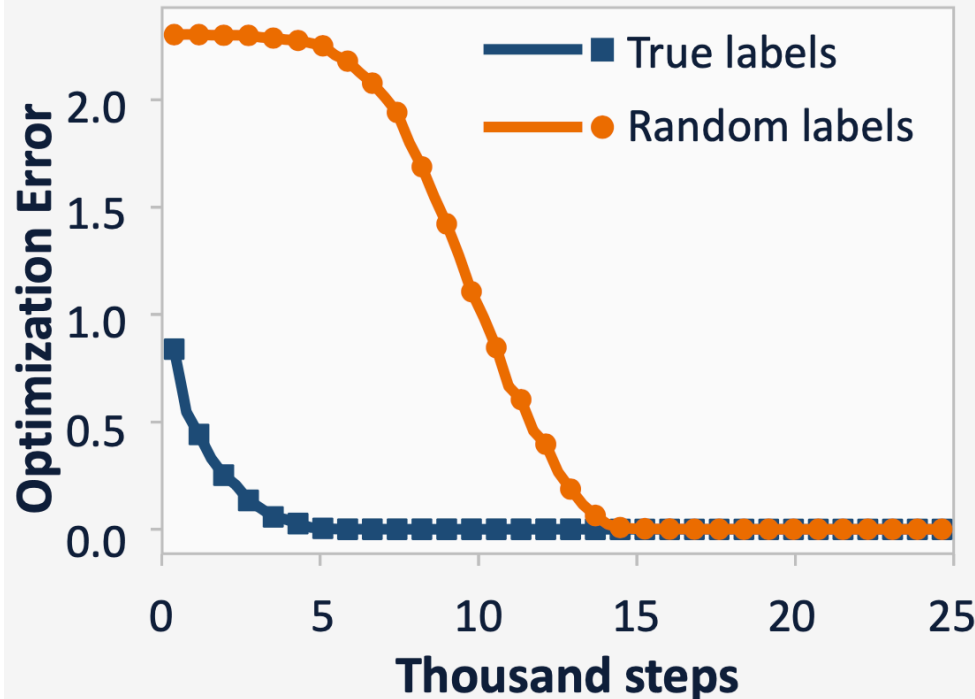
W

Gradient descent finds global minima

Practice: gradient descent

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

parameter
→ U



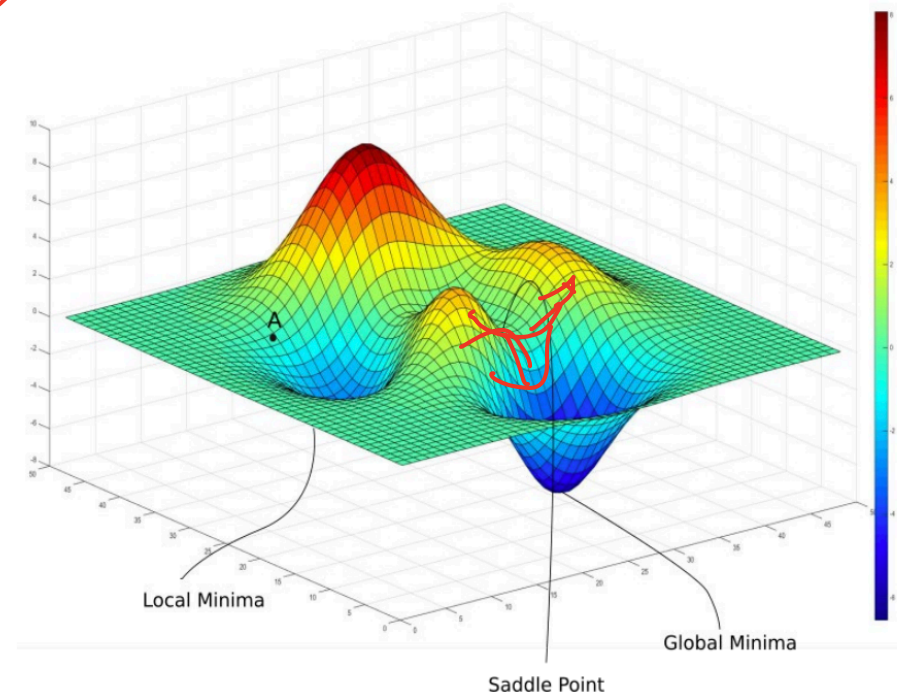
Optimization error $\rightarrow 0$ for both *true labels* and *random labels* !

Zhang Bengio Hardt Recht Vinyals 2017

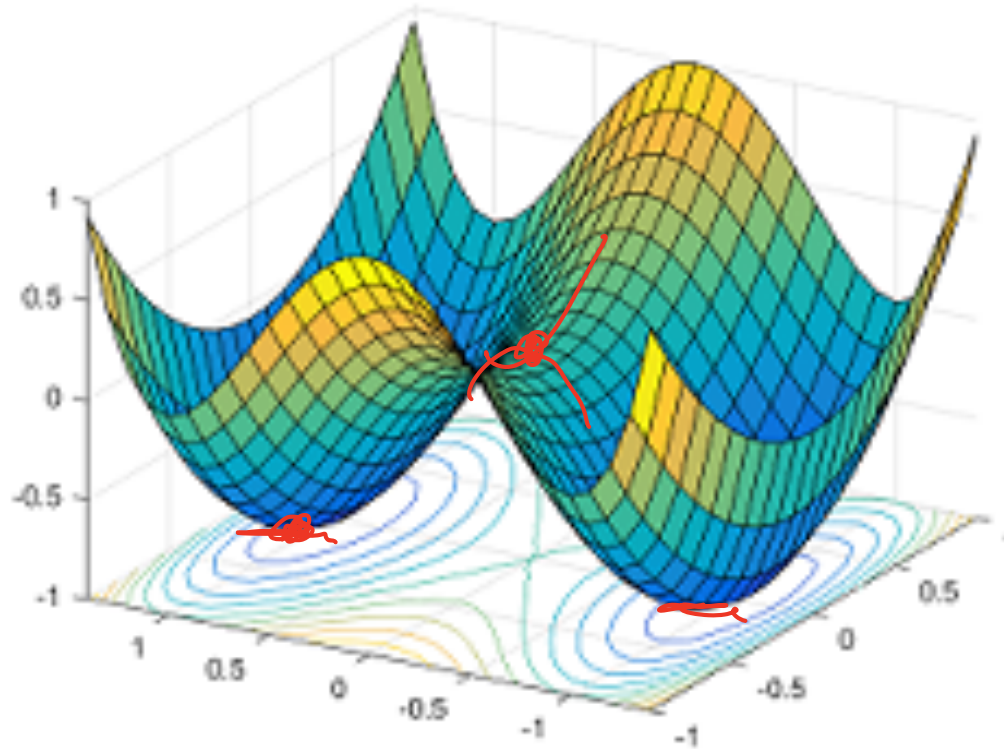
Understanding DL Requires Rethinking Generalization

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



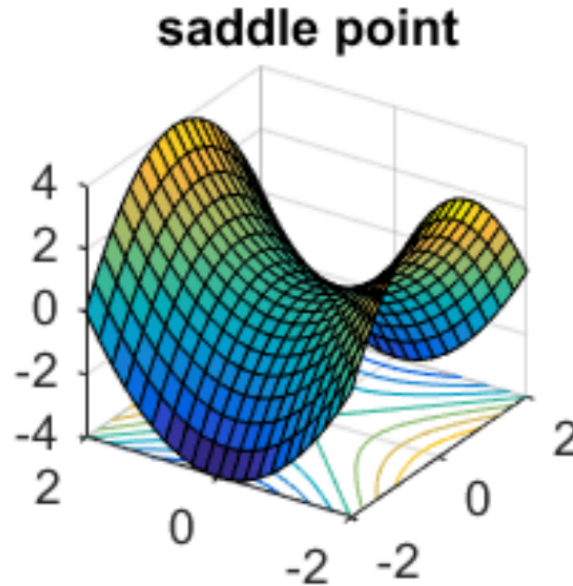
hope

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



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

$$V^T A = \Lambda \text{diag} V$$



$$\nabla f(x) = 0$$

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

$$\min f(x) = \frac{1}{2} x^T A x, \quad \lambda_{\min}(A) < 0, \quad \text{eigen vector } v; \text{ unit vector}$$

$$x=0, \quad \nabla f(x) = Ax = 0, \quad \nabla^2 f(x) = A$$

$$x_0 \text{ close to } 0, \quad x_{t+1} = x_t - \eta \nabla f(x_t)$$

$$\|x_{t+1}\|_2 \approx \left| v^T x_{t+1} \right| = \left| v^T (x_t - \eta A x_t) \right| = \left| v^T x_t - \eta \lambda_{\min}(A) v^T x_t \right| = \underbrace{\left| 1 - \eta \lambda_{\min}(A) \right|}_{> 1} \cdot \left| v^T x_t \right|$$

suppose η sufficiently small

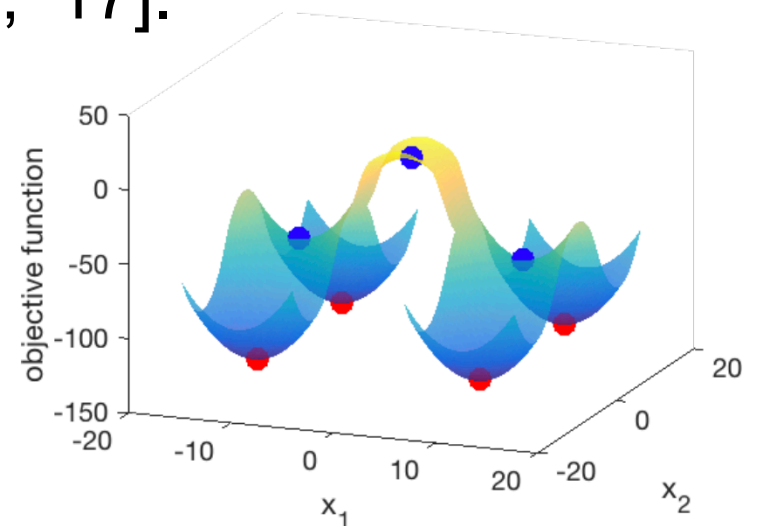
exp ↑

Escaping Strict Saddle Points

$$X_{t+\eta} = X_t - \eta \nabla f(X_t) + \eta \cdot \varepsilon$$
$$\varepsilon \sim \mathcal{N}(0, \Sigma)$$

- **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. *$X_0 \sim$ randomly initialized*
- 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

$$\min_{U, V} \|UV^T - A\|_F^2$$

- Matrix sensing

$$\sum_{i=1}^n (\langle S_i, UV^T \rangle - y_i)^2$$

- Matrix completion

$$\begin{matrix} \begin{matrix} \# & \# \\ \# & \# \end{matrix} \\ \hline \sqrt{\otimes} \end{matrix} \quad \text{low-rank}$$

- Tensor factorization

- Two-layer neural network with quadratic activation

$$f(x_i) = \sum_{j=1}^m \langle x_i, w_j \rangle^2$$

What about neural networks?

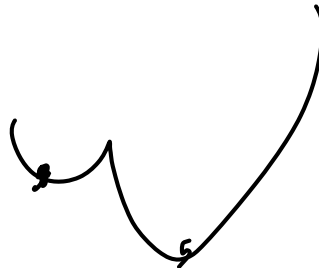
$$x \in \mathbb{R}^d, d \text{ large}$$

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

$$\exists x, \partial \min_y (Df(x)) = 0 \quad \min_{w_1, \dots, w_{L-1}} (w_{L-1} w_{L-2} \dots w_1 x - y)^2, \quad \min_w (w x - y)^2$$

- 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].



Global convergence of gradient descent



Global convergence of gradient descent

(convex loss)

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.

ϵ - global min

$\text{poly}(n) \log\left(\frac{1}{\epsilon}\right)$
for quadratic

Neural Tangent Kernel

proof for a two-layer NN

Gradient Flow: a Kernel Point of View

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(f(\theta, x_i), y_i)$$

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

$$\text{GF: } \frac{d\theta(t)}{dt} = - \frac{\partial L(\theta)}{\partial \theta}$$

if $L(\theta)$ strongly convex, \exists unique θ^* , $\theta(t) \rightarrow \theta^*$
for NN, # of parameters, $\dim(\theta) > n$

we want show, $t \rightarrow \infty$, $f(\theta(t), x_i) \rightarrow y_i$

Gradient Flow: a Kernel Point of View

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

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

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

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

$$= l'(u_i(t), y_i)$$

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

$$\left(\left\langle \frac{\partial u_1(t)}{\partial \theta(t)}, \frac{\partial u_1(t)}{\partial \theta(t)} \right\rangle, \dots, \left\langle \frac{\partial u_n(t)}{\partial \theta(t)}, \frac{\partial u_n(t)}{\partial \theta(t)} \right\rangle \right)$$

$$[H(t)]_{ij}$$

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

$$\frac{du(t)}{dt} = -\frac{1}{n} H(t) \cdot l'(u(t), y)$$

Gradient Flow: a Kernel Point of View

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

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

$$\frac{d(u(t), y)}{dt} = -\frac{1}{\eta} H(t) (u(t) - y)$$

If $H(t)$ is always positive definite
 $\forall t, \lambda_{\min}(H(t)) \geq \lambda_0, \lambda_0 > 0$

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

$$\text{Of: } \frac{d\left(\frac{1}{2} \|u(t) - y\|_2^2\right)}{dt} = -\frac{1}{\eta} (u(t) - y)^T H(t) (u(t) - y) \leq -\frac{\lambda_0}{\eta} \|u(t) - y\|_2^2$$

H p.d.

$\lambda_{\min}(H) \geq \lambda_0$
any vector \checkmark

$u^T H u$
 $\geq \lambda_0 \|u\|_2^2$

Gradient Flow: a Kernel Point of View

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

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

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

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

$$\Rightarrow \frac{1}{2} \|u(t) - y\|_2^2 \leq C \cdot \exp\left(-\frac{\lambda_0 t}{\eta}\right)$$

$$\log\left(\frac{1}{\epsilon}\right)$$

$$t \rightarrow \infty, \text{ loss} \rightarrow 0, u(t) - y$$

Gradient Flow: a Kernel Point of View

$$f(\theta, x) = \frac{1}{\sqrt{m}} \sum_{r=1}^m a_r \cdot \sigma(W_r^T x),$$

m : width, $x \in \mathbb{R}^d$, $a_r \in \mathbb{R}$, $W_r \in \mathbb{R}^d$, $\sigma(-)$: ReLU

• Initialization: $a_r \sim \text{unit } \{1, -1\}$ for simplicity

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

• Training: only train w_1, \dots, w_m

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

$$u_i(t) = f(x_i, a, w(t))$$

$$\frac{du(t)}{dt} = -\frac{1}{n} H(t) (u(t) - y)$$

H^* : NTK

Idea: $H(t)$ stays the same for $\forall t$

$$H(t) \approx H^*$$

$$H_{ij}^* = \lim_{m \rightarrow \infty} \frac{1}{m} \sum_{r=1}^m \left\langle \frac{\partial f(x_i, a, w_r)}{\partial \theta}, \frac{\partial f(x_j, a, w_r)}{\partial \theta} \right\rangle$$

Gradient Flow: a Kernel Point of View

$$H_{ij}'(t) = \left\langle \frac{\partial u_i(t)}{\partial w(t)}, \frac{\partial u_j(t)}{\partial w(t)} \right\rangle, \quad w \in \mathcal{D}^{m \times d}$$

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

$$\frac{\partial u_i(t)}{\partial w_r(t)} = \frac{1}{\sum_m} a_r \cdot x_i = \mathbb{1}_{\{w_r^T x_i \geq 0\}}$$

$$H_{ij}'(t) = \sum_{r=1}^m \frac{1}{m} \mathbb{1}_{\{w_r^T x_i \geq 0\}} \mathbb{1}_{\{w_r^T x_j \geq 0\}}$$

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

To show: $H(t) \approx H^*$, (1) $H(0) \approx H^*$
 (2) $H(t) \approx H(0)$, $\forall t$