

Important Techniques in Neural Network Training



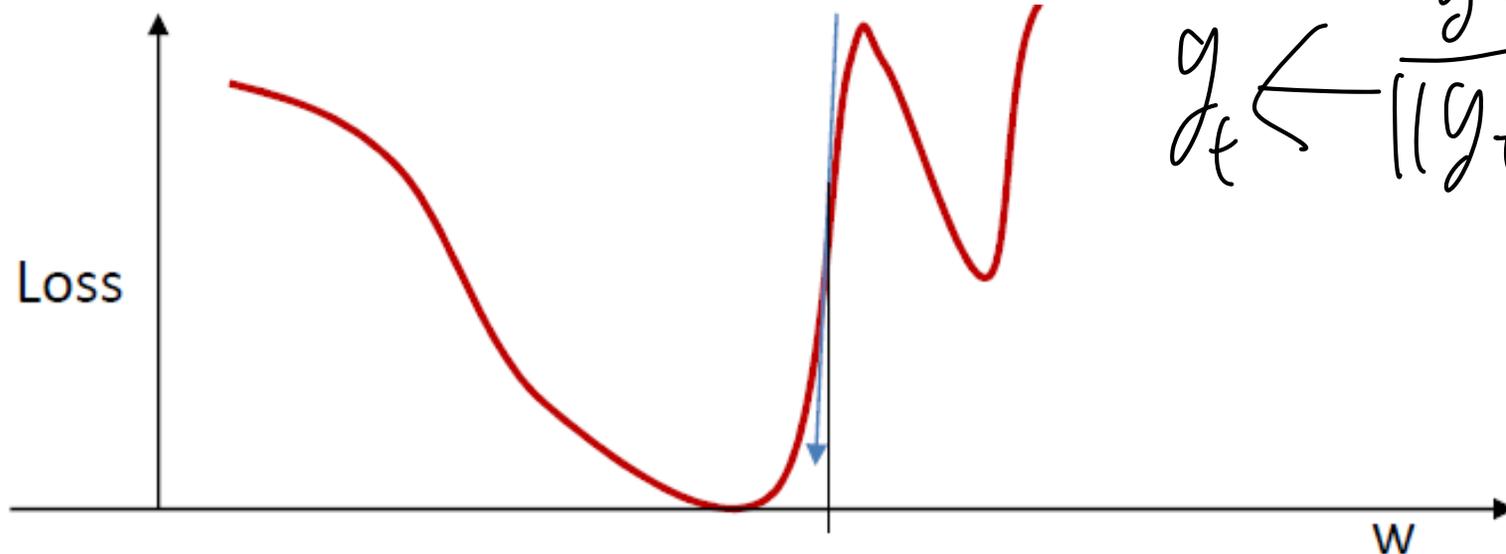
Gradient Clipping

$$X_t = X_t - \eta \nabla f(x_t)$$
$$g_t = \nabla f(x_t)$$

- The loss can occasionally lead to a steep descent
- This result in immediate instability
- If gradient norm bigger than a threshold, set the gradient to the threshold.

if $\|g_t\|_2 > \text{threshold}$

$$g_t \leftarrow \frac{g_t}{\|g_t\|_2} \cdot \text{threshold}$$



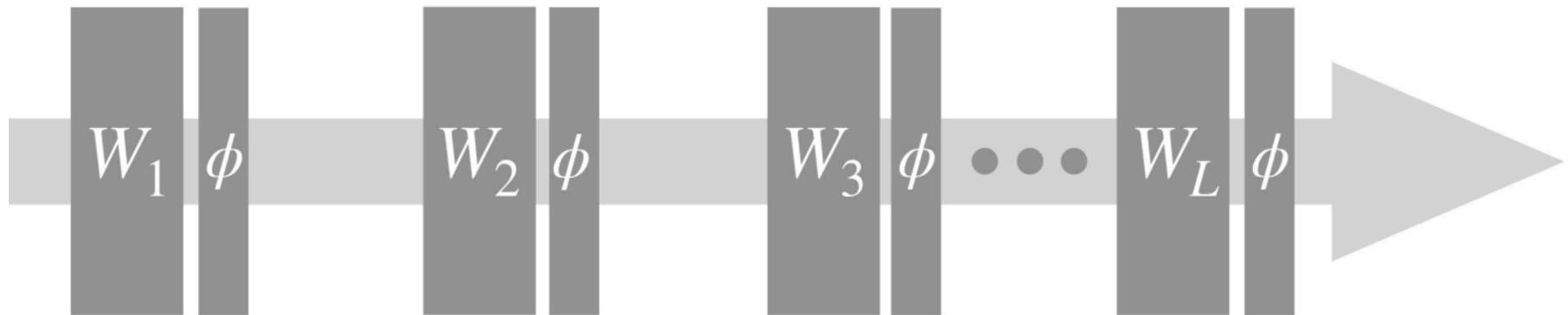
Batch Normalization (Ioffe & Szegedy, '14)

$x \rightarrow \frac{x - \mu}{\sigma}$ (mean = 0, variance = 1)

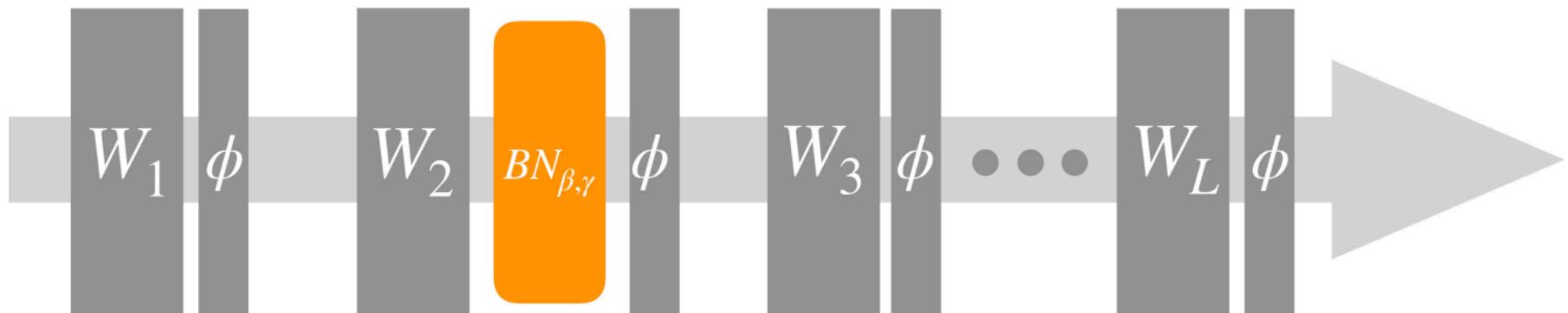
- **Normalizing/whitening** (mean = 0, variance = 1) the inputs is generally useful in machine learning.
 - Could normalization be useful at the level of hidden layers?
 - **Internal covariate shift**: the calculations of the neural networks change the distribution in hidden layers even if the inputs are normalized
- **Batch normalization** is an attempt to do that:
 - Each unit's **pre-activation** is normalized (mean subtraction, std division)
 - During training, mean and std is computed for each minibatch (can be backproped!)

Batch Normalization (Ioffe & Szegedy, '14)

Standard Network

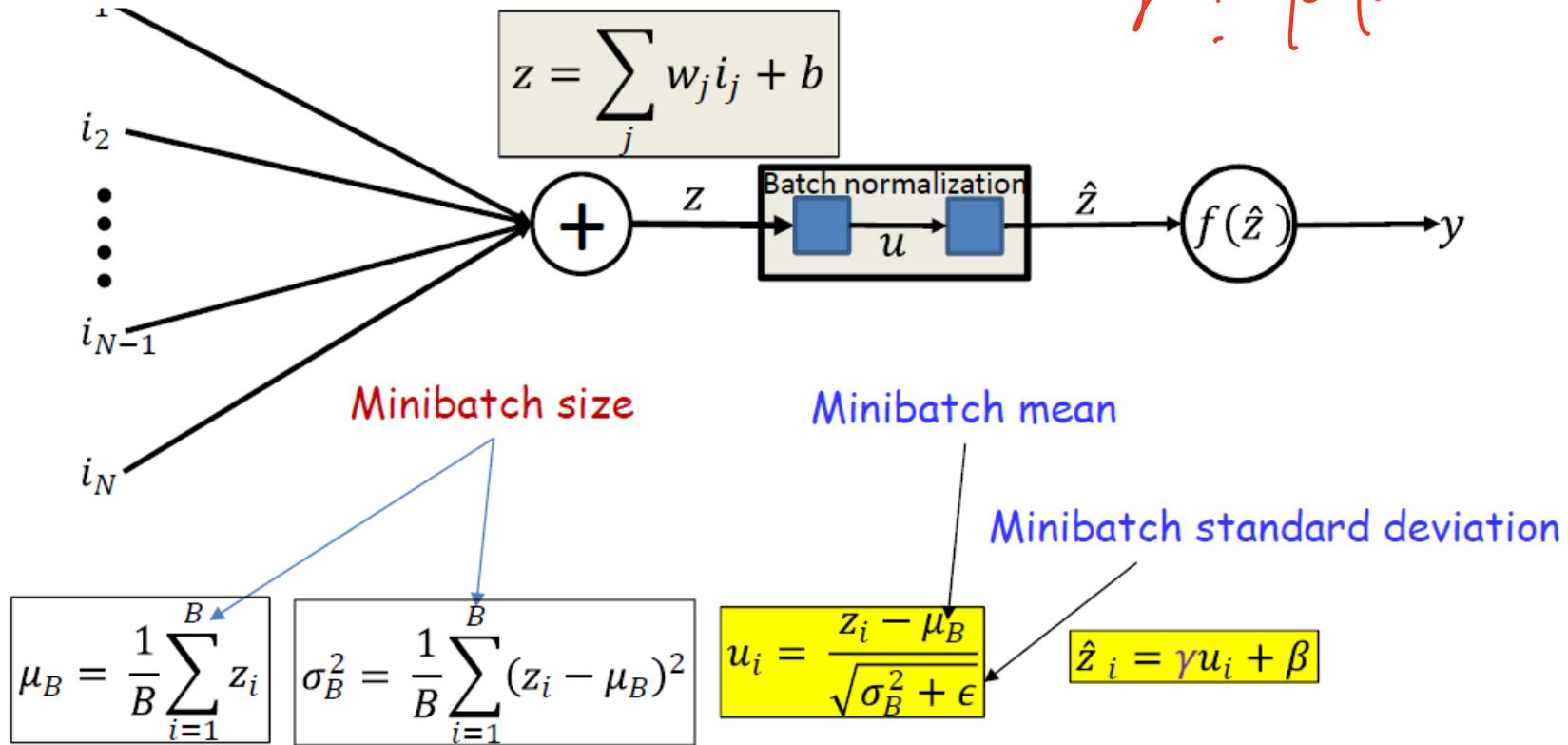


Adding a BatchNorm layer (between weights and activation function)



Batch Normalization (Ioffe & Szegedy, '14)

β : population mean
 γ : population std



z_1, \dots, z_B
 B : batch size

ϵ : to prevent denominator > 0

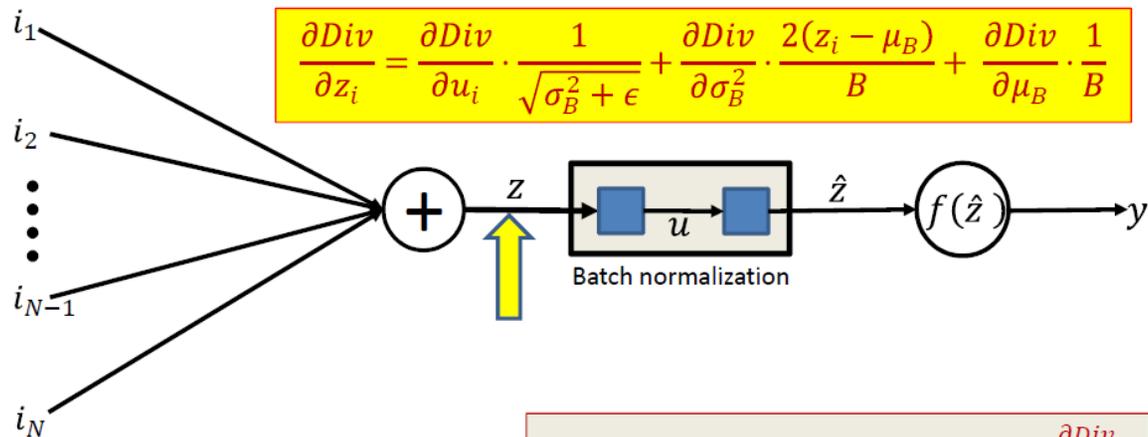
Batch Normalization (Ioffe & Szegedy, '14)

- BatchNorm at training time
 - Standard backprop performed for each single training data
 - Now backprop is performed over entire batch.

β, γ

$$\frac{\partial Div}{\partial \sigma_B^2} = \frac{-1}{2} (\sigma_B^2 + \epsilon)^{-3/2} \sum_{i=1}^B \frac{\partial Div}{\partial u_i}$$

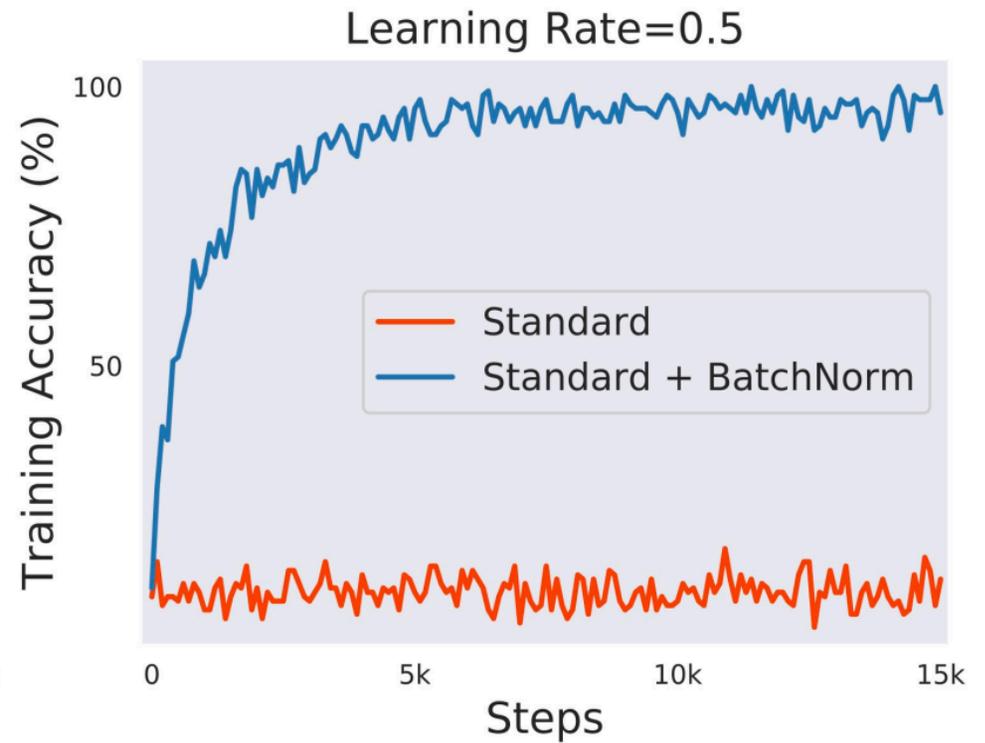
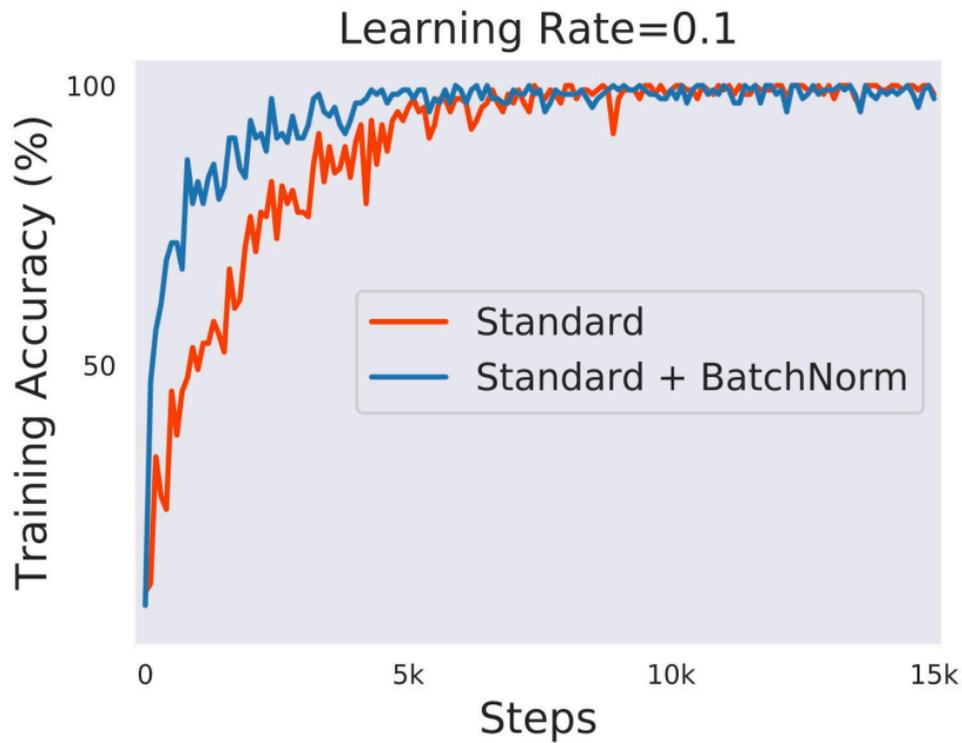
$$\frac{\partial Div}{\partial \mu_B} = \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \sum_{i=1}^B \frac{\partial Div}{\partial u_i}$$



$$\frac{\partial Div}{\partial z_i} = \frac{\partial Div}{\partial u_i} \cdot \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial Div}{\partial \sigma_B^2} \cdot \frac{2(z_i - \mu_B)}{B} + \frac{\partial Div}{\partial \mu_B} \cdot \frac{1}{B}$$

The rest of backprop continues from $\frac{\partial Div}{\partial z_i}$

Batch Normalization (Ioffe & Szegedy, '14)



What is BatchNorm actually doing?

- May not be due to covariate shift (Santurkar et al. '18):
 - Inject non-zero mean, non-standard covariance Gaussian noise after BN layer: removes the whitening effect
 - Still performs well.
- Only training β, γ with random convolution kernels gives non-trivial performance (Frankle et al. '20)
- BN can use exponentially increasing learning rate! (Li & Arora '19)

constant

$$\frac{1}{\sqrt{t}}$$

More normalizations

- Layer normalization (Ba, Kiros, Hinton, '16)
 - Batch-independent
 - Suitable for RNN, MLP
- Weight normalization (Salimans, Kingma, '16)
 - Suitable for meta-learning (higher order gradients are needed)
- Instant normalization (Ulyanov, Vedaldi, Lempitsky, '16)
 - Batch-independent, suitable for generation tasks
- Group normalization (Wu & He, '18)
 - Batch-independent, improve BatchNorm for small batch size

NP-hard

Non-convex

Optimization Landscape

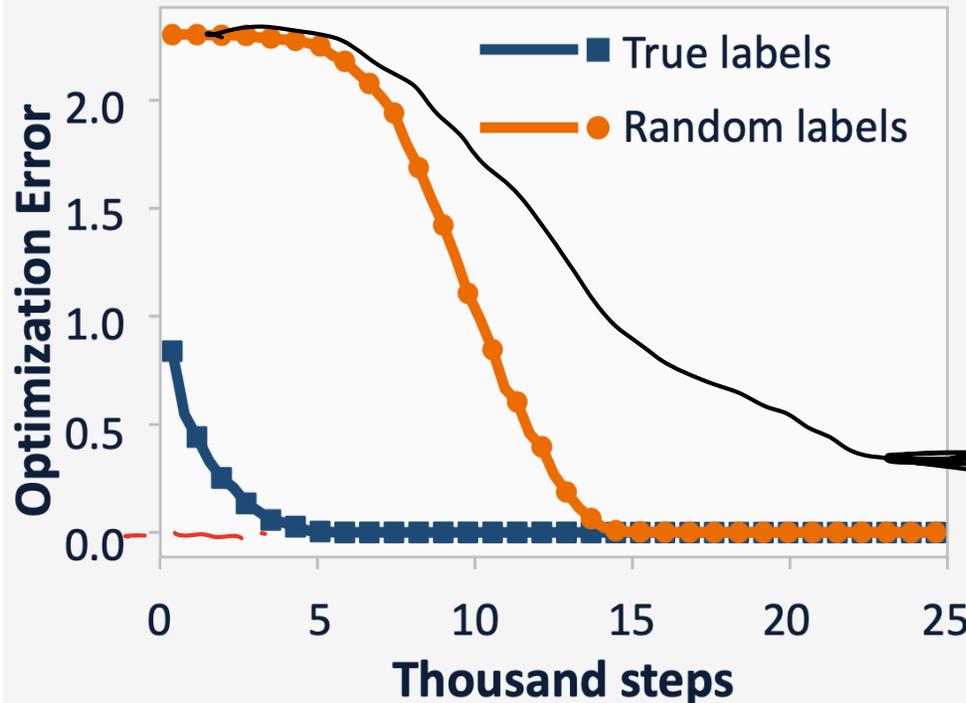
W

Gradient descent finds global minima

$$\min L(\theta)$$

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** !

\otimes overparameterization

NN

$\dim(\theta)$

$\gg n$

n : # of data

Zhang Bengio Hardt Recht Vinyals 2017

Understanding DL Requires Rethinking Generalization

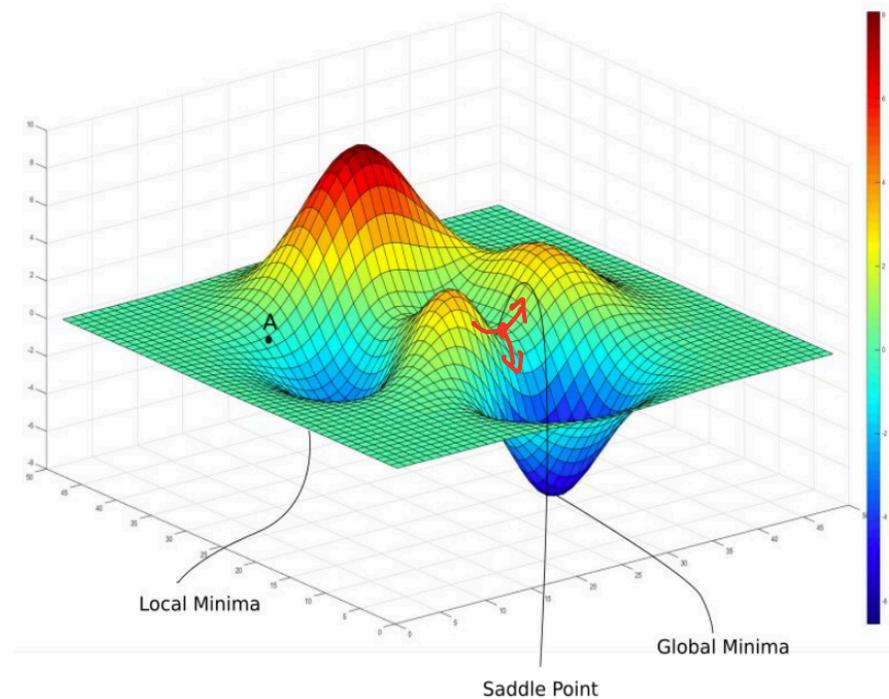
Types of stationary points

- 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

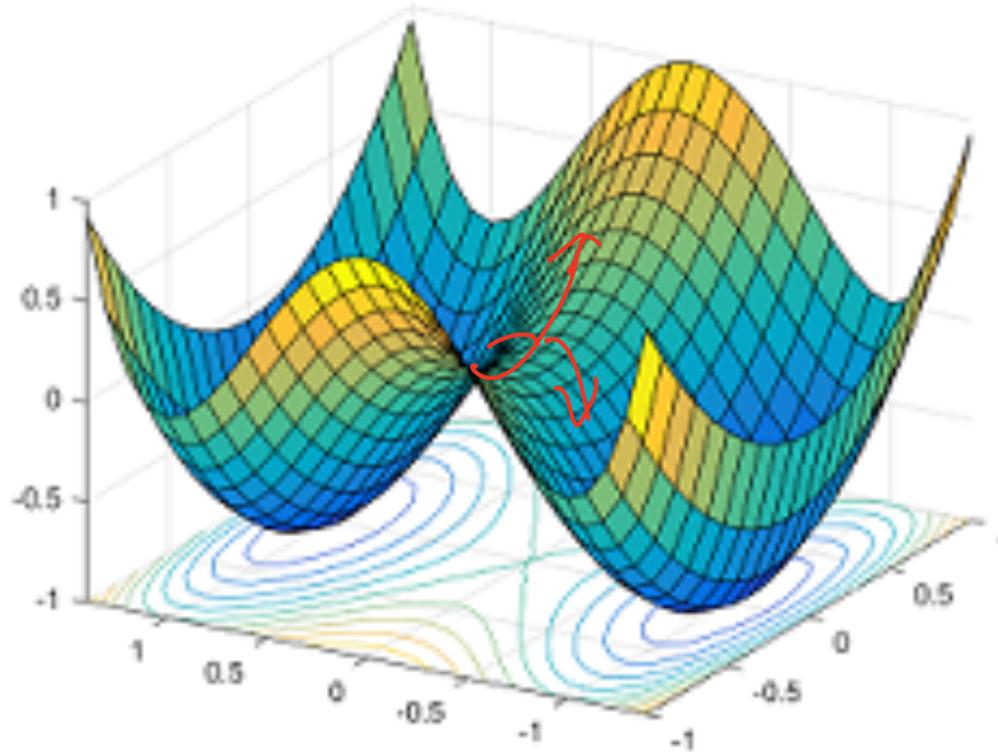
$$\nabla^2 f(x)$$

$$f(x)$$

$$\| \nabla f(x) \| \leq \epsilon, o\left(\frac{1}{\epsilon}\right)$$



Landscape Analysis

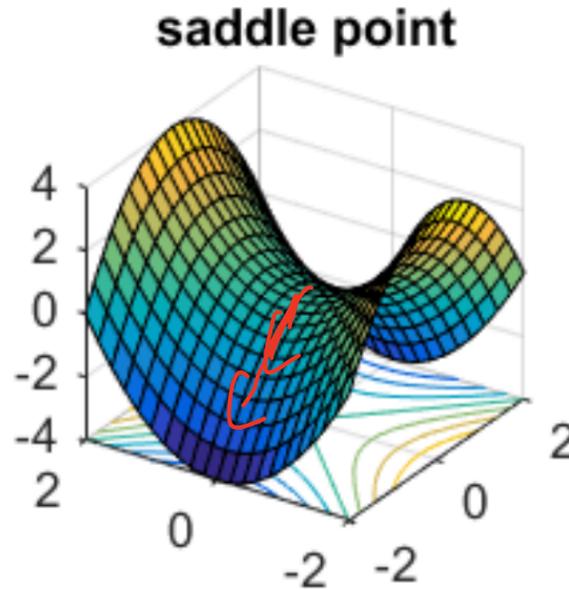


- All local minima are global!
- Gradient descent can escape saddle points.

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

$$f(x) = \frac{1}{2} x^T A x$$

$$x = \begin{pmatrix} x^1 \\ x^2 \end{pmatrix}$$



$$A = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$$

$$\text{if } x=0, Ax=0$$

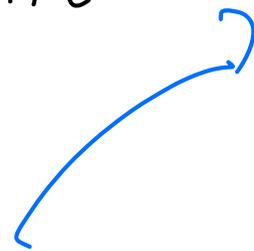
$$\nabla^2 f(x) = A$$

$$\lambda_{\min}(\nabla^2 f(x)) = -1$$

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

$$x_{t+1} = x_t - \eta_t A x_t$$

$$\begin{aligned} x_{t+1}^2 &= x_t^2 + \eta_t x_t^2 \\ &= (I + \eta_t A) x_t^2 \\ &= (I + \eta_t A)^{t+1} x_0^2 \end{aligned}$$



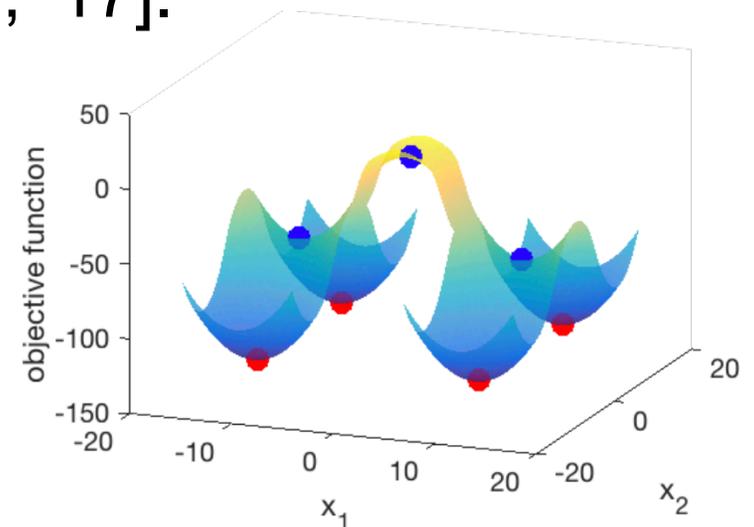
$$\text{if } x_0^2 \neq 0$$

$$|x_{t+1}^2| = (I + \eta_t A)^{t+1} |x_0^2|$$

Escaping Strict Saddle Points

- **Noise-injected** gradient descent can escape strict saddle points in polynomial time [Ge et al., '15, Jin et al., '17].
$$x_{t+1} = x_t - \eta \nabla f(x_t) + \eta \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, I)$$
- Randomly initialized gradient descent can escape all strict saddle points asymptotically [Lee et al., '15]. **X noise**
 - 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

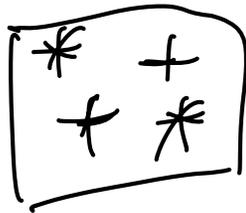
$$\min_{U, V} \|UV^T - X\|_F^2$$

- Matrix sensing

$$y_i = \langle A_i, X \rangle, \quad X: \text{low-rank}$$

$$\min_{U, V} \sum_{i=1}^m (\langle A_i, UV^T \rangle - y_i)^2$$

- Matrix completion



- Tensor factorization

- Two-layer neural network with quadratic activation

$$f(x) = \sum_{j=1}^m \langle W_j, X \rangle^2$$

What about neural networks?

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

$$f(\theta) = W_{H+1} W_H \dots W_1 x \Leftrightarrow f(x) = W^T x$$

$$\min_w \sum_{i=1}^n (w^T x_i - y_i)^2$$

- 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,

$$y = \mathcal{NN}(x) \quad \{(x_i, y_i)\}_{i=1}^n, \quad y_i = \sigma(w^T x_i)$$

$$\min_{w_1, \dots, w_{H+1}} \sum_{i=1}^n (f(x_i) - y_i)^2$$

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

$\exists x$, local but not global

Trajectory-based Analysis for Non-convex optimization



Gradient Flow: a Kernel Point of View

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

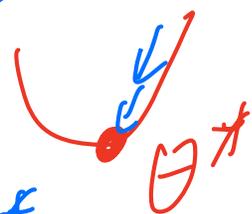
gradient flow:
$$\frac{d\theta(t)}{dt} = - \frac{\partial L(\theta)}{\partial \theta}$$

if $L(\theta)$ is strongly convex, \exists unique θ^*

$$\theta(t) \rightarrow \theta^*$$

for NN, $\dim(\theta) > n$, $\exists \theta^*$

\Rightarrow study $\frac{f(\theta_t, x_i)}{\rightarrow y_i}$ for each i as $t \rightarrow \infty$



Gradient Flow: a Kernel Point of View

$$u_i(t) = f(\theta t, X_i)$$

$$u(t) = \begin{pmatrix} u_1(t) \\ u_2(t) \\ \vdots \\ u_n(t) \end{pmatrix} \in \mathbb{R}^n$$

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

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

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

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

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

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

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

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

$$\frac{d u(t)}{d t} = -\frac{1}{\eta} H(t) l'(u(t), y)$$

Gradient Flow: a Kernel Point of View

If l quadratic, $l(u(t), y) = \frac{1}{2} \|u(t) - y\|_2^2$

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

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

If $H(t)$ is always P.D.
 $\lambda_{\min}(H(t)) > 0$

$$\Rightarrow u(t) \rightarrow y$$