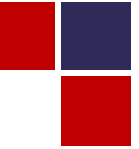

AIL721: Deep Learning

「Instructor: James Arambam」

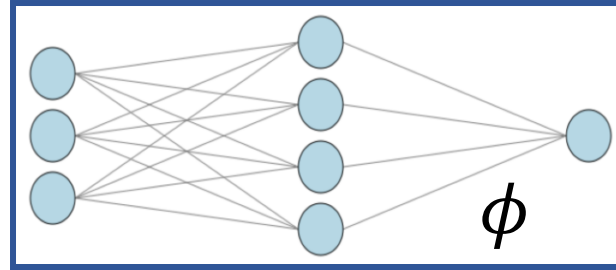


ScAI

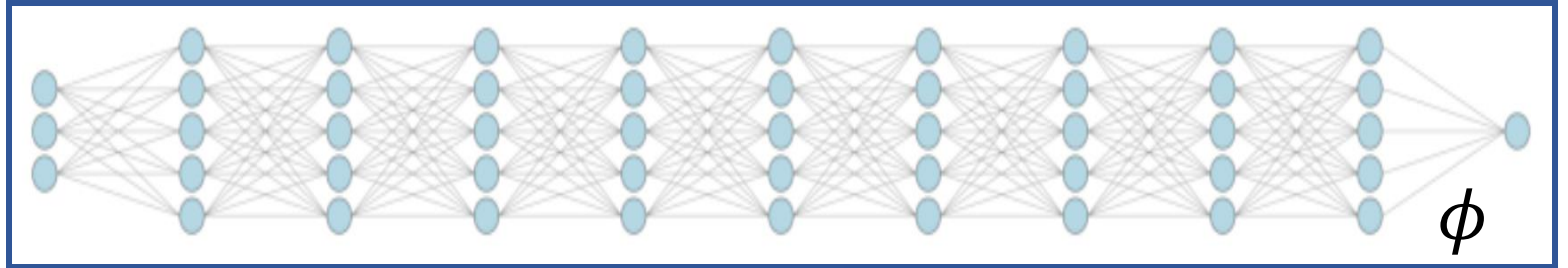
Yardi School of Artificial Intelligence
Indian Institute of Technology Delhi



❑ Shallow Neural Network



❑ Deep Neural Network



❑ Loss Function

**Least Square Error
Loss:**

$$L[\phi] = \sum_{i=1}^I (f[x_i, \phi] - y_i)^2$$

**Negative Loglikelihood
Loss:**

$$L[\phi] = \left[- \sum_{i=1}^I \log \left[Pr(\mathbf{y}_i | \mathbf{f}[\mathbf{x}_i, \phi]) \right] \right]$$

Cross Entropy Loss:

Next Step: Optimization

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[L(\phi) \right]$$

□ Optimization

$$\hat{\phi} = \operatorname{argmin}_{\phi} [L(\phi)]$$

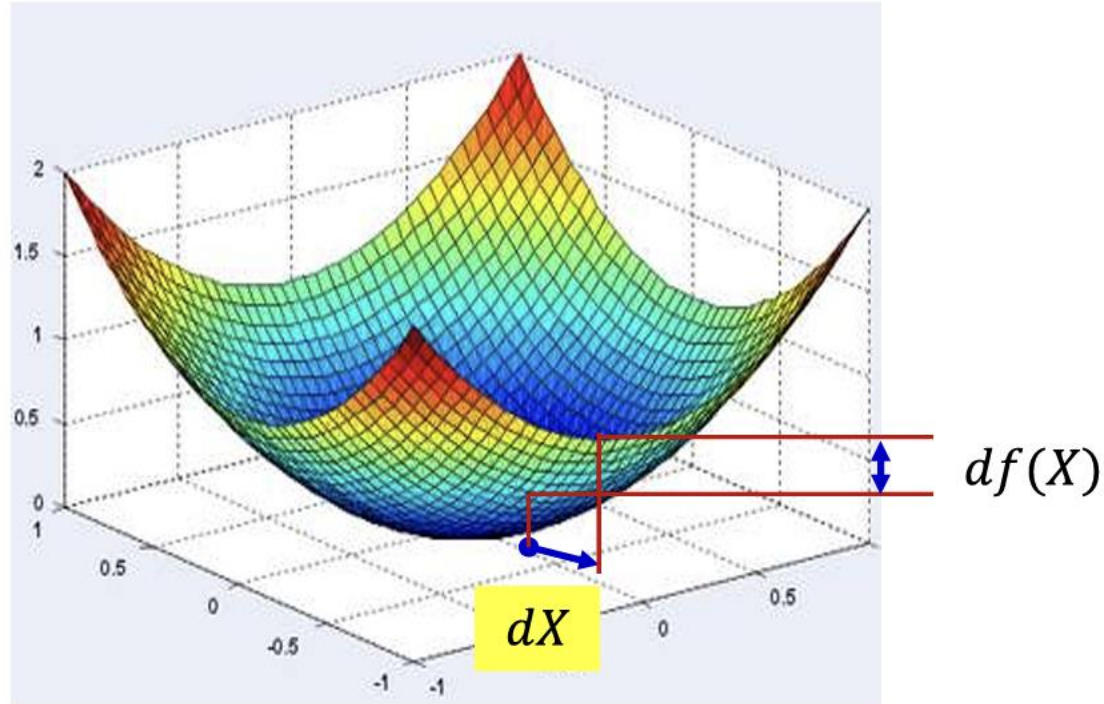
- Find the parameter that **minimizes** the loss function.
- This is known as the **learning or training**.

Training data: $(x_1, y_1), (x_2, y_2), \dots, (x_I, y_I)$

$$\begin{aligned} L[\phi] &= \sum_{i=1}^I (f[x_i, \phi] - y_i)^2 \\ &= \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2 \end{aligned}$$

Gradient

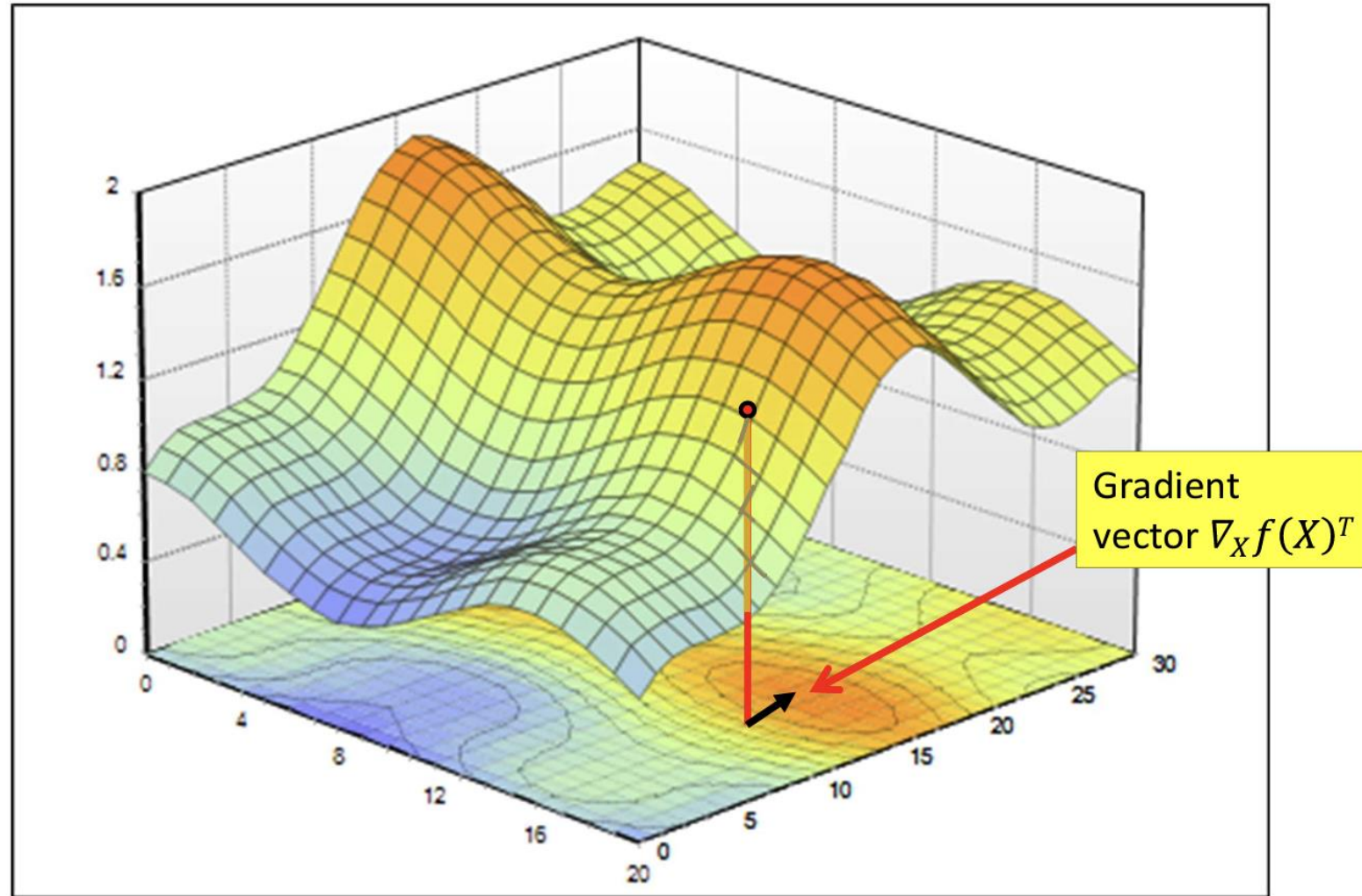
□ Gradient

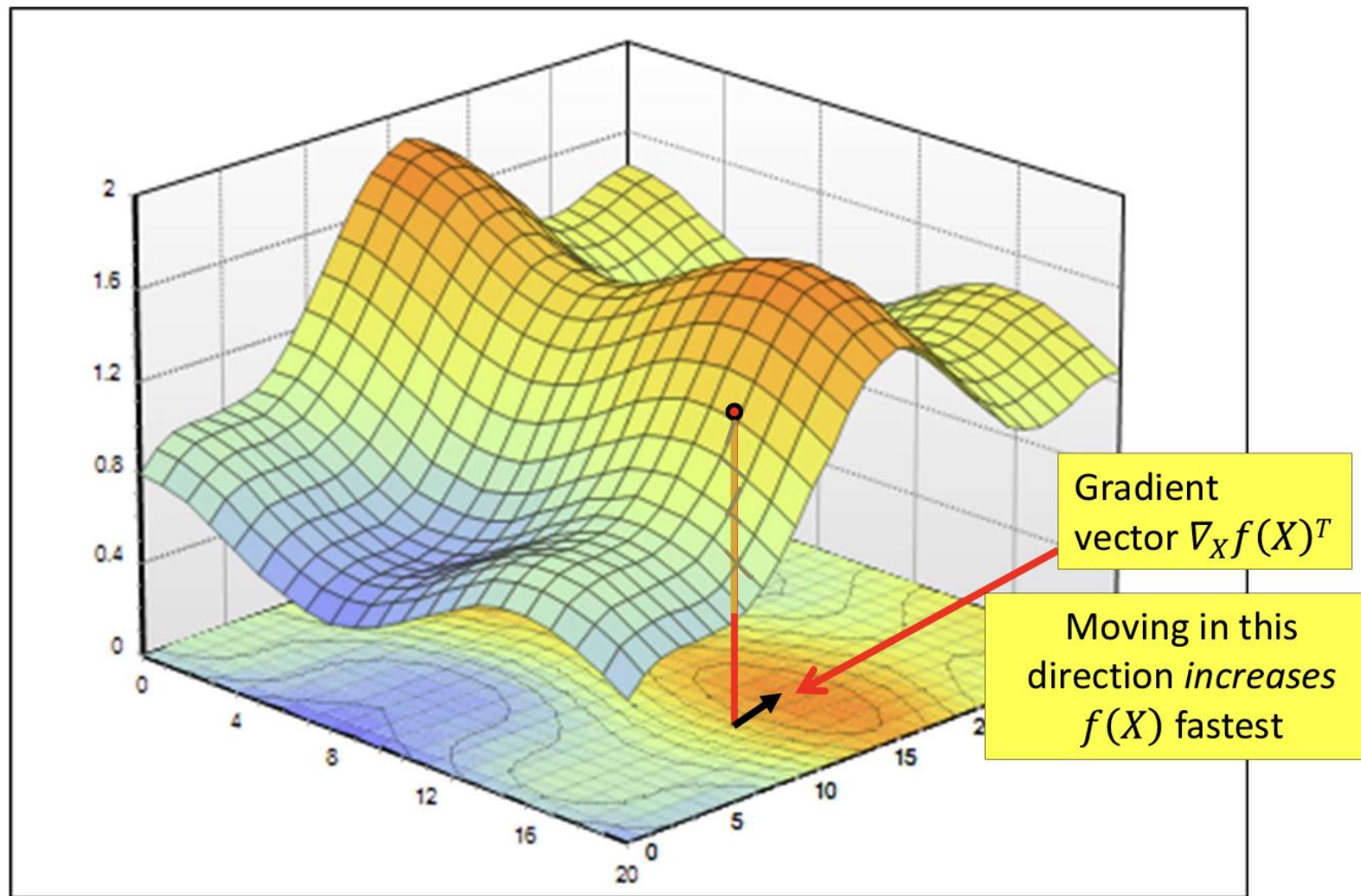


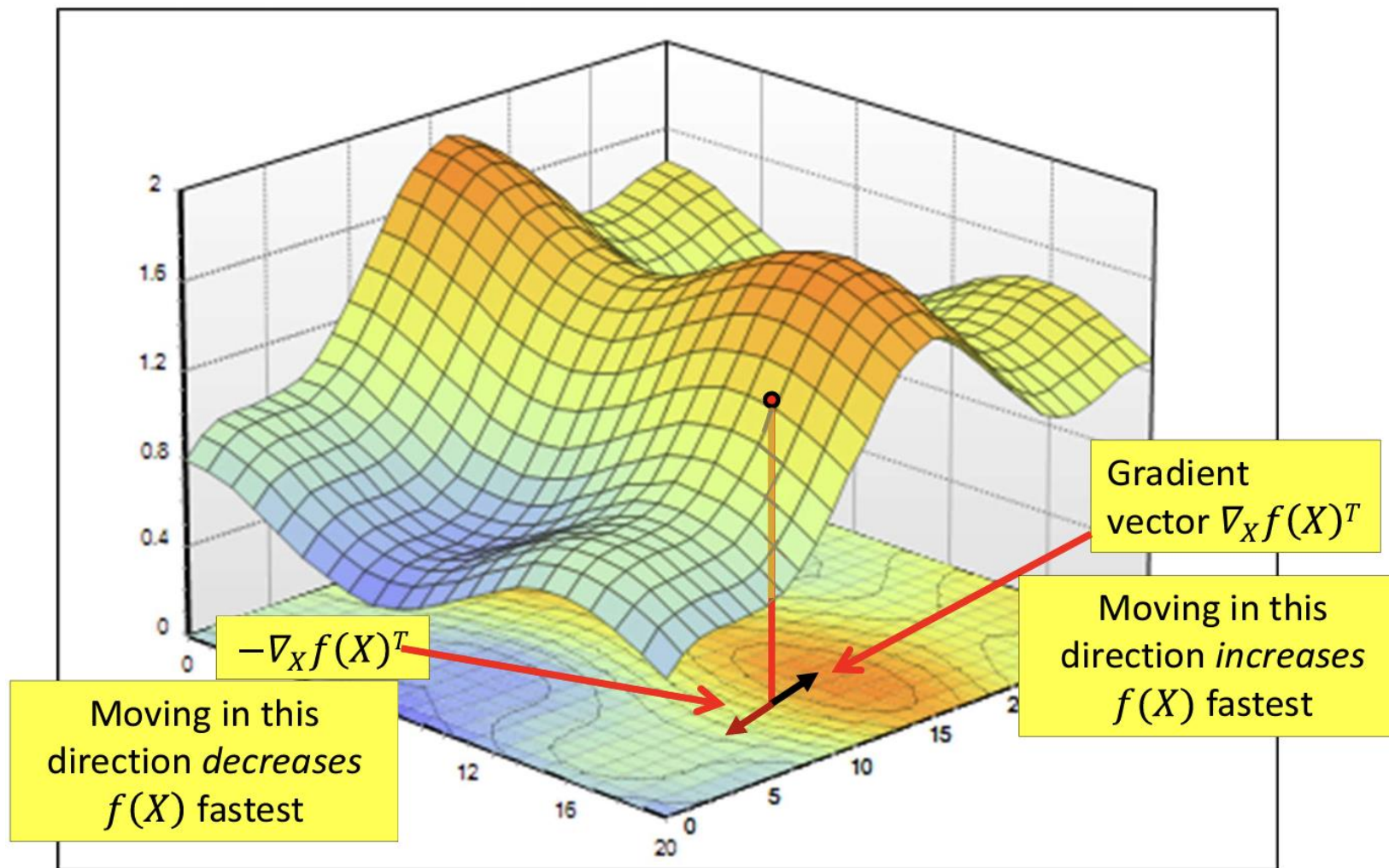
$$- \nabla_X f(X) = \begin{bmatrix} \frac{\partial f(X)}{\partial x_1} & \frac{\partial f(X)}{\partial x_2} & \dots & \frac{\partial f(X)}{\partial x_n} \end{bmatrix}$$

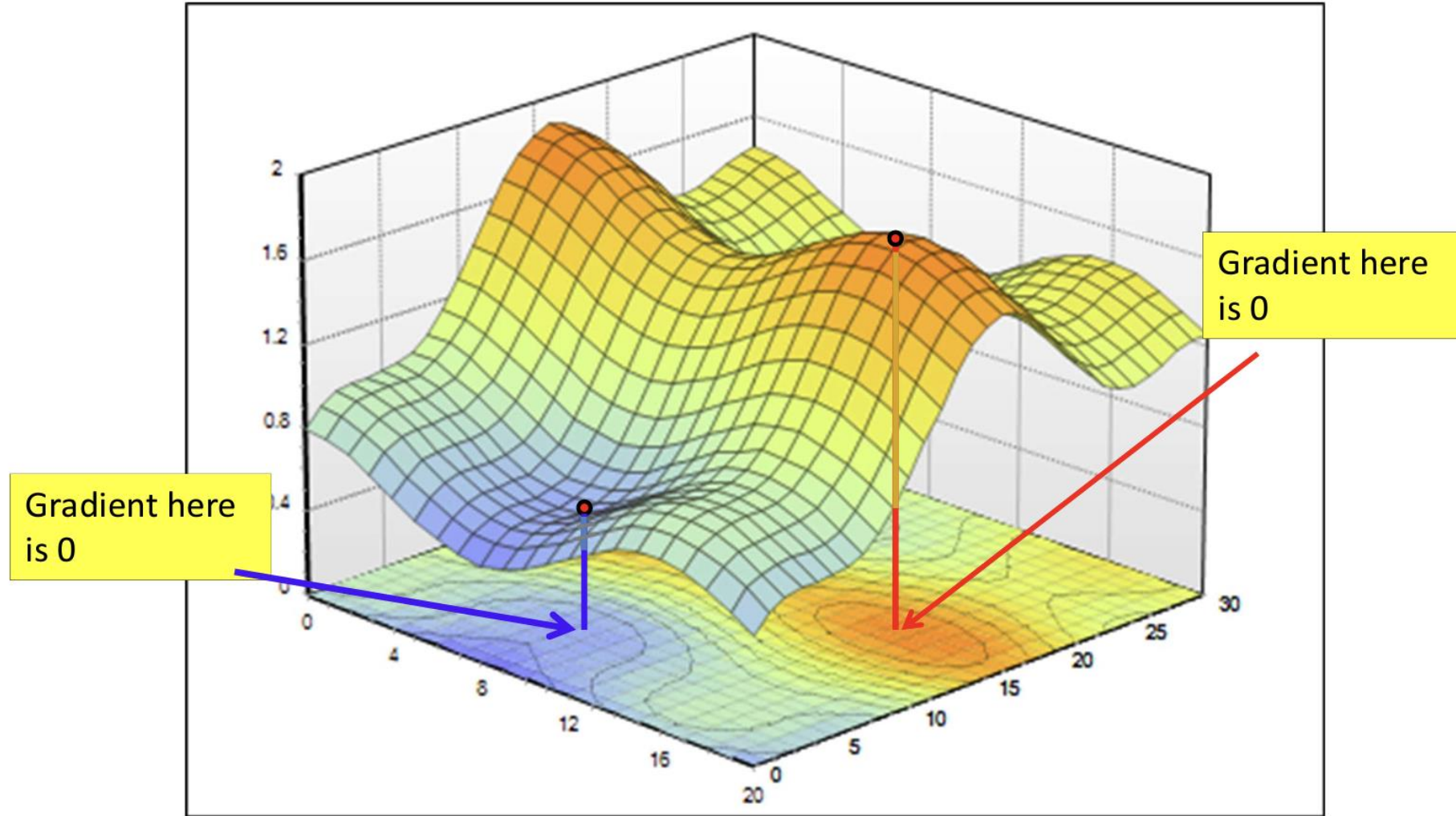
The **gradient** is the transpose of the derivative $\nabla_X f(X)^T$

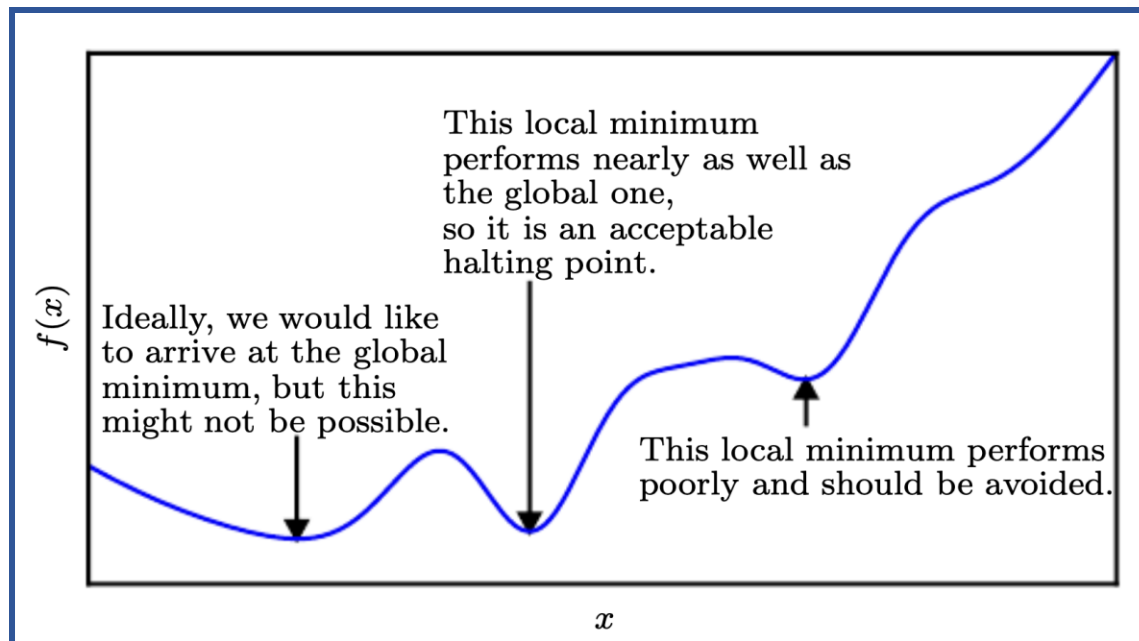
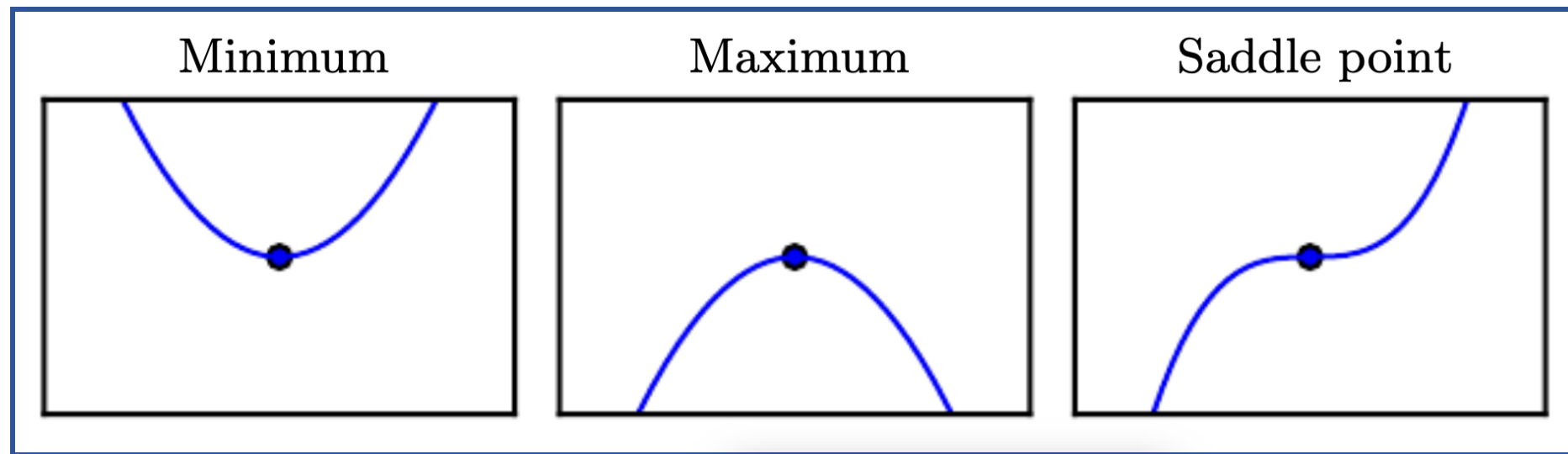
- A column vector of the same dimensionality as X











Gradient = 0



Line search
(like gradient ascent)



Trust region

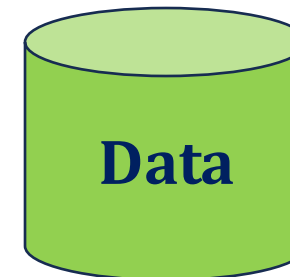


$$\begin{aligned} \max_{s \in \mathbb{R}^n} \quad & m_k(s) \\ \text{s.t.} \quad & \|s\| \leq \delta \end{aligned}$$

Optimization

□ Optimization

$$\hat{\phi} = \operatorname{argmin}_{\phi} [L(\phi)]$$



Step 1: Initialize the parameters ϕ_0

Step 2: Compute the derivatives of loss with respect to the parameters.

$$\frac{\partial L}{\partial \phi} = \begin{bmatrix} \frac{\partial L}{\partial \phi_0} \\ \frac{\partial L}{\partial \phi_1} \\ \vdots \\ \frac{\partial L}{\partial \phi_N} \end{bmatrix}$$

Gradient Descent

Step 3: Update the parameters.

$$\phi \leftarrow \phi - \alpha \cdot \frac{\partial L}{\partial \phi}$$

Learning rate

Optimization

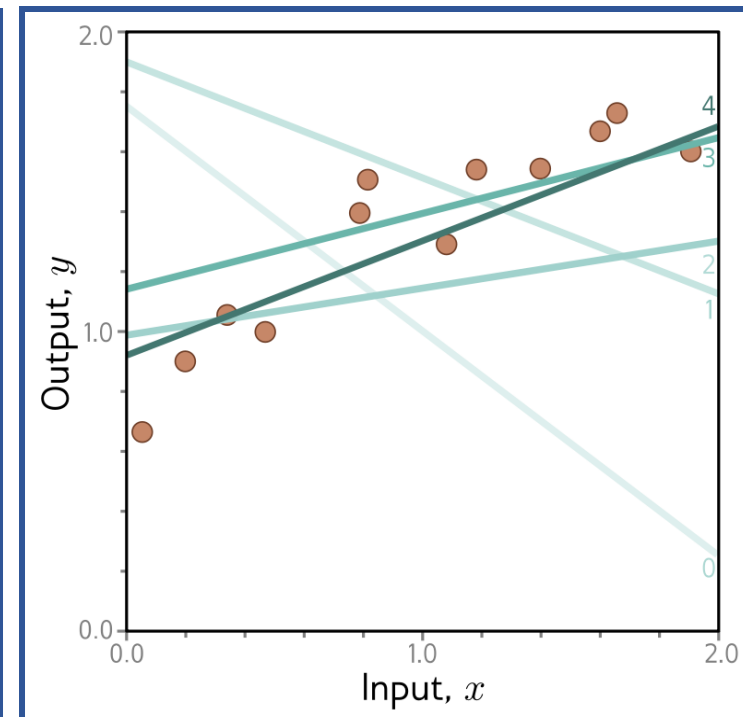
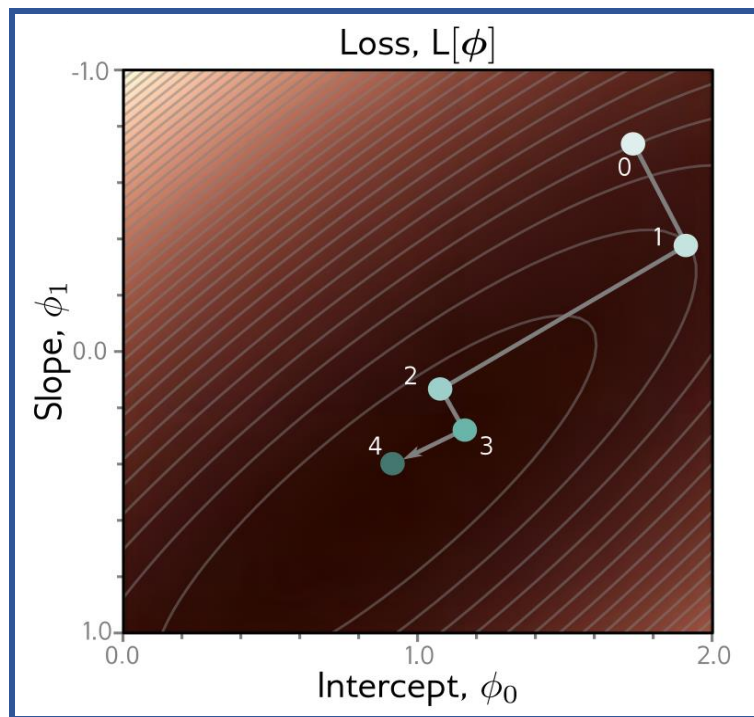
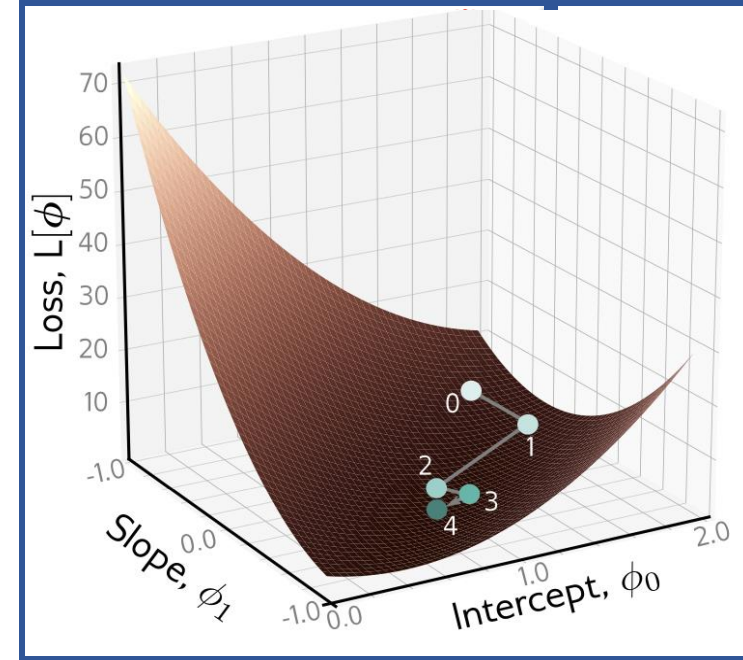
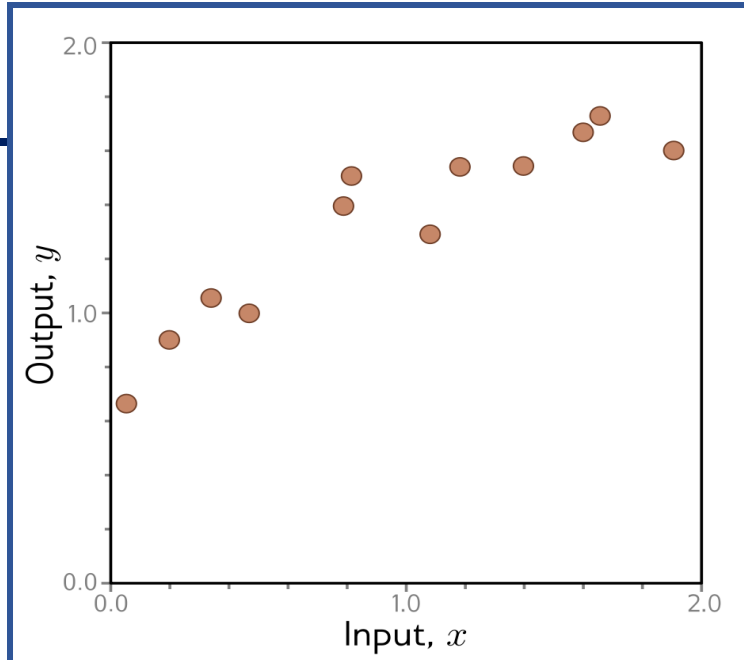
Training data:

$$(x_1, y_1), (x_2, y_2), \dots, (x_I, y_I)$$

$$\begin{aligned} L[\phi] &= \sum_{i=1}^I (f[x_i, \phi] - y_i)^2 \\ &= \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2 \end{aligned}$$

$$\hat{\phi} = \operatorname{argmin}_{\phi} [L(\phi)]$$

Gradient Descent



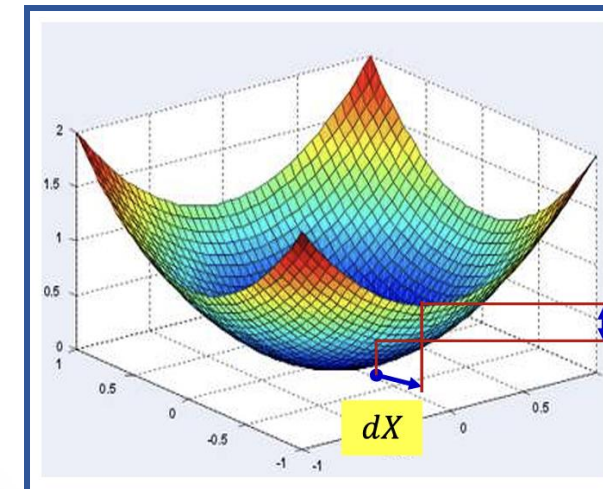
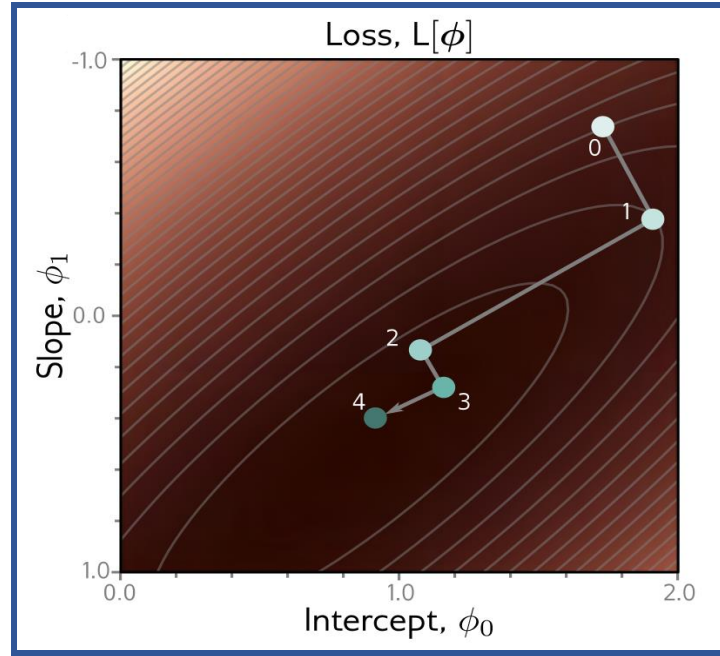
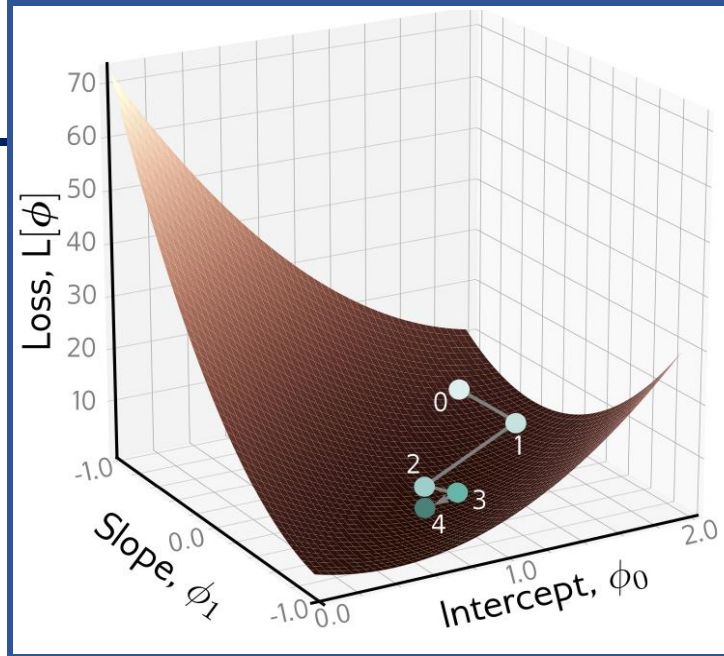
Optimization

Training data:

$$(x_1, y_1), (x_2, y_2), \dots, (x_I, y_I)$$

$$\begin{aligned} L[\phi] &= \sum_{i=1}^I (f[x_i, \phi] - y_i)^2 \\ &= \sum_{i=1}^I (\phi_0 + \phi_1 x_i - y_i)^2 \end{aligned}$$

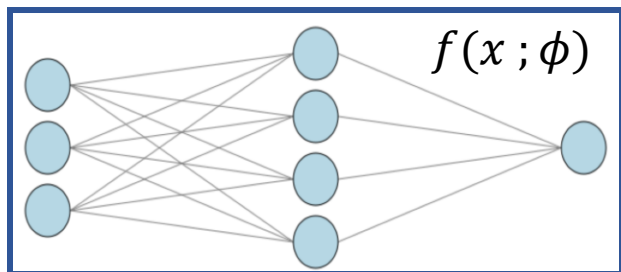
$$\hat{\phi} = \operatorname{argmin}_{\phi} [L(\phi)]$$



Do we reach a global minimum?

Optimization

Non-linear



$$\hat{\phi} = \operatorname{argmin}_{\phi} \left[L(\phi) \right]$$

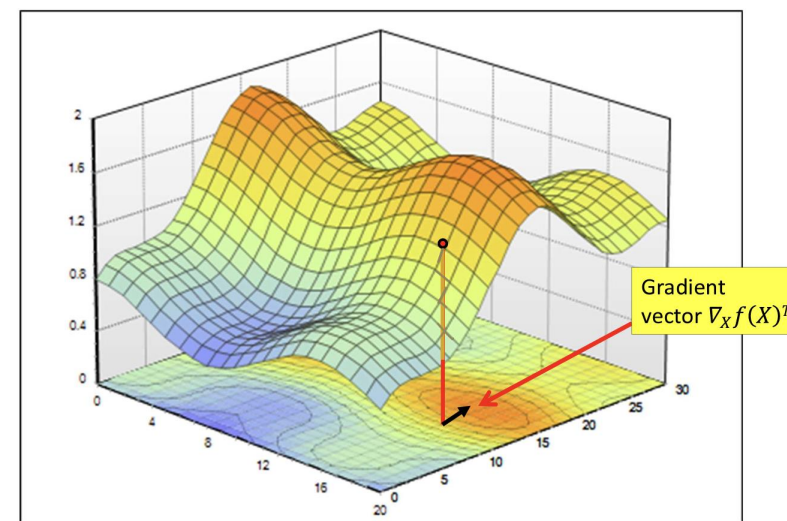


May become a **non-convex optimization**

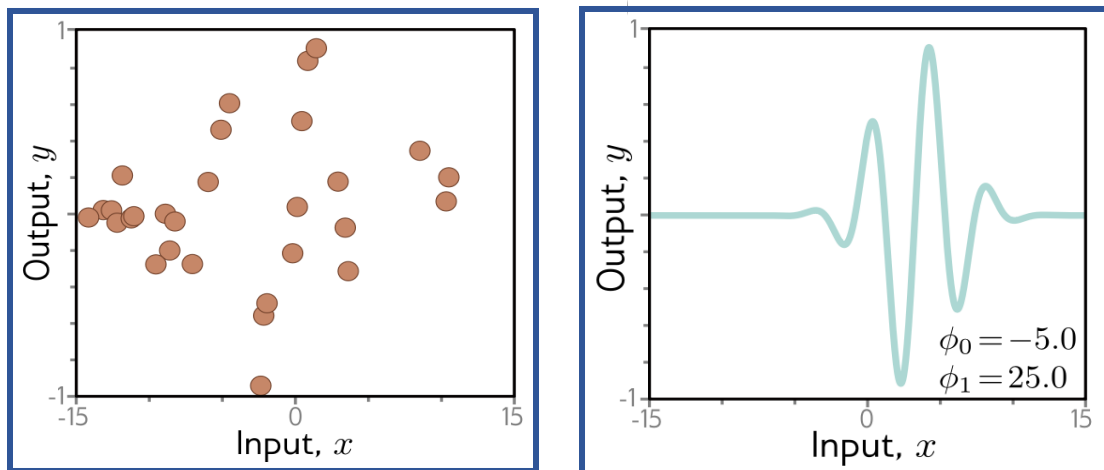
Do we still reach a global minimum?

Note:

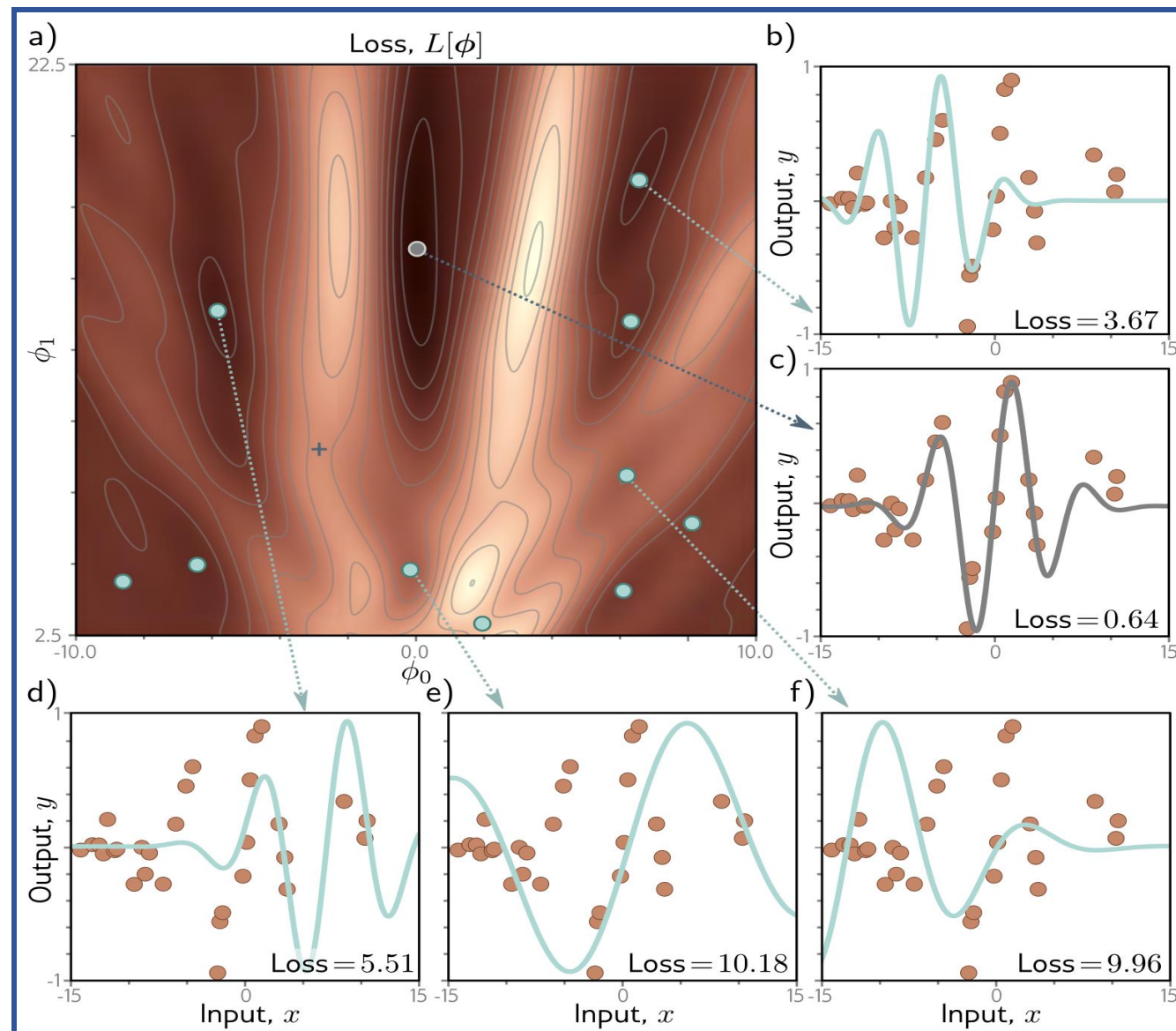
Not all non-linear functions are non-convex.



❑ Non-linear function

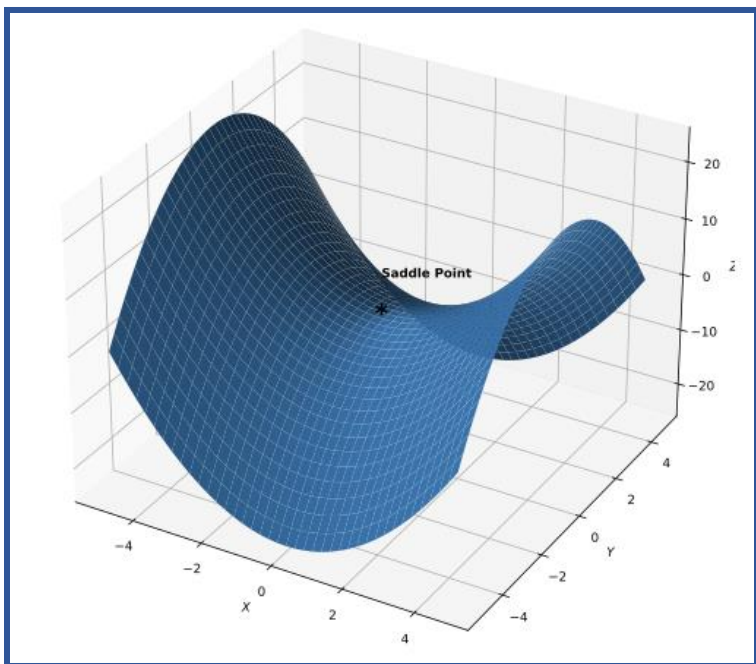


$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} [L(\phi)]$$

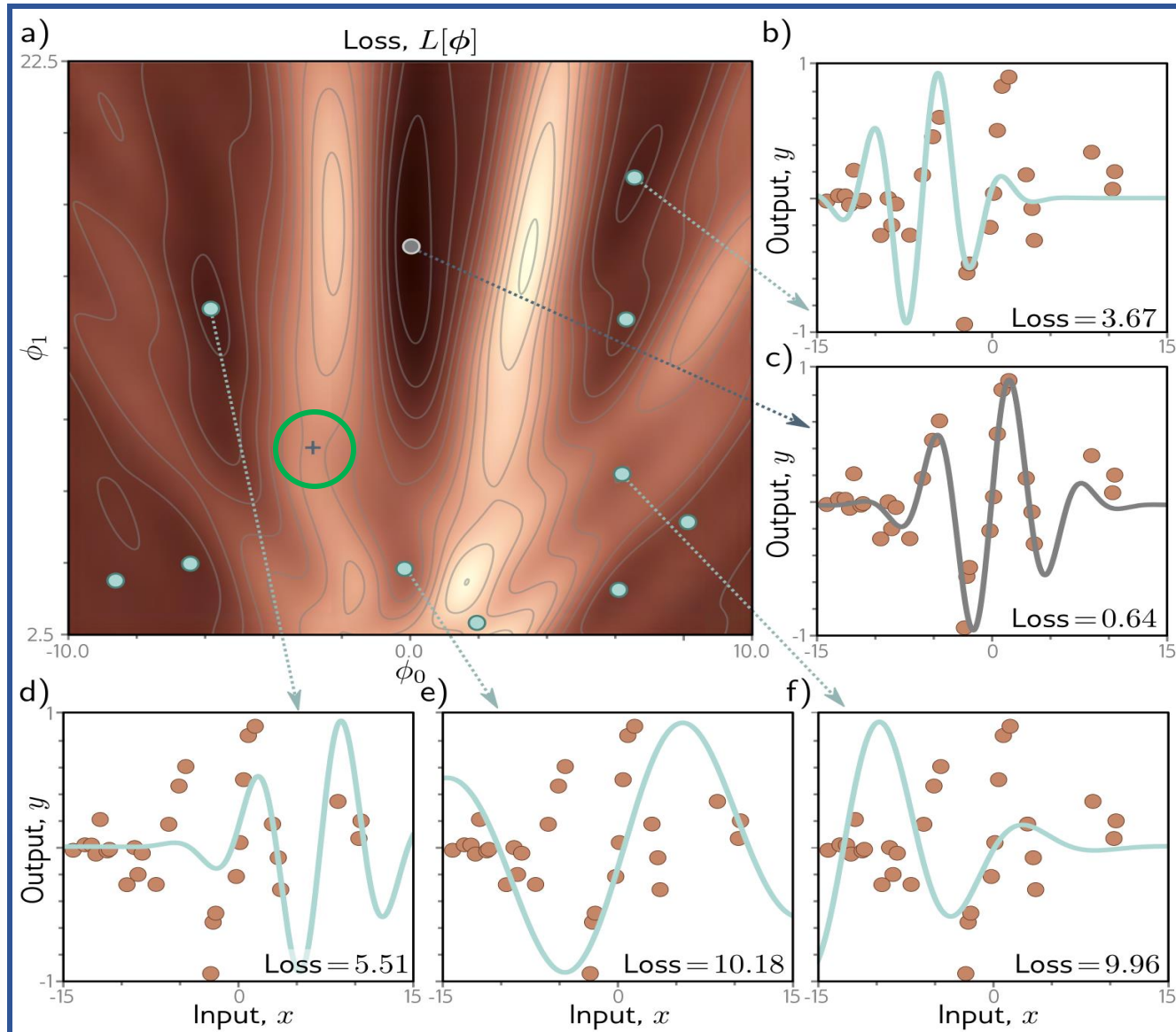


Optimization

❑ Saddle Points

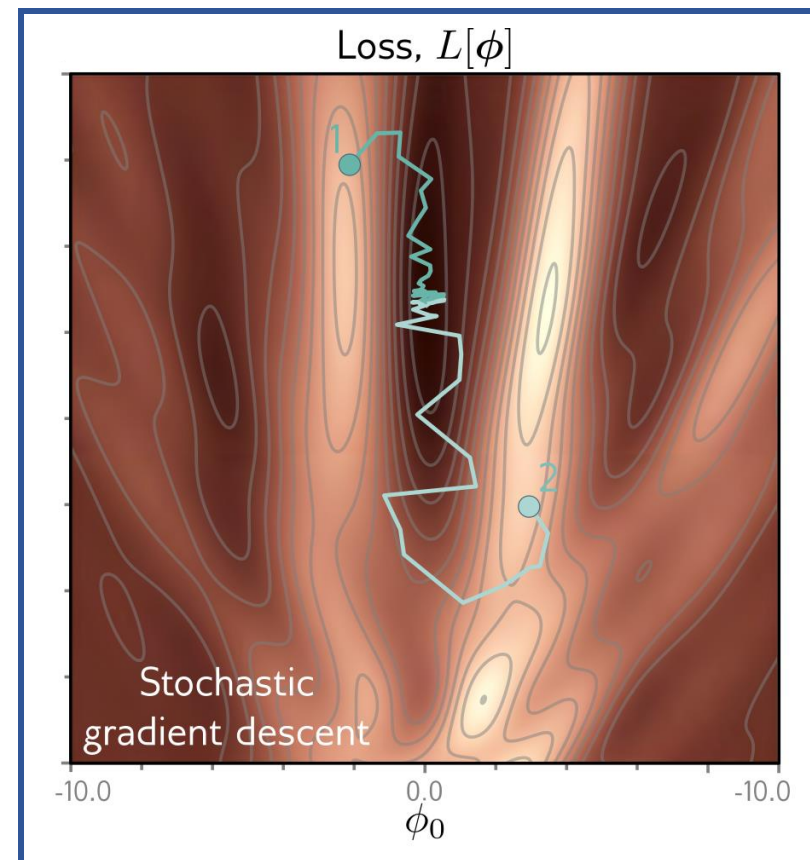
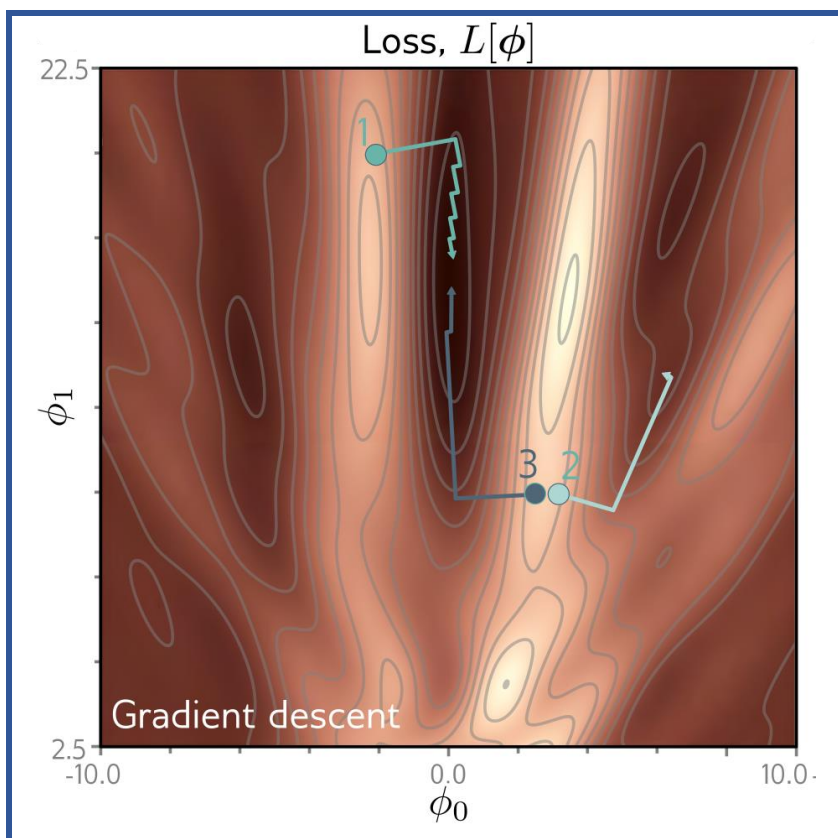


How to escape from saddle points and bad local minima?



❑ Stochastic Gradient Descent (SGD)

- Introduces some noise in the optimization process.



Optimization

❑ Stochastic Gradient Descent (SGD)

$$\hat{\phi} = \operatorname{argmin}_{\phi} [L(\phi)]$$

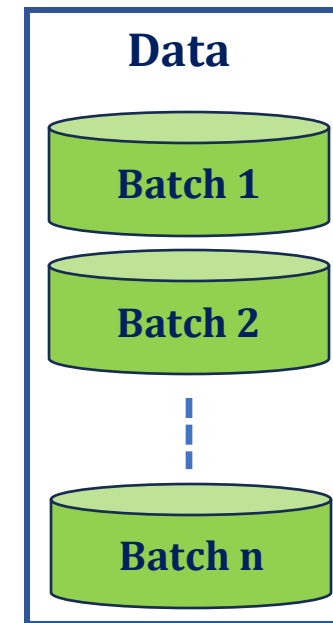
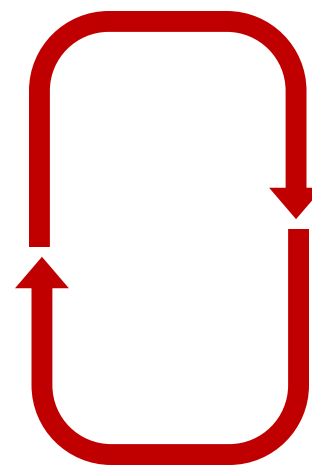
Step 1: Initialize the parameters ϕ_0

Step 2: Compute the derivatives of loss with respect to the parameters.

$$\sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\phi_t]}{\partial \phi}$$

Step 3: Update the parameters.

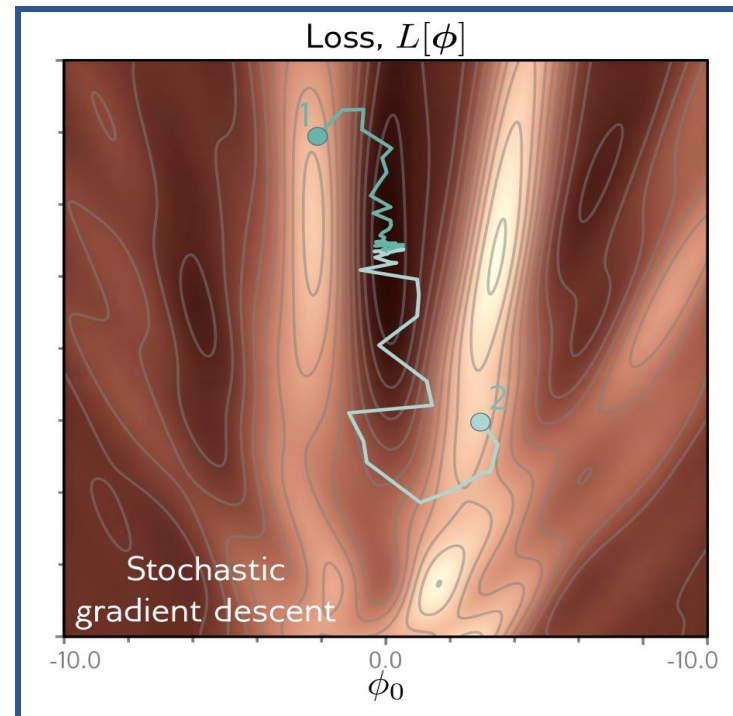
$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\phi_t]}{\partial \phi}$$



Batch or Mini-batch

□ Stochastic Gradient Descent (SGD)

- Computationally **less expensive** to compute the batch gradient.
- In principle, it can escape **local minima and saddle points**.



Any Major Drawback in Gradient Descent?

Optimization

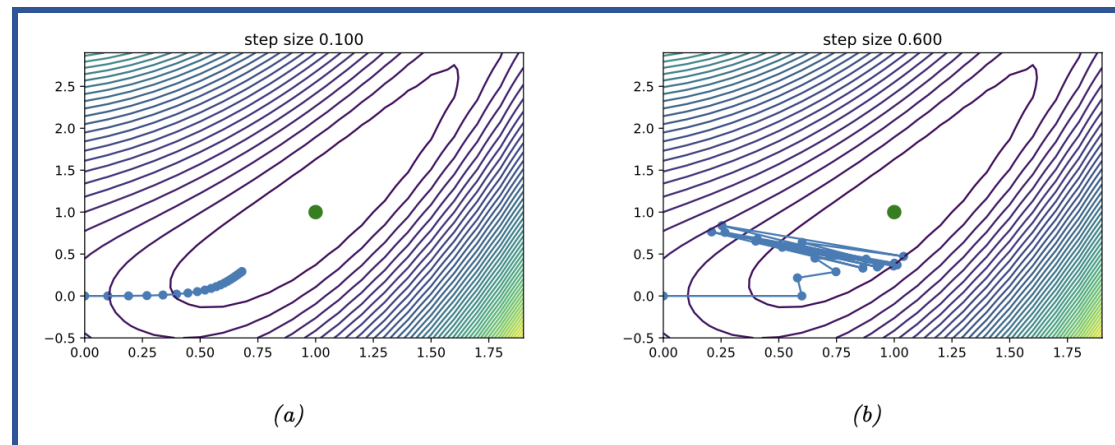
□ Momentum

- Gradient descent move very slowly.
- Simple heuristic: **Momentum or heavy ball**
 - Move fast in the previously good direction.
 - Slow down if the gradient suddenly changes.

Like a rolling ball!

$$\begin{aligned} \mathbf{m}_t &= \beta \mathbf{m}_{t-1} + \mathbf{g}_{t-1} \\ \boldsymbol{\theta}_t &= \boldsymbol{\theta}_{t-1} - \eta_t \mathbf{m}_t \end{aligned}$$

- If gradients are different:



$$\mathbf{m}_t = \beta \mathbf{m}_{t-1} + \mathbf{g}_{t-1} = \beta^2 \mathbf{m}_{t-2} + \beta \mathbf{g}_{t-2} + \mathbf{g}_{t-1} = \dots = \sum_{\tau=0}^{t-1} \beta^\tau \mathbf{g}_{t-\tau-1}$$

- If gradients are same:

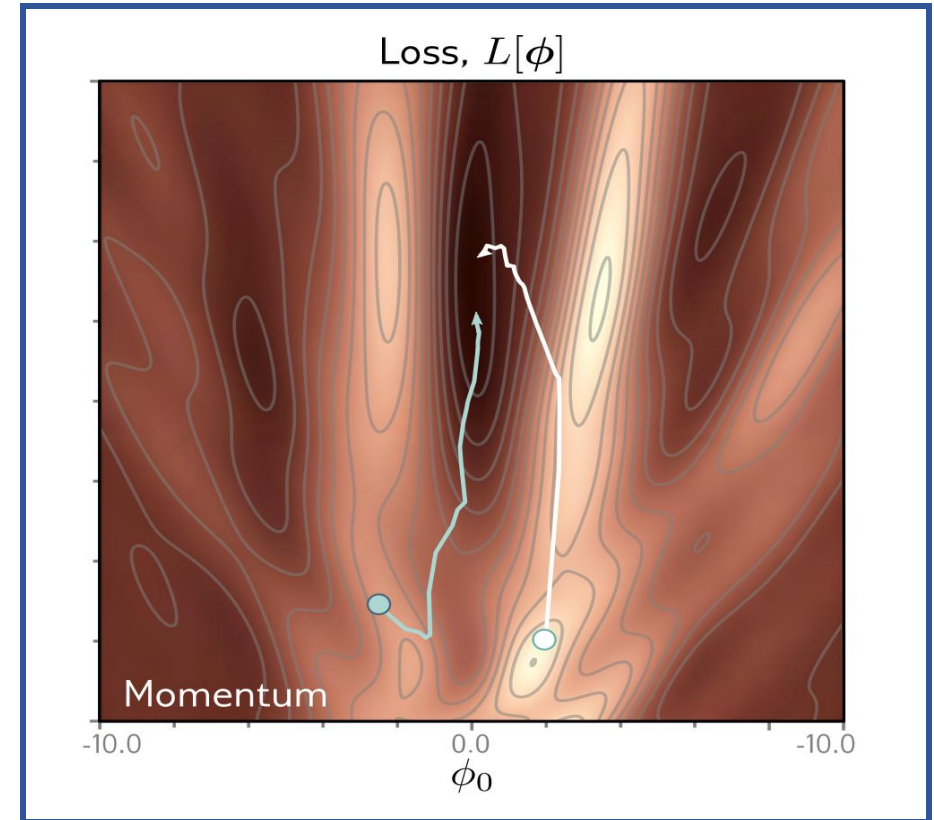
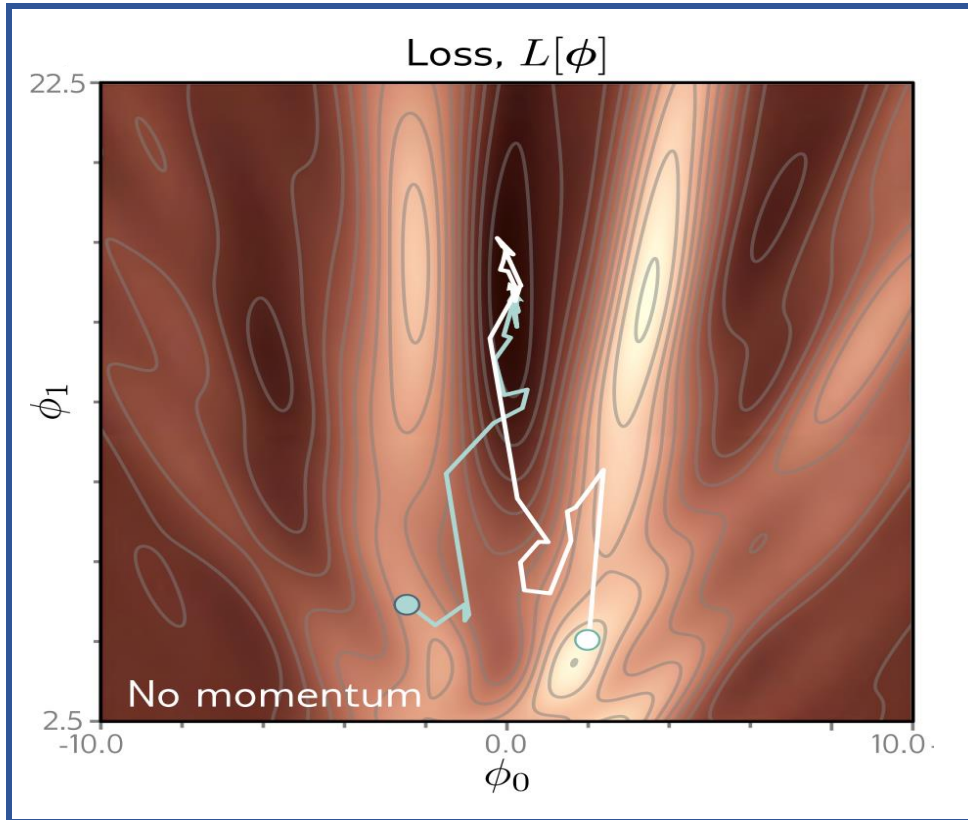
$$\mathbf{m}_t = \mathbf{g} \sum_{\tau=0}^{t-1} \beta^\tau$$

$$1 + \beta + \beta^2 + \dots = \sum_{i=0}^{\infty} \beta^i = \frac{1}{1 - \beta}$$

If, $\beta = 0.9$ We are **scaling** the gradient by a factor of **10**.

Move Fast!

❑ Momentum



□ Learning Rate

■ Parameter Update:

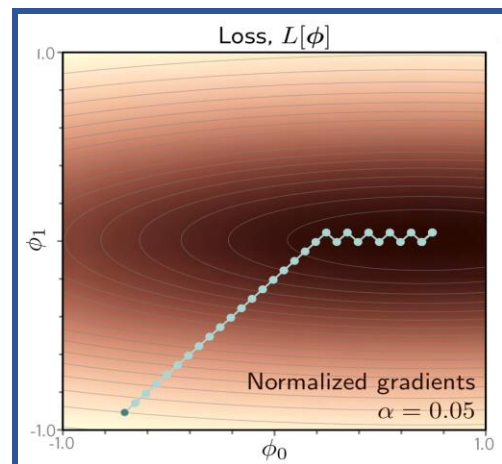
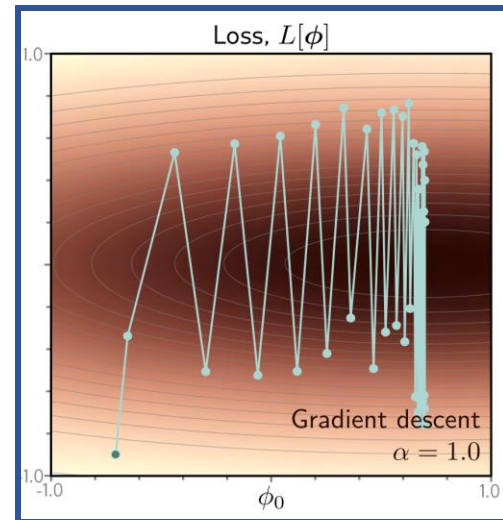
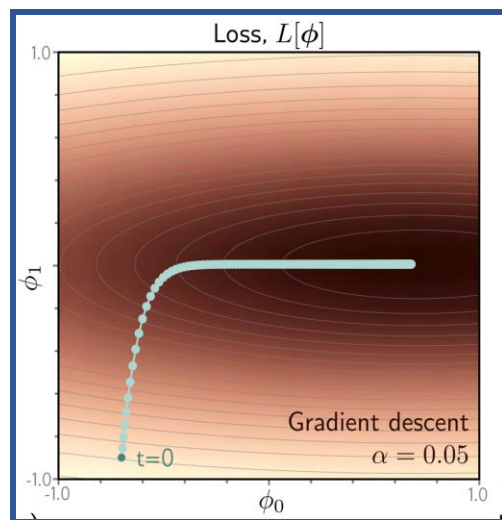
$$\phi \leftarrow \phi - \alpha \cdot \frac{\partial L}{\partial \phi}$$

Learning rate

■ Normalize the gradient:

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1}} + \epsilon}$$

$$\begin{aligned} \mathbf{m}_{t+1} &\leftarrow \frac{\partial L[\phi_t]}{\partial \phi} \\ \mathbf{v}_{t+1} &\leftarrow \left(\frac{\partial L[\phi_t]}{\partial \phi} \right)^2 \end{aligned}$$

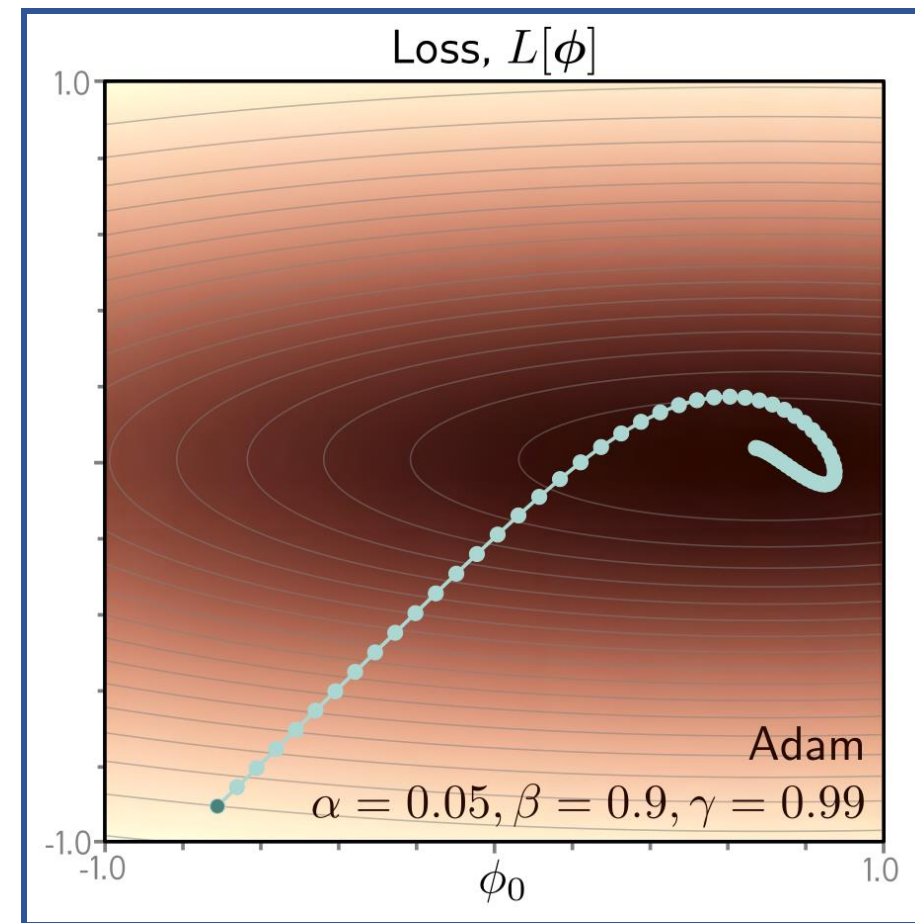


□ Adaptive Momentum Estimation (Adam) Optimizer

- Include Momentum:

$$\begin{aligned}\mathbf{m}_{t+1} &\leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta) \frac{\partial L[\phi_t]}{\partial \phi} \\ \mathbf{v}_{t+1} &\leftarrow \gamma \cdot \mathbf{v}_t + (1 - \gamma) \left(\frac{\partial L[\phi_t]}{\partial \phi} \right)^2\end{aligned}$$

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1}} + \epsilon}$$



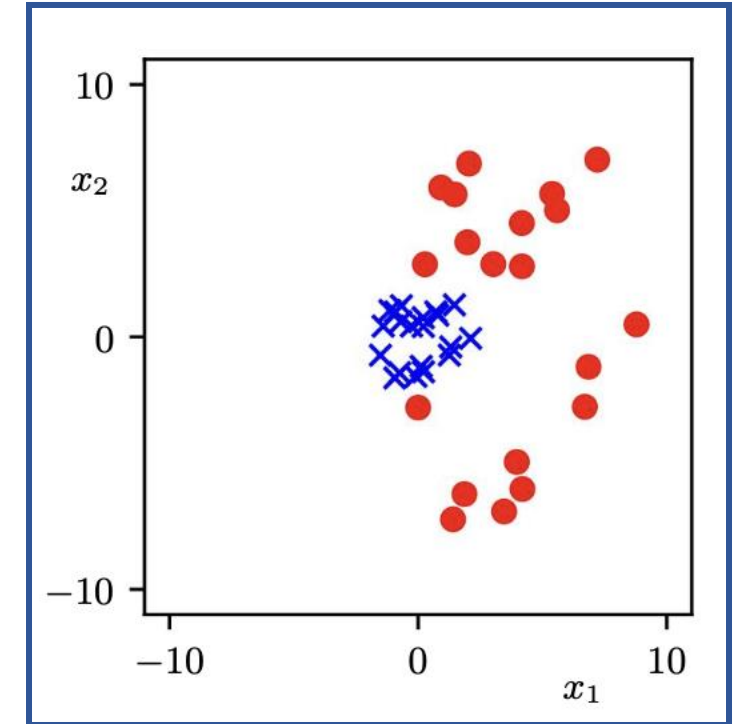
□ Input Data Normalization

- Input features may have a wide range of values.
- Done prior to training.

$$\mu_i = \frac{1}{N} \sum_{n=1}^N x_{ni}$$
$$\sigma_i^2 = \frac{1}{N} \sum_{n=1}^N (x_{ni} - \mu_i)^2,$$

$$\tilde{x}_{ni} = \frac{x_{ni} - \mu_i}{\sigma_i}$$

- Zero mean and unit variance.



Optimization

□ Batch Normalization

- Normalizing the hidden layers.
- Vanishing gradients and exploding gradients.

$$\frac{\partial E}{\partial w_i} = \sum_m \cdots \sum_l \sum_j \frac{\partial z_m^{(1)}}{\partial w_i} \cdots \frac{\partial z_j^{(K)}}{\partial z_l^{(K-1)}} \frac{\partial E}{\partial z_j^{(K)}}$$

$$z_i = h(a_i)$$

- Any problem with the normalization?

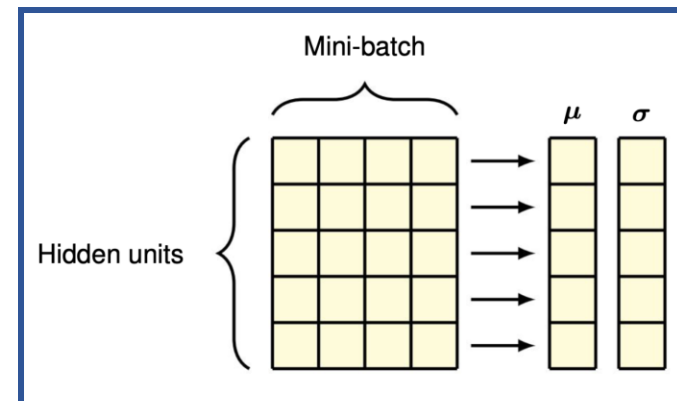
- We reduce the degree of freedom/ representation capacity of that hidden layer.

$$\tilde{a}_{ni} = \gamma_i \hat{a}_{ni} + \beta_i$$

Stdv

Mean

- β_i and γ_i are adaptive learnable parameters.
- Trained using gradient descent, same as other parameters.



$$\mu_i = \frac{1}{K} \sum_{n=1}^K a_{ni}$$

$$\sigma_i^2 = \frac{1}{K} \sum_{n=1}^K (a_{ni} - \mu_i)^2$$

$$\hat{a}_{ni} = \frac{a_{ni} - \mu_i}{\sqrt{\sigma_i^2 + \delta}}$$

Batch Normalization Layer

END