



Introduction to Machine Learning [Fall 2022]

Diving Deeper into Gradients

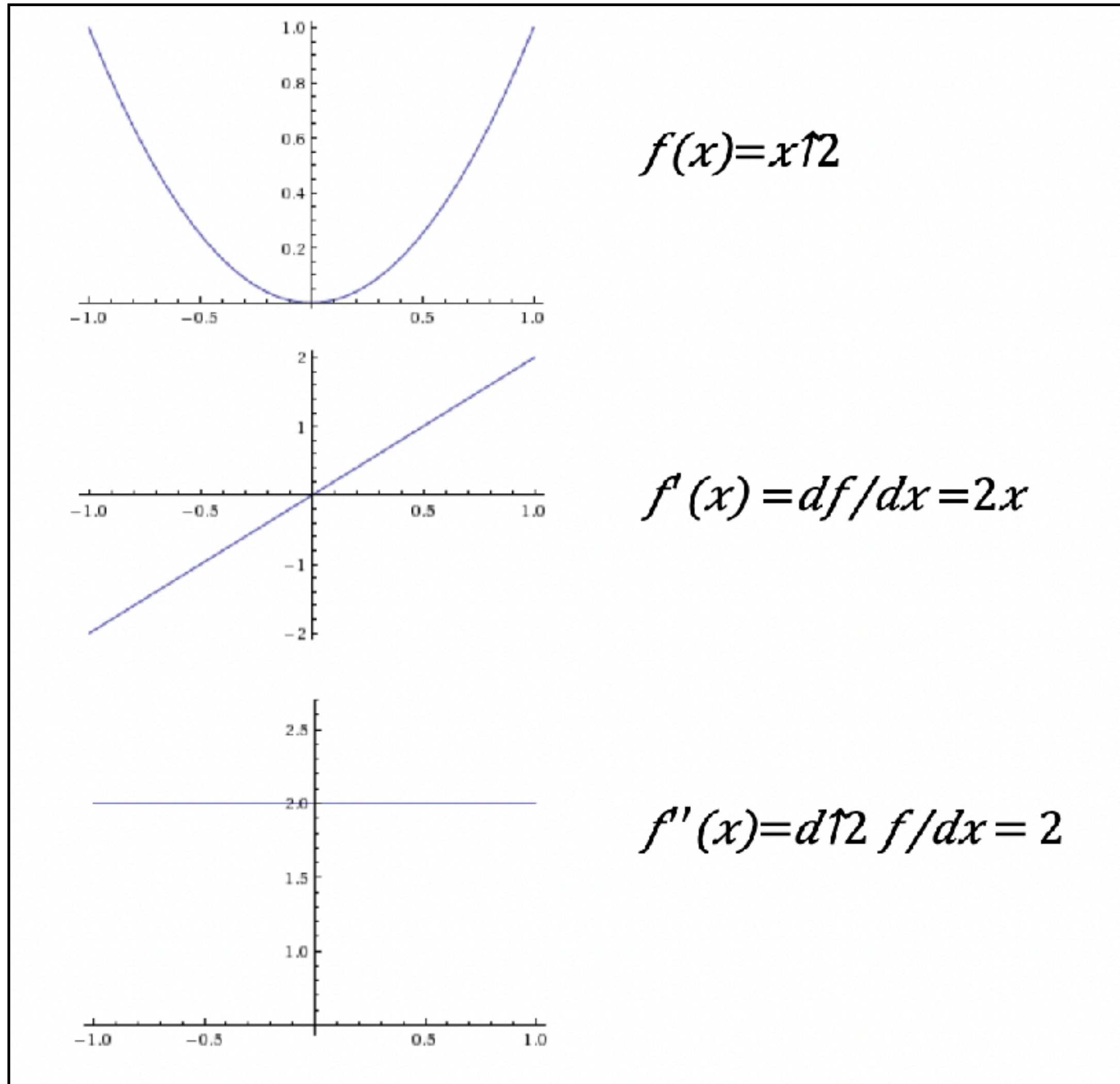
October 18, 2022

Lerrel Pinto

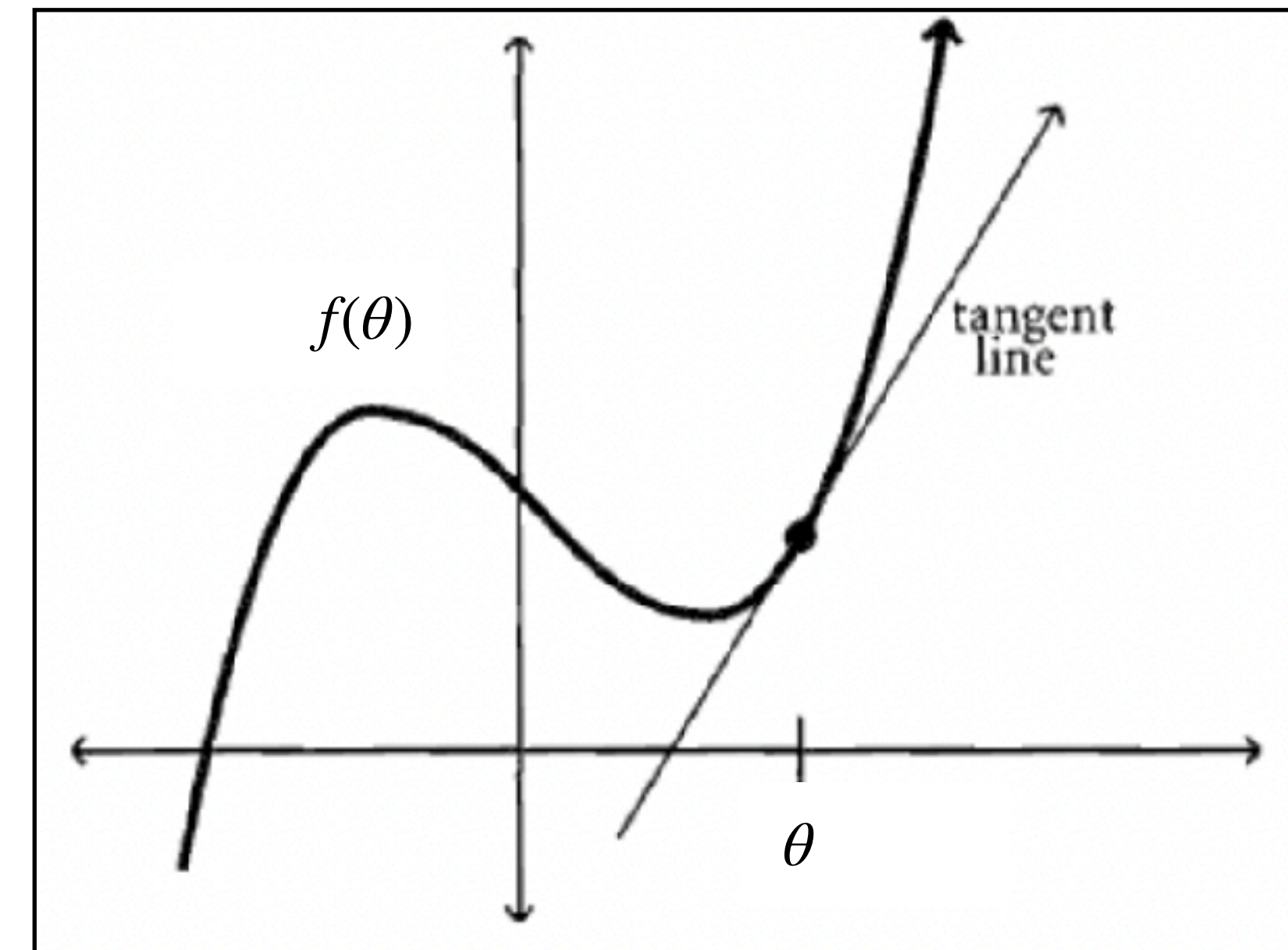
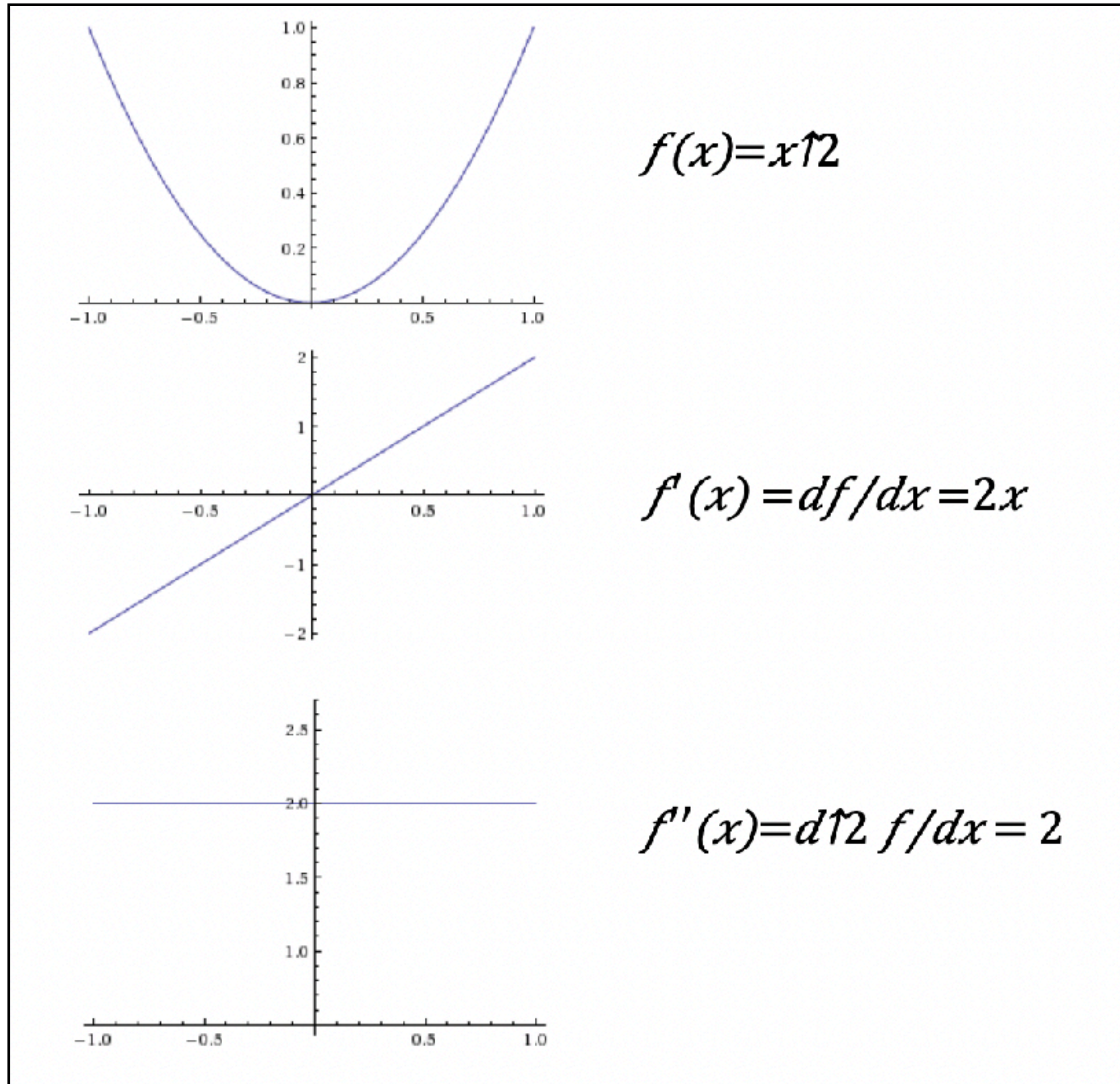
Topics for today

- Fundamentals for gradients
- Gradients with high-dimensional inputs
- Examples of gradient descent

Recap: Derivatives



Recap: Derivatives



Gradient Descent Algorithm

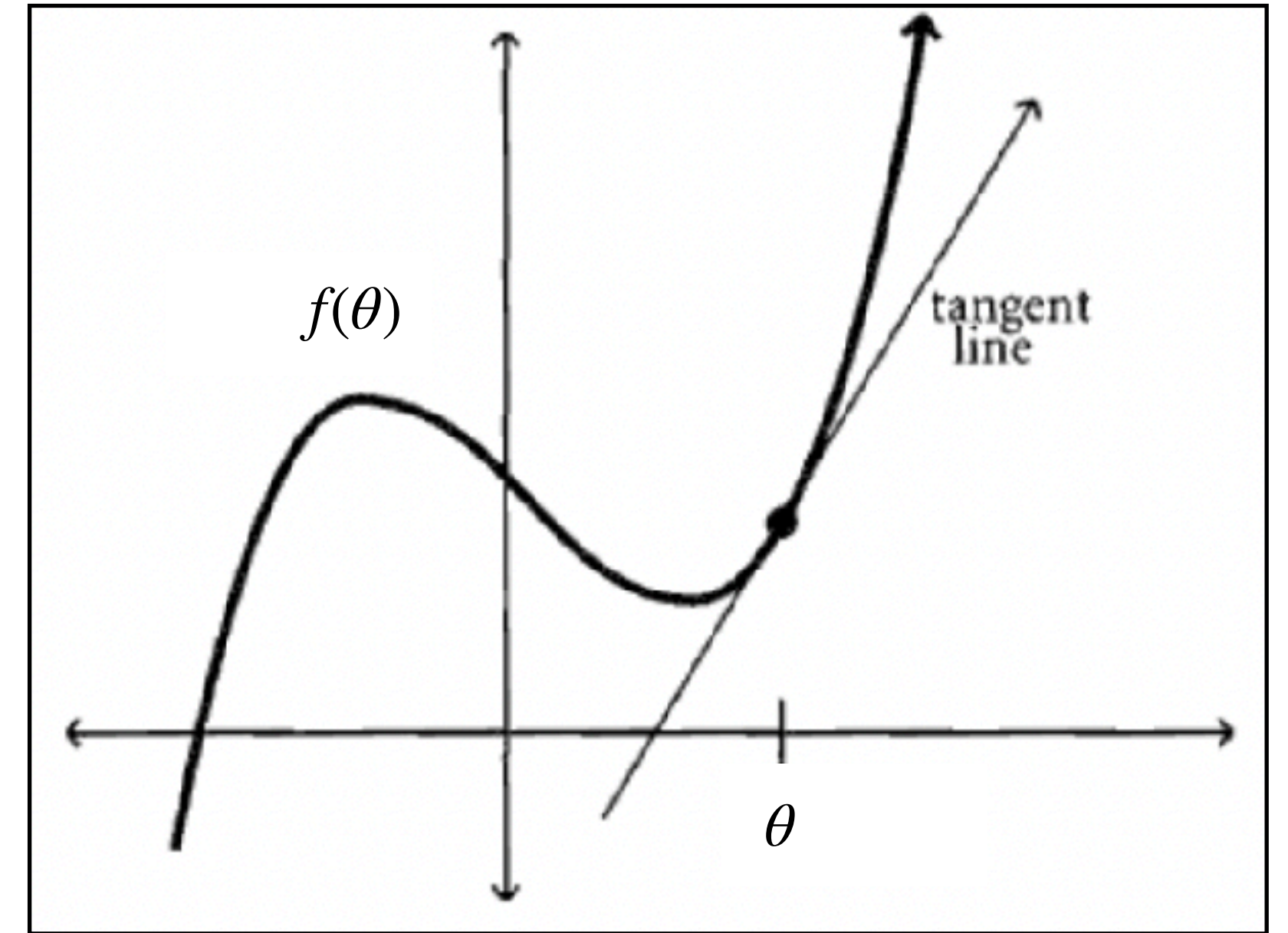
- Given: cost / loss/ objective function $f(\theta)$.
- Goal: find θ^* such that $f(\theta^*) = \min_{\theta} f(\theta)$.

Gradient Descent Algorithm

- Given: cost / loss/ objective function $f(\theta)$.
- Goal: find θ^* such that $f(\theta^*) = \min_{\theta} f(\theta)$.
- Gradient descent solution:
 - Start from initial guess θ^0 and learning rate α
 - Update $\theta^{i+1} \leftarrow \theta^i - \alpha \frac{df(\theta)}{d\theta}$
 - Repeat until change in θ is small, or maximum number of steps reached.

Gradient Descent Algorithm

- Given: cost / loss/ objective function $f(\theta)$.
- Goal: find θ^* such that $f(\theta^*) = \min_{\theta} f(\theta)$.
- Gradient descent solution:

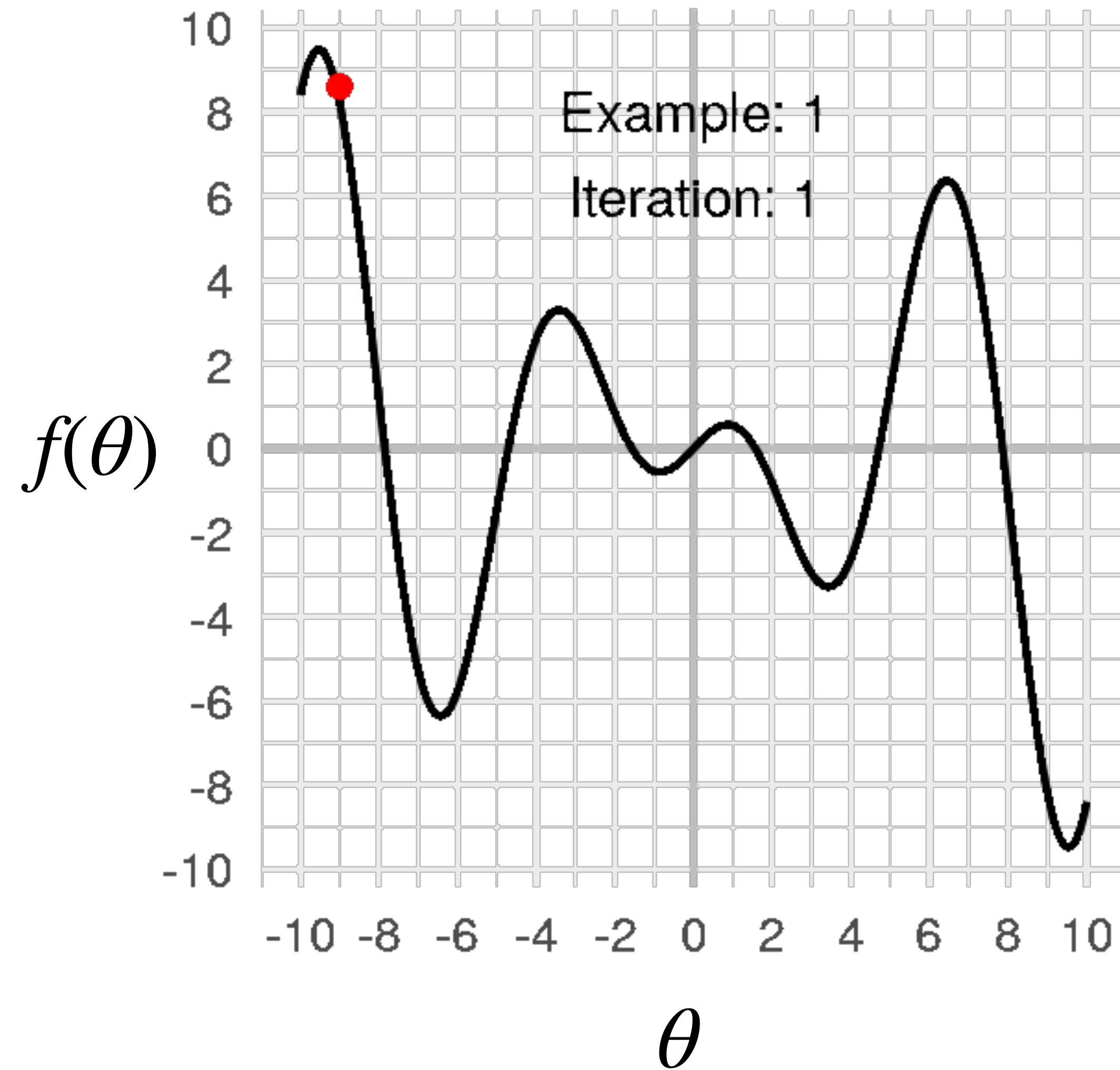


- Start from initial guess θ^0 and learning rate α

- Update $\theta^{i+1} \leftarrow \theta^i - \alpha \frac{df(\theta)}{d\theta}$

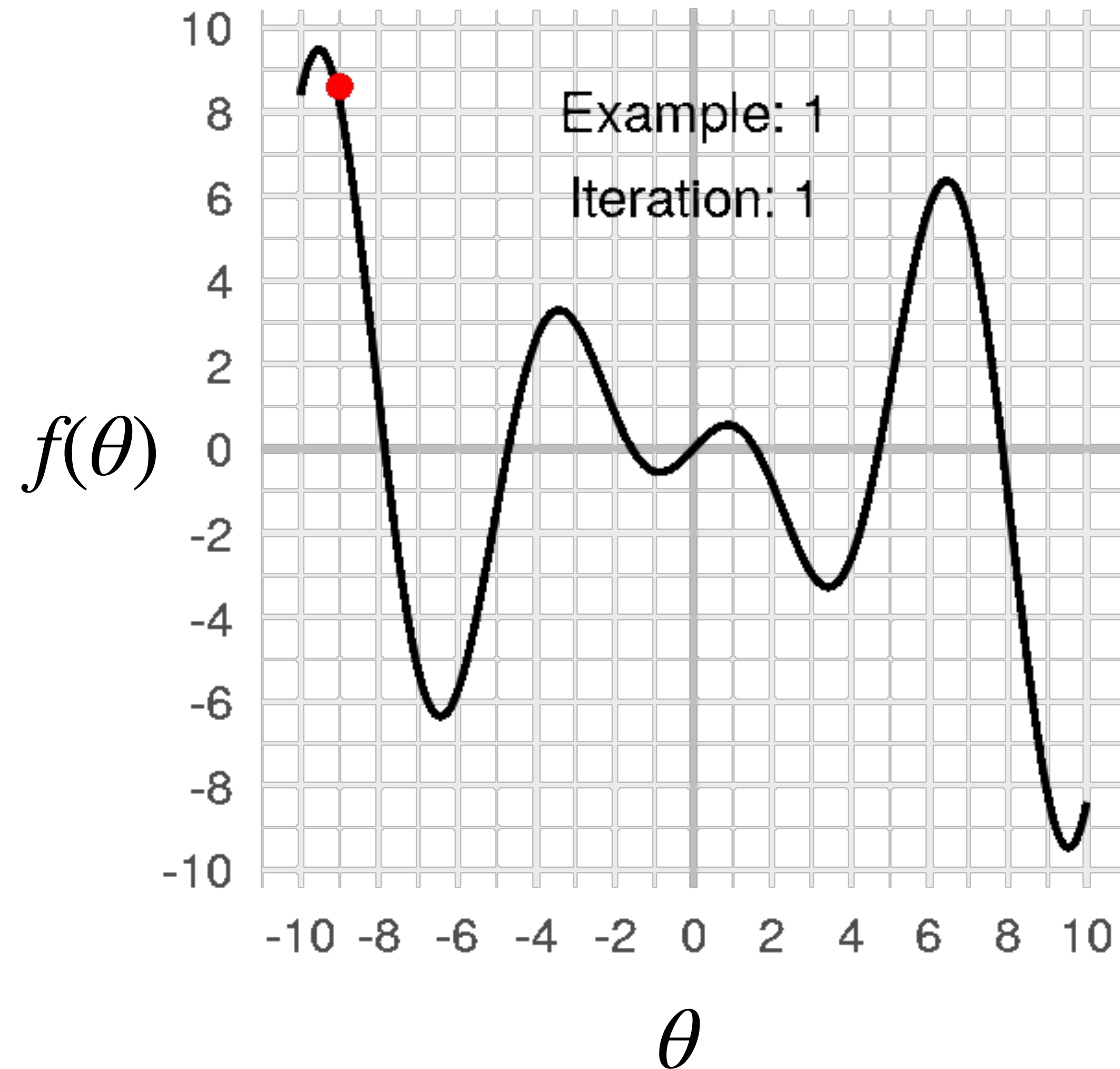
- Repeat until change in θ is small, or maximum number of steps reached.

Gradient Descent Algorithm



Credits: Charles Bordet

Gradient Descent Algorithm

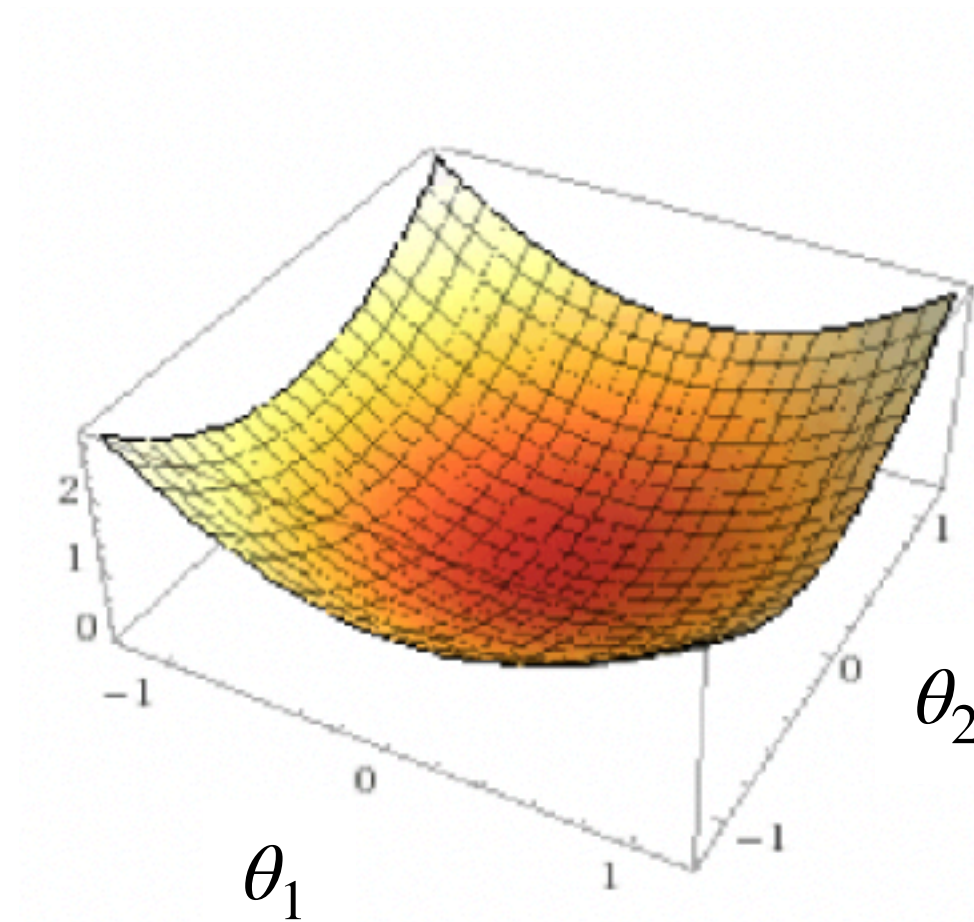


Credits: Charles Bordet

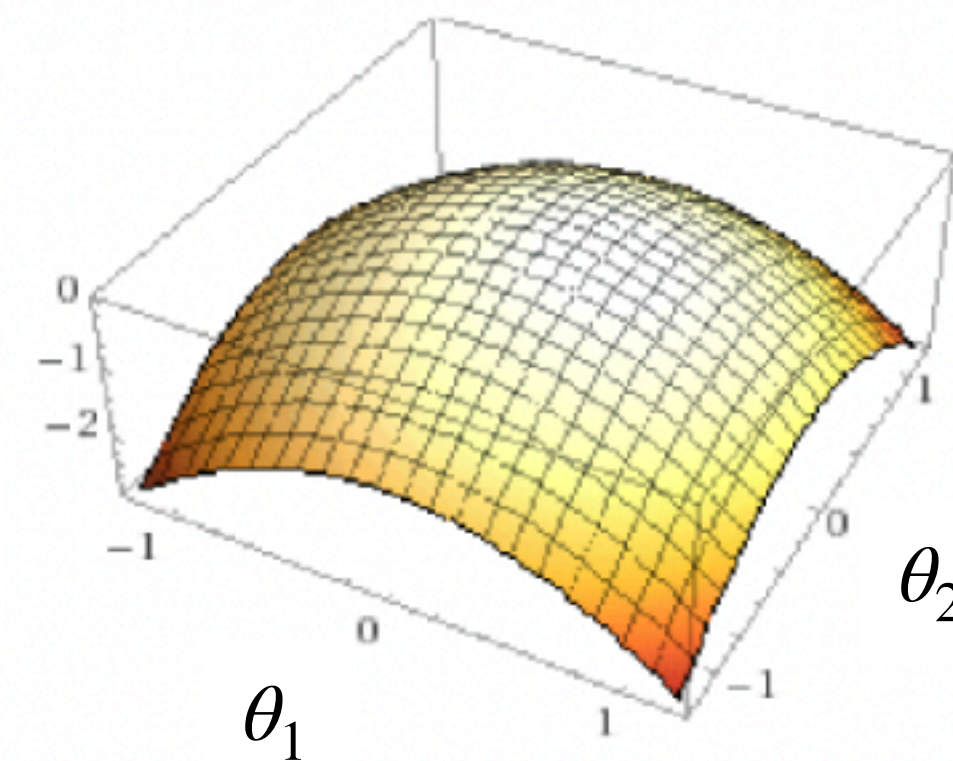
Interactive demo:

<https://uclaacm.github.io/gradient-descent-visualiser/>

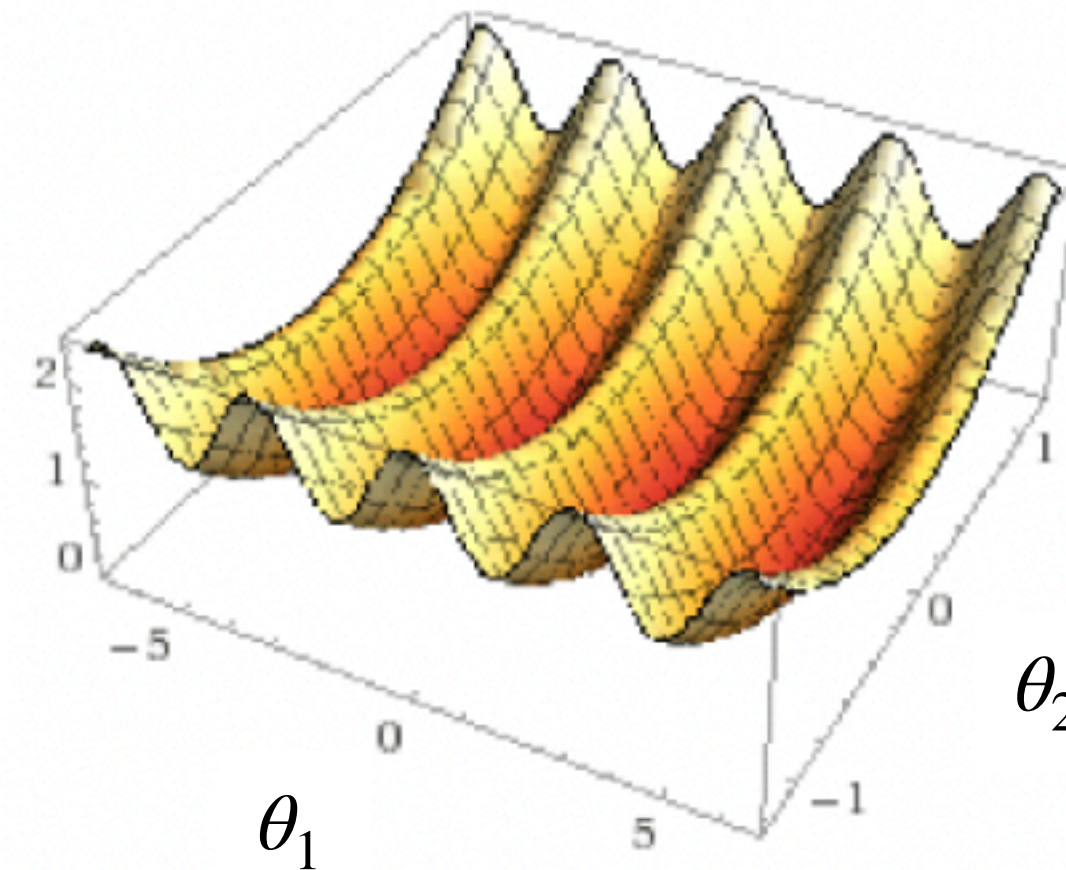
What happens when number of dimensions is high?



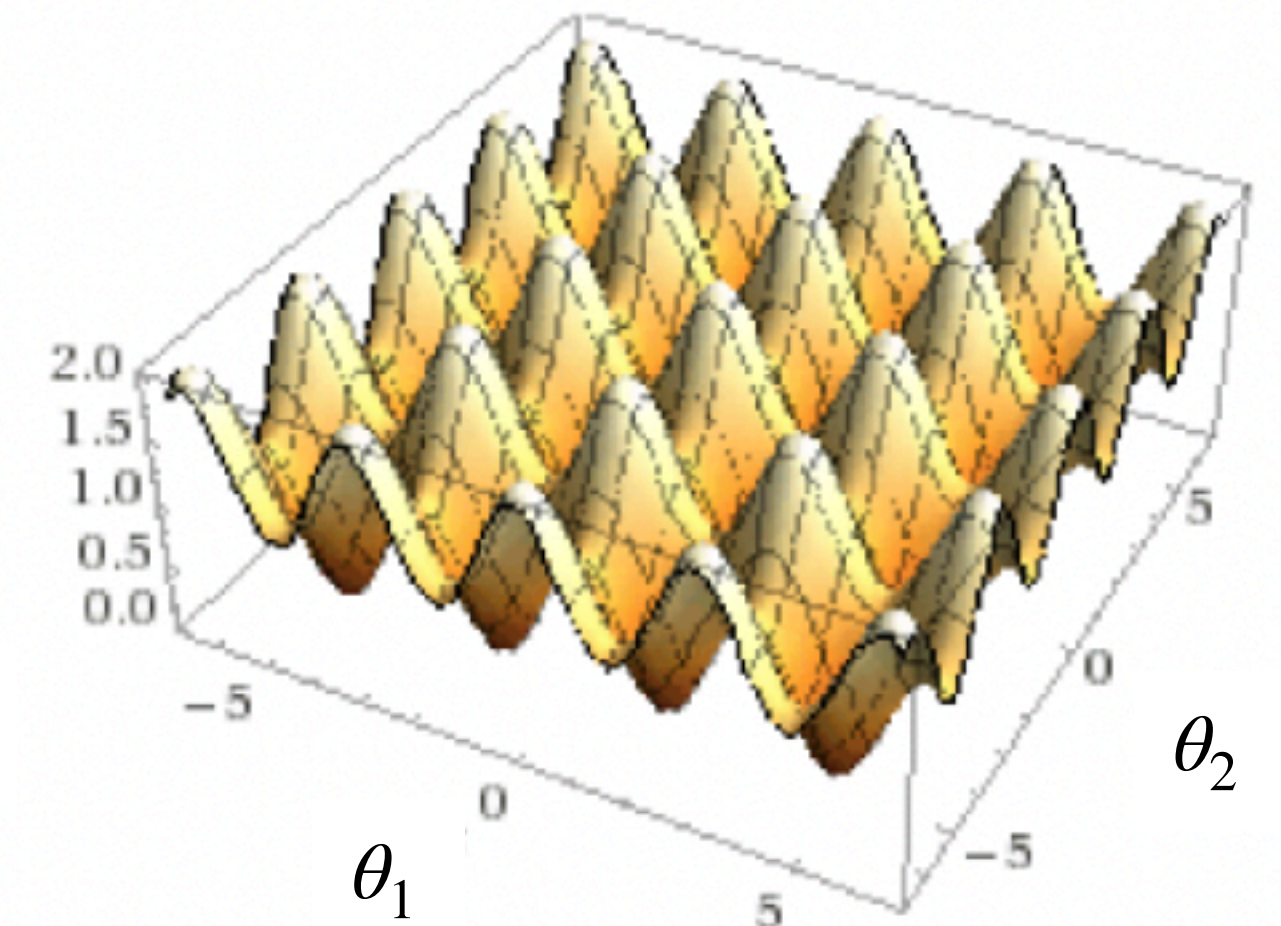
$$f(\theta_1, \theta_2) = \theta_1^2 + \theta_2^2$$



$$f(\theta_1, \theta_2) = -\theta_1^2 - \theta_2^2$$



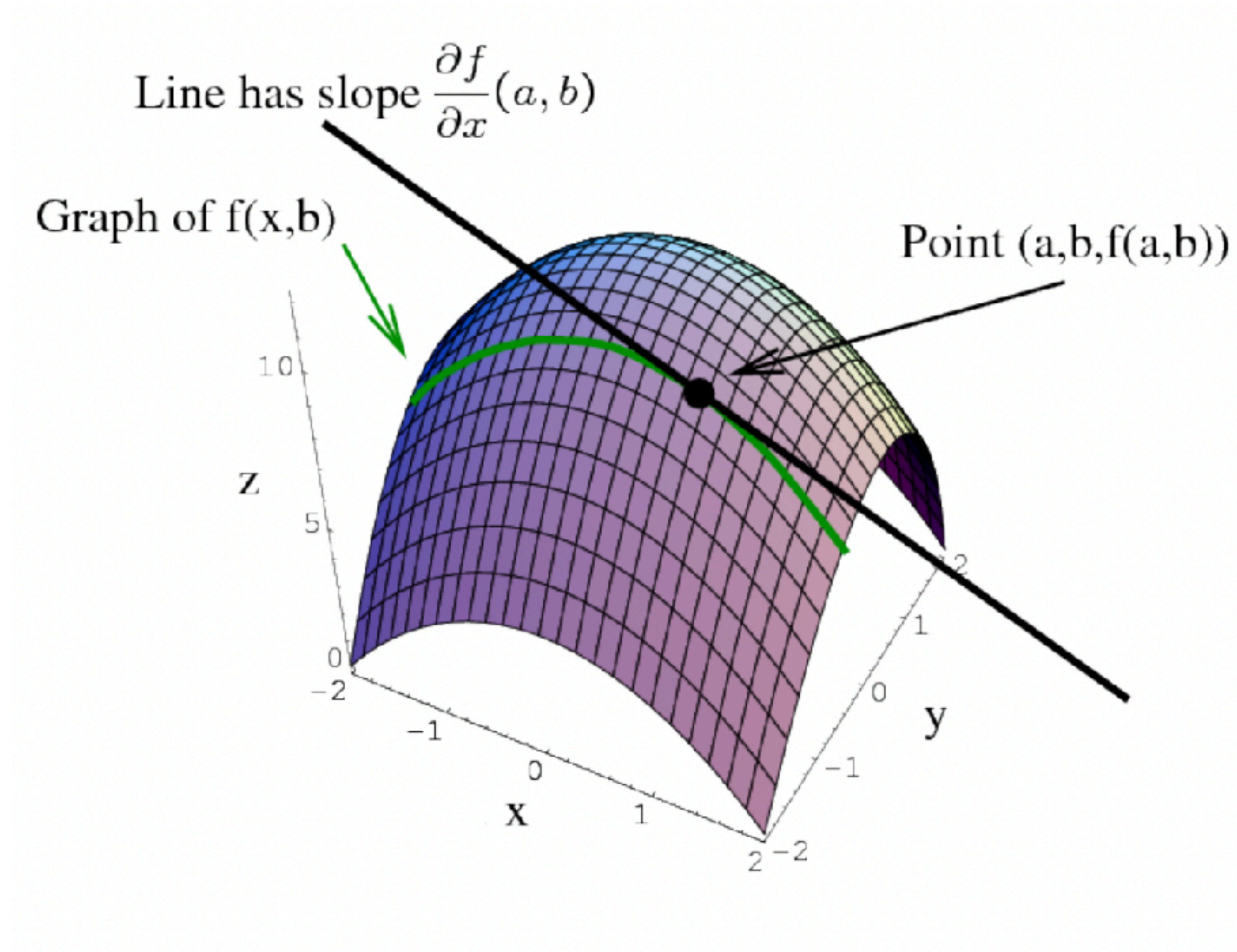
$$f(\theta_1, \theta_2) = \cos(\theta_1)^2 + \theta_2^2$$



$$f(\theta_1, \theta_2) = \cos(\theta_1)^2 + \cos(\theta_2)^2$$

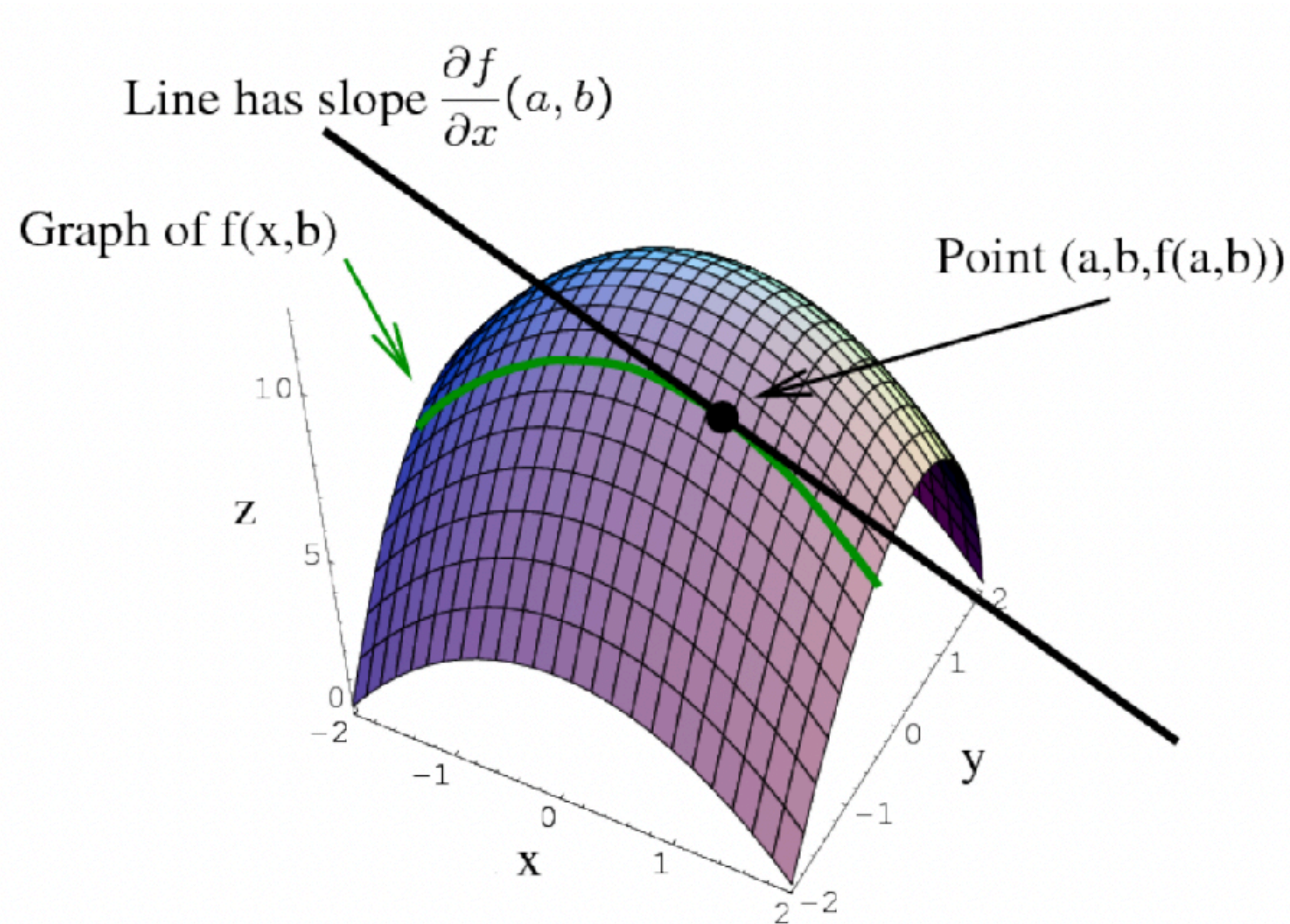
Credits: Michail Michailidis & Patrick Maiden

What happens when number of dimensions is high?



Credits: Michail Michailidis & Patrick Maiden

What happens when number of dimensions is high?



$$f(x, y) = 9 - x^2 - y^2$$

$$f(x, y) = 9 - x^2 - c^2$$

$$f(x, y) = 9 - c^2 - y^2$$

$$\frac{\partial f(x, y)}{\partial x} = -2x$$

$$\frac{\partial f(x, y)}{\partial y} = -2y$$

Gradient vector:

$$\nabla f(x, y) = \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right)$$

$$\nabla f(x, y) = (-2x, -2y)$$

Credits: Michail Michailidis & Patrick Maiden

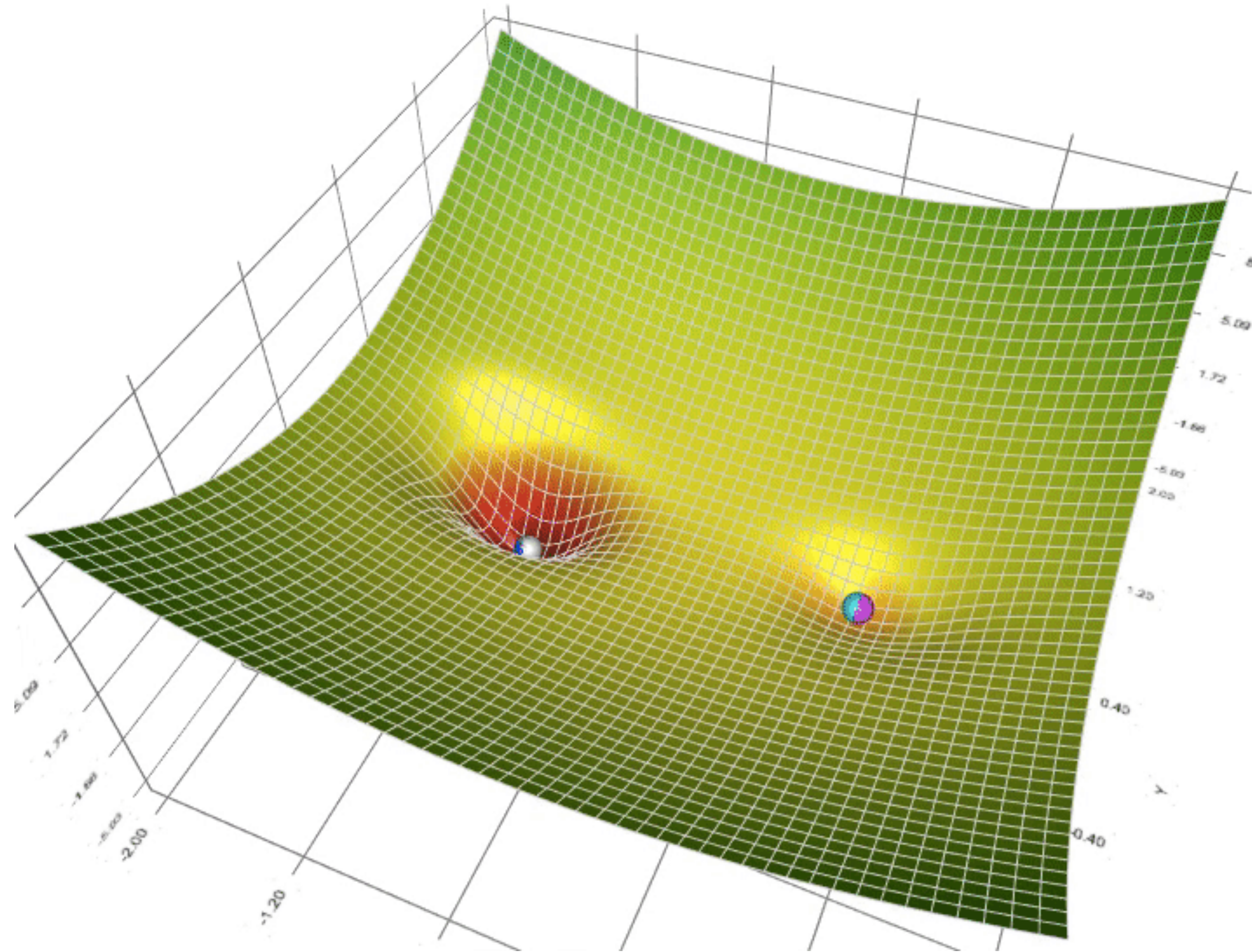
Gradient Descent Algorithm with multiple params

- Given: cost / loss/ objective function $f(\vec{\theta})$. Where $\vec{\theta} \in \mathbb{R}^d$.
- Goal: find $\vec{\theta}^*$ such that $f(\vec{\theta}^*) = \min_{\vec{\theta}} f(\vec{\theta})$.

Gradient Descent Algorithm with multiple params

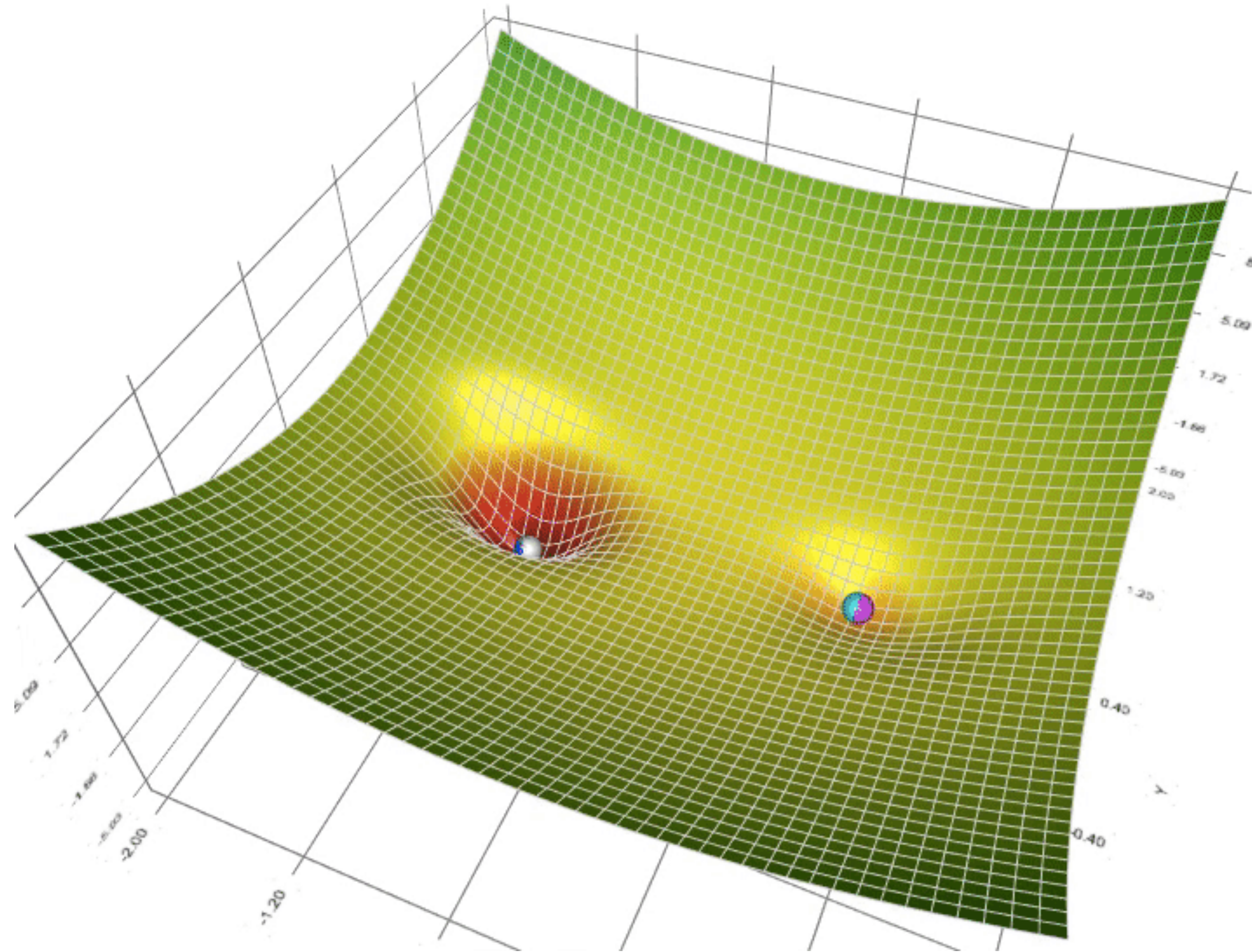
- Given: cost / loss/ objective function $f(\vec{\theta})$. Where $\vec{\theta} \in \mathbb{R}^d$.
- Goal: find $\vec{\theta}^*$ such that $f(\vec{\theta}^*) = \min_{\vec{\theta}} f(\vec{\theta})$.
- Gradient descent solution:
 - Start from initial guess $\vec{\theta}^0$ and learning rate α
 - Update $\vec{\theta}^{i+1} \leftarrow \vec{\theta}^i - \alpha \nabla f(\vec{\theta})$
 - Repeat until change in θ is small, or maximum number of steps reached.

Gradient Descent Algorithm Visualization



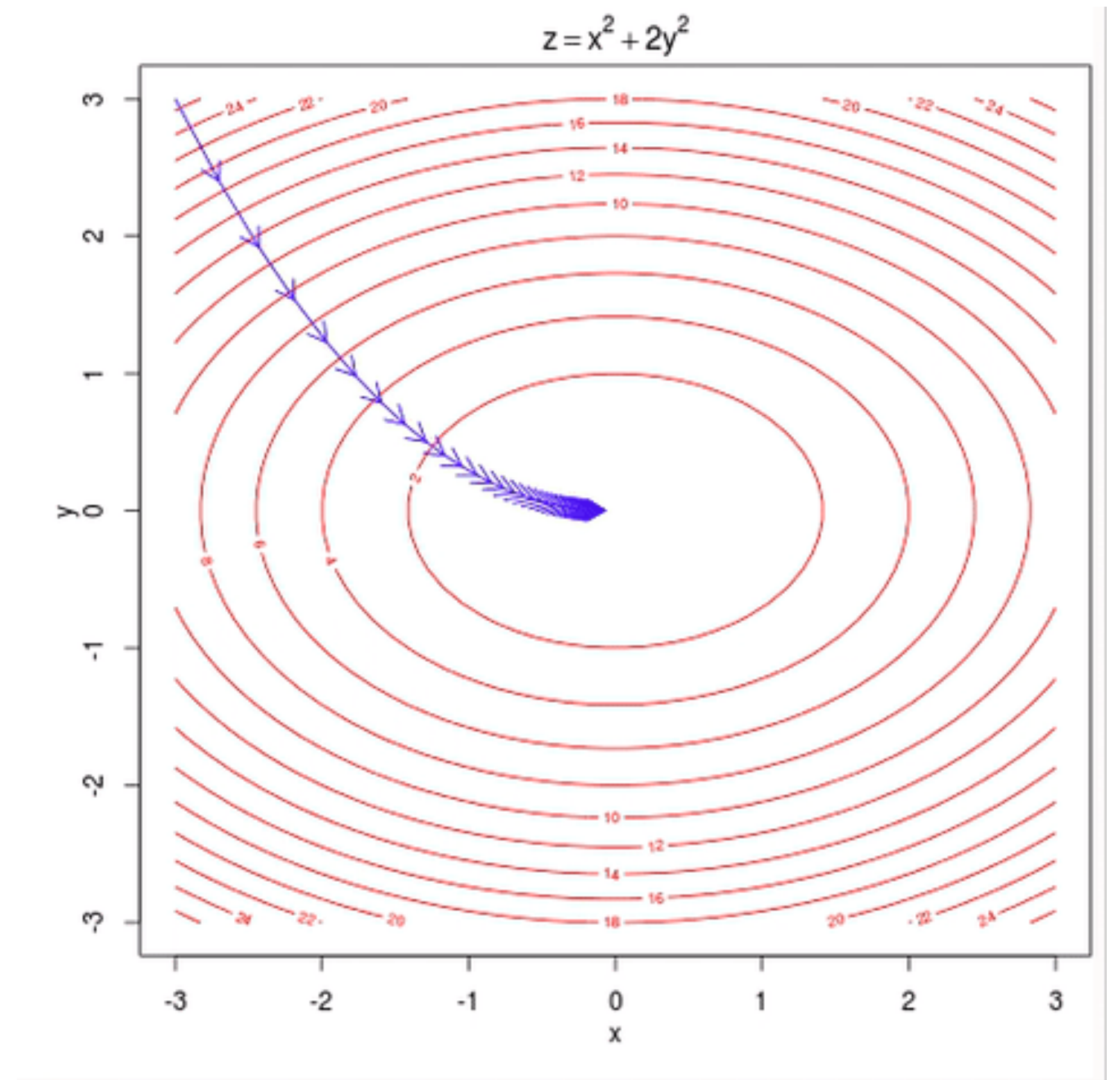
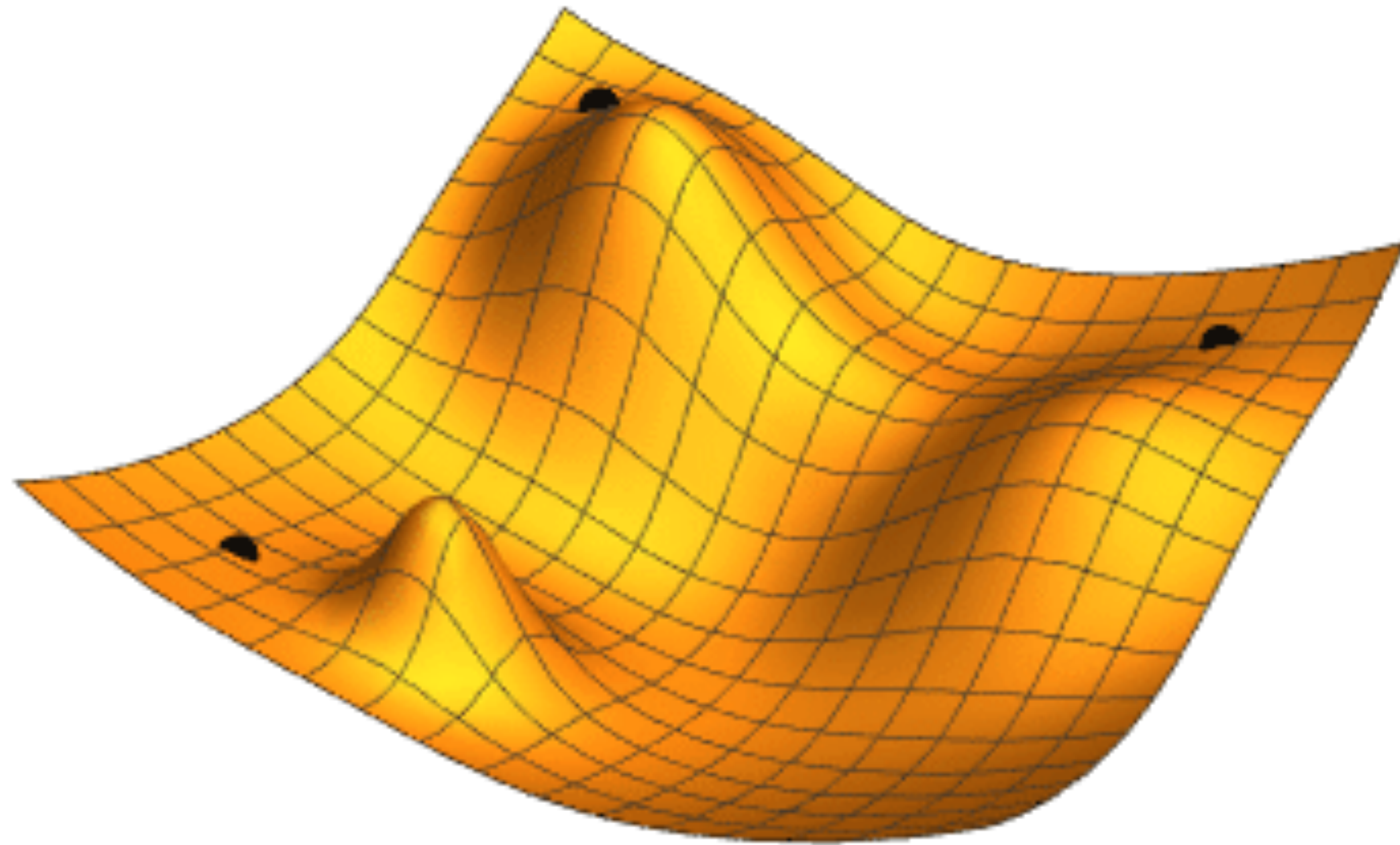
Credits: Lili Jiang

Gradient Descent Algorithm Visualization



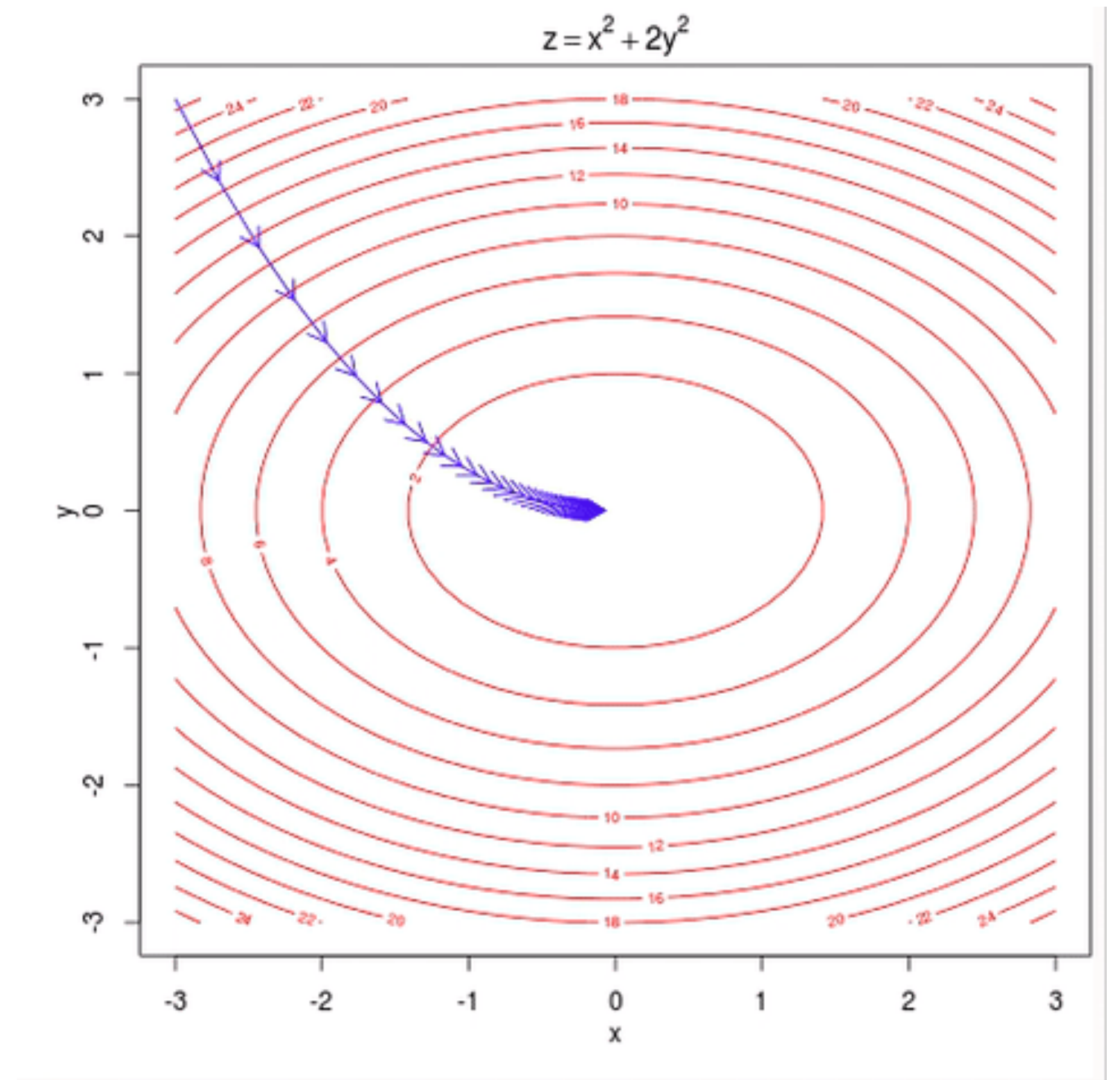
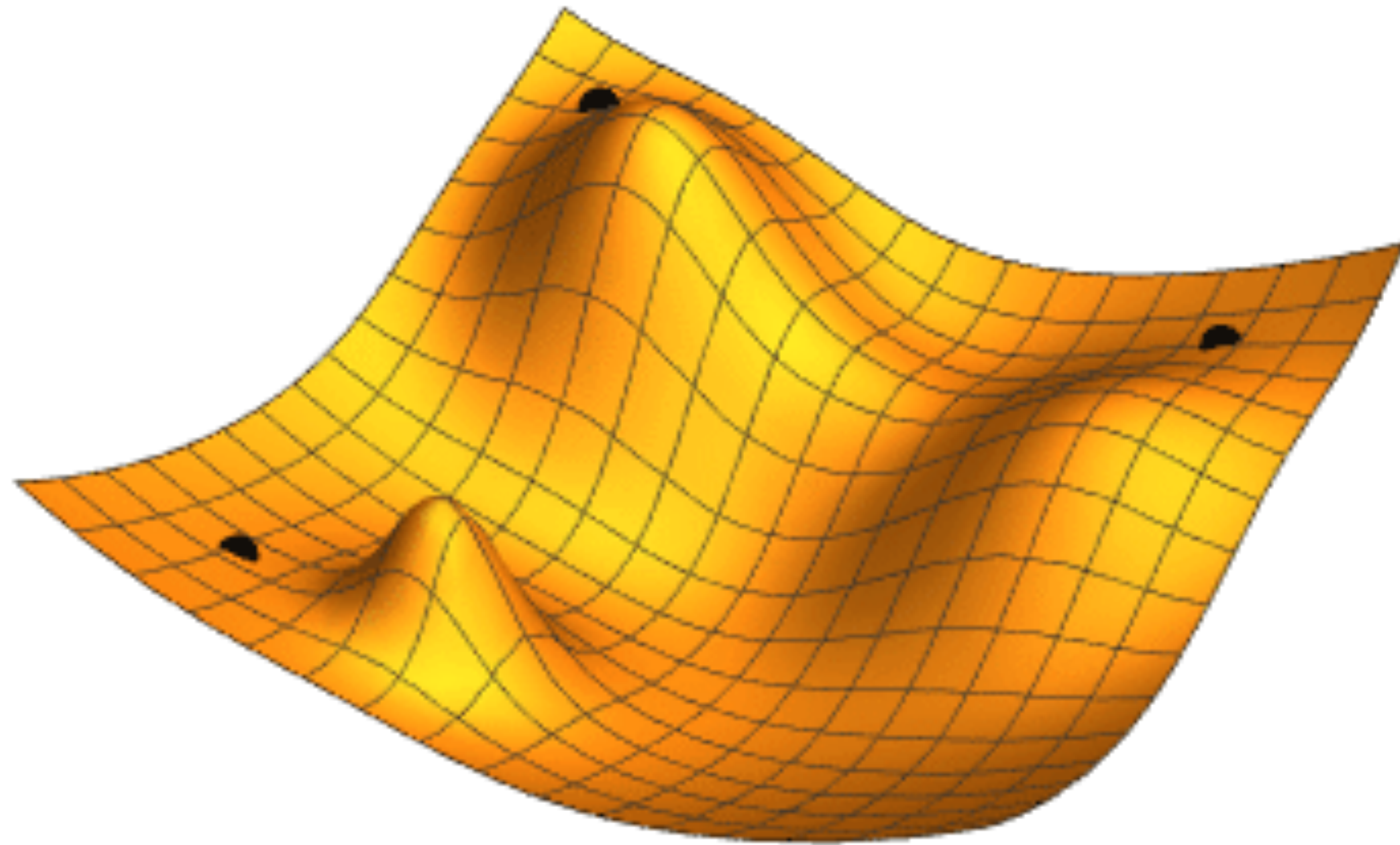
Credits: Lili Jiang

Gradient Descent Algorithm Visualization



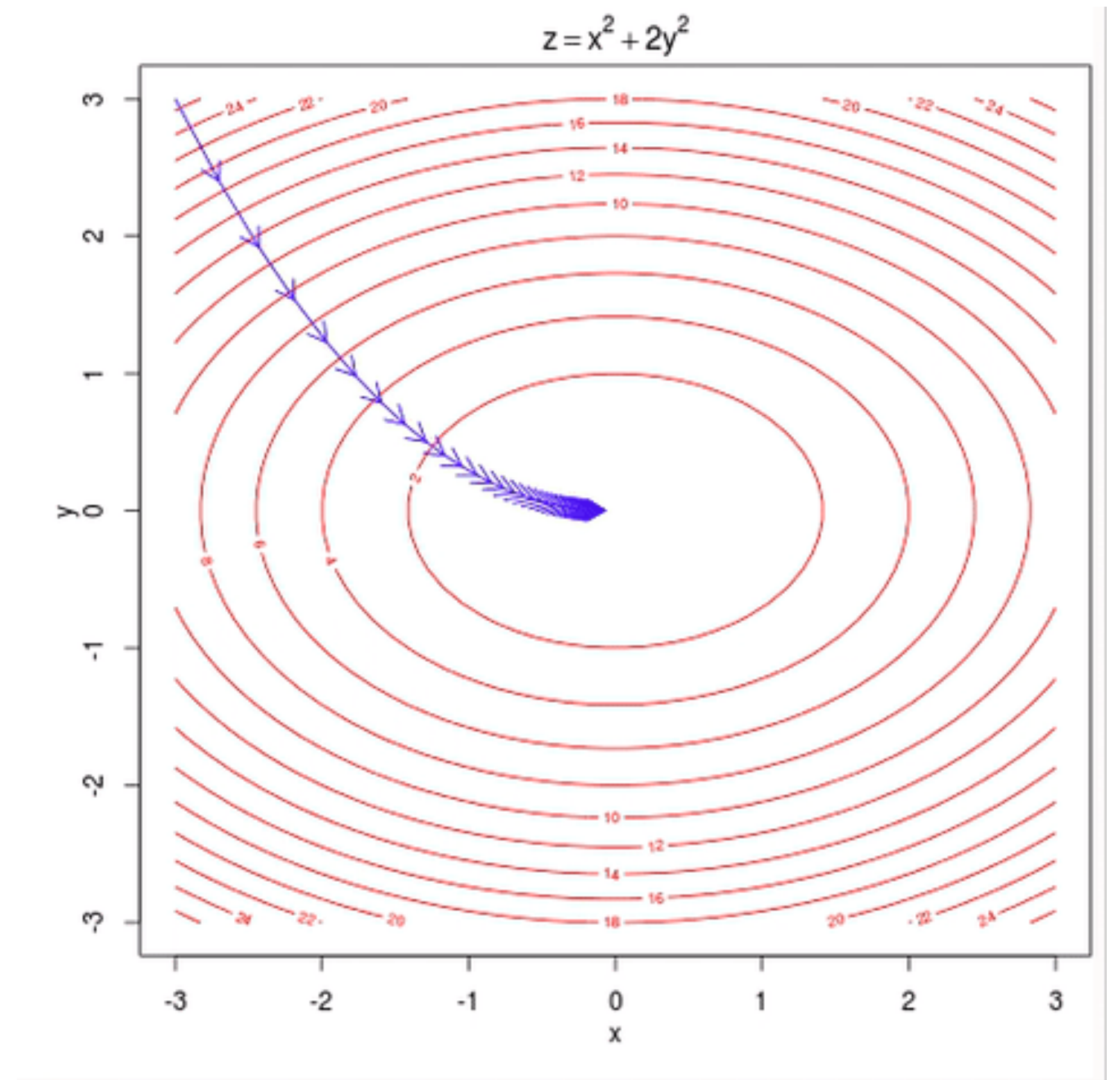
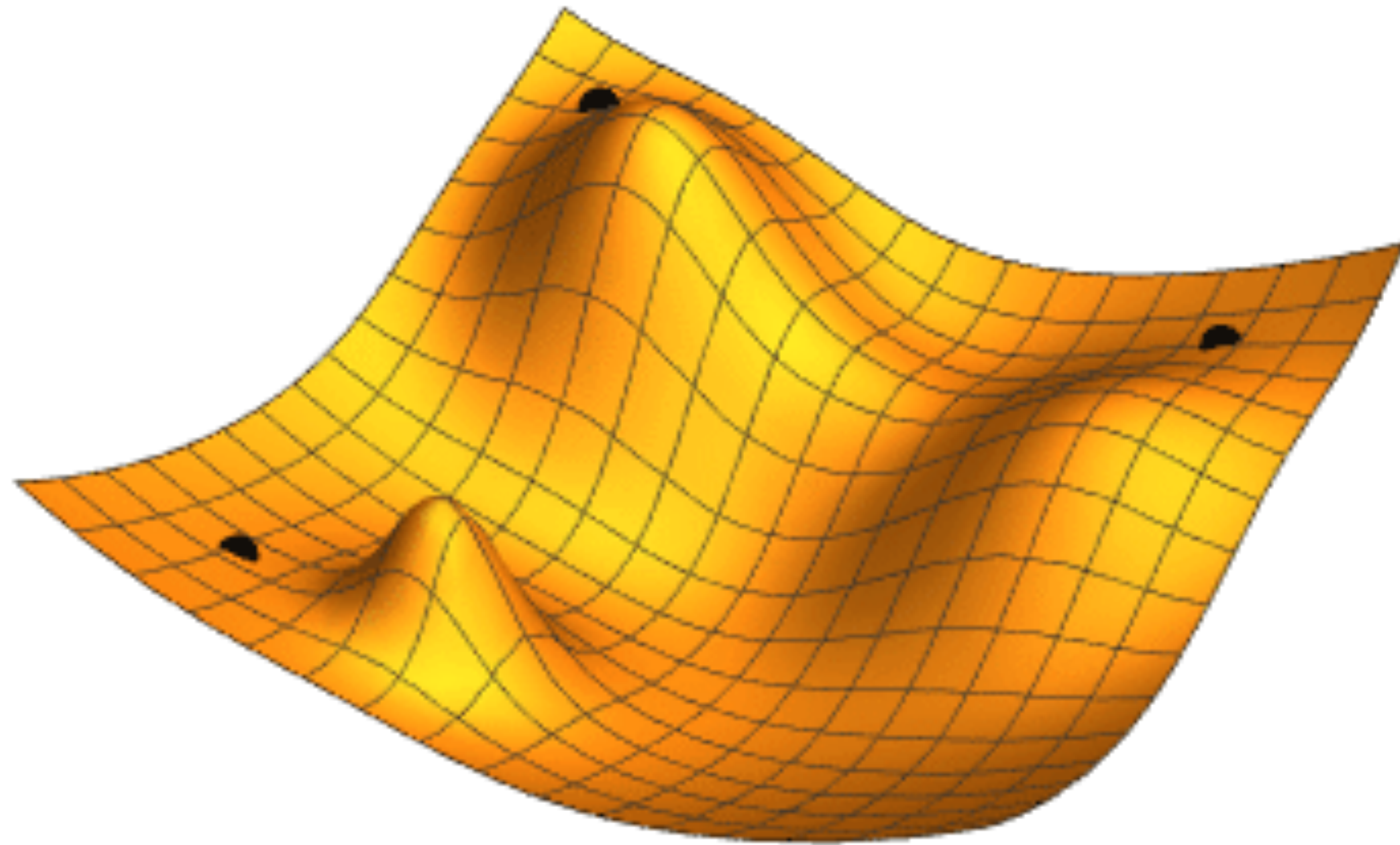
Credits: Wikimedia, Hoang Duong

Gradient Descent Algorithm Visualization



Credits: Wikimedia, Hoang Duong

Gradient Descent Algorithm Visualization



Credits: Wikimedia, Hoang Duong

Linear Regression with Gradient Descent

- Input data: $X \in \mathbb{R}^{d \times n}$, $Y \in \mathbb{R}^n$, where $(\vec{x} \in \mathbb{R}^d, y \in \mathbb{R}^1)$ corresponds to a data point.
 - $n \rightarrow \#$ of data points, $d \rightarrow \#$ of features / input dim.
- Goal: to find $\vec{w} \in \mathbb{R}^d$ such that $\langle \vec{w}, \vec{x} \rangle = y$
 - Minimize $\|X^T \vec{w} - Y\|^2$
- Analytic Solution: $\vec{w} = (XX^T)^{-1}XY$

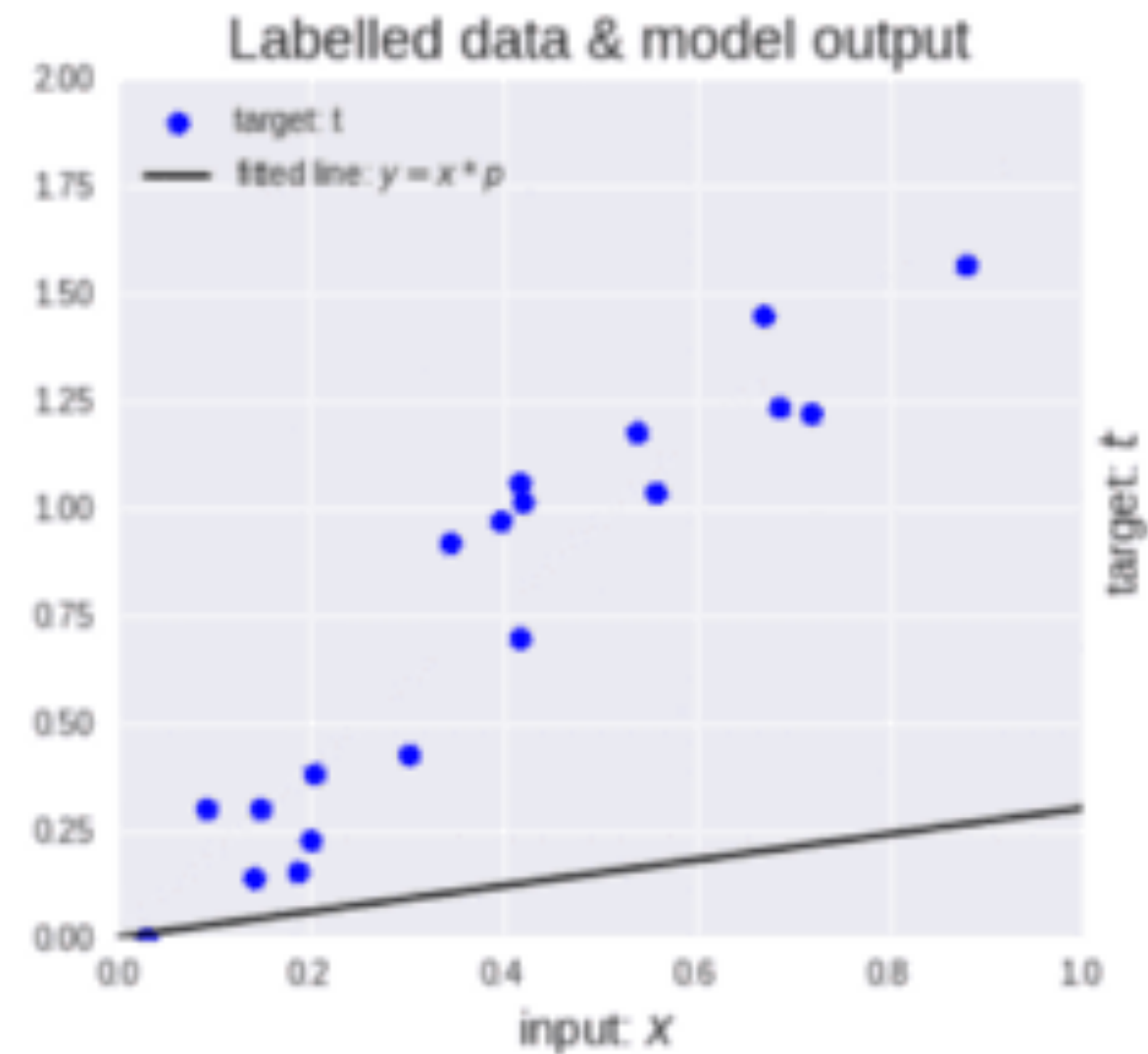
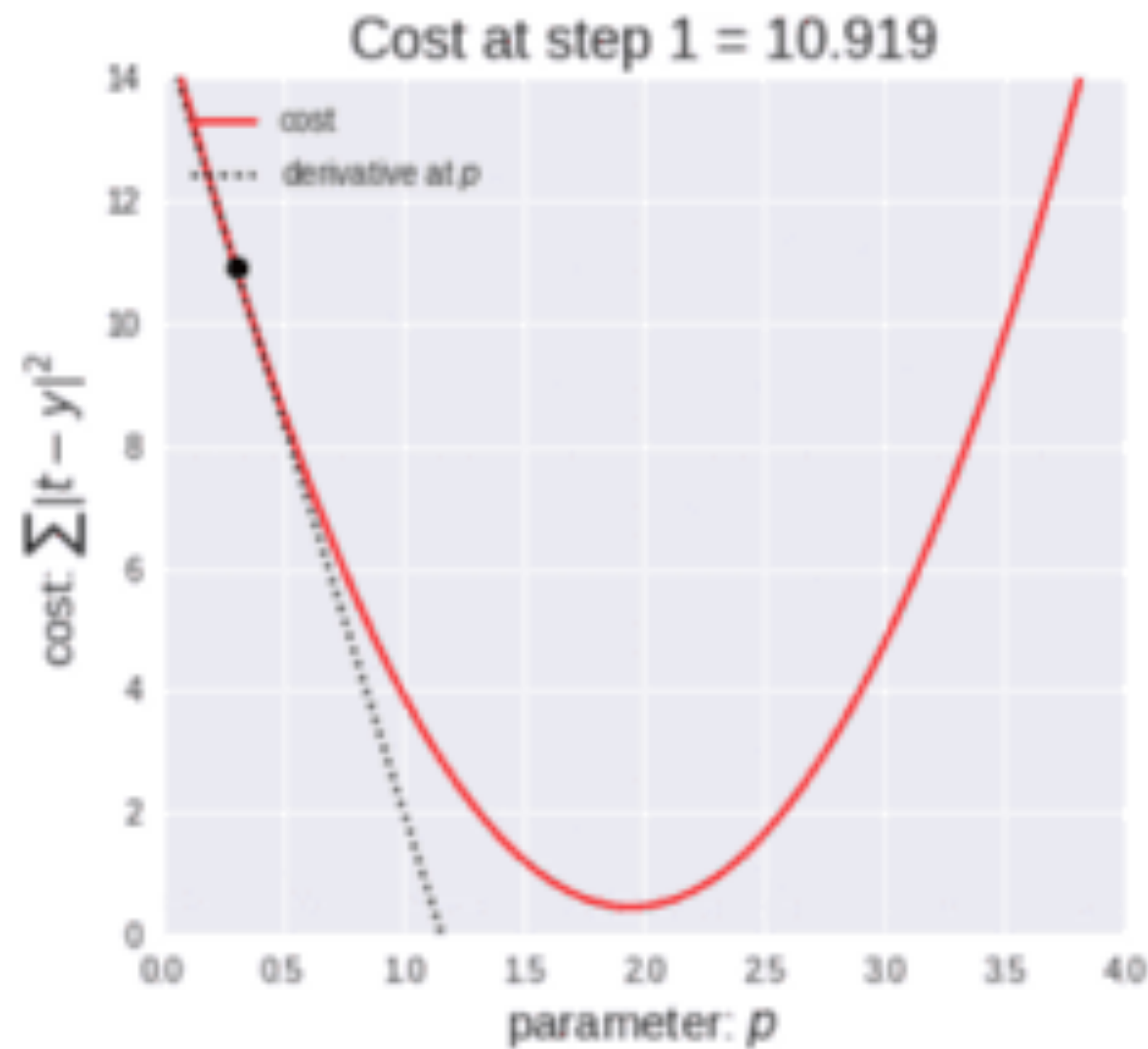
Linear Regression with Gradient Descent

- Input data: $X \in \mathbb{R}^{d \times n}$, $Y \in \mathbb{R}^n$, where $(\vec{x} \in \mathbb{R}^d, y \in \mathbb{R}^1)$ corresponds to a data point.
 - $n \rightarrow \#$ of data points, $d \rightarrow \#$ of features / input dim.
- Goal: to find $\vec{w} \in \mathbb{R}^d$ such that $\langle \vec{w}, \vec{x} \rangle = y$
 - Minimize $\|X^T \vec{w} - Y\|^2$
- Loss function $f(\vec{w}) = \|X^T \vec{w} - Y\|^2 = (X^T \vec{w} - Y)^T (X^T \vec{w} - Y)$

Linear Regression with Gradient Descent

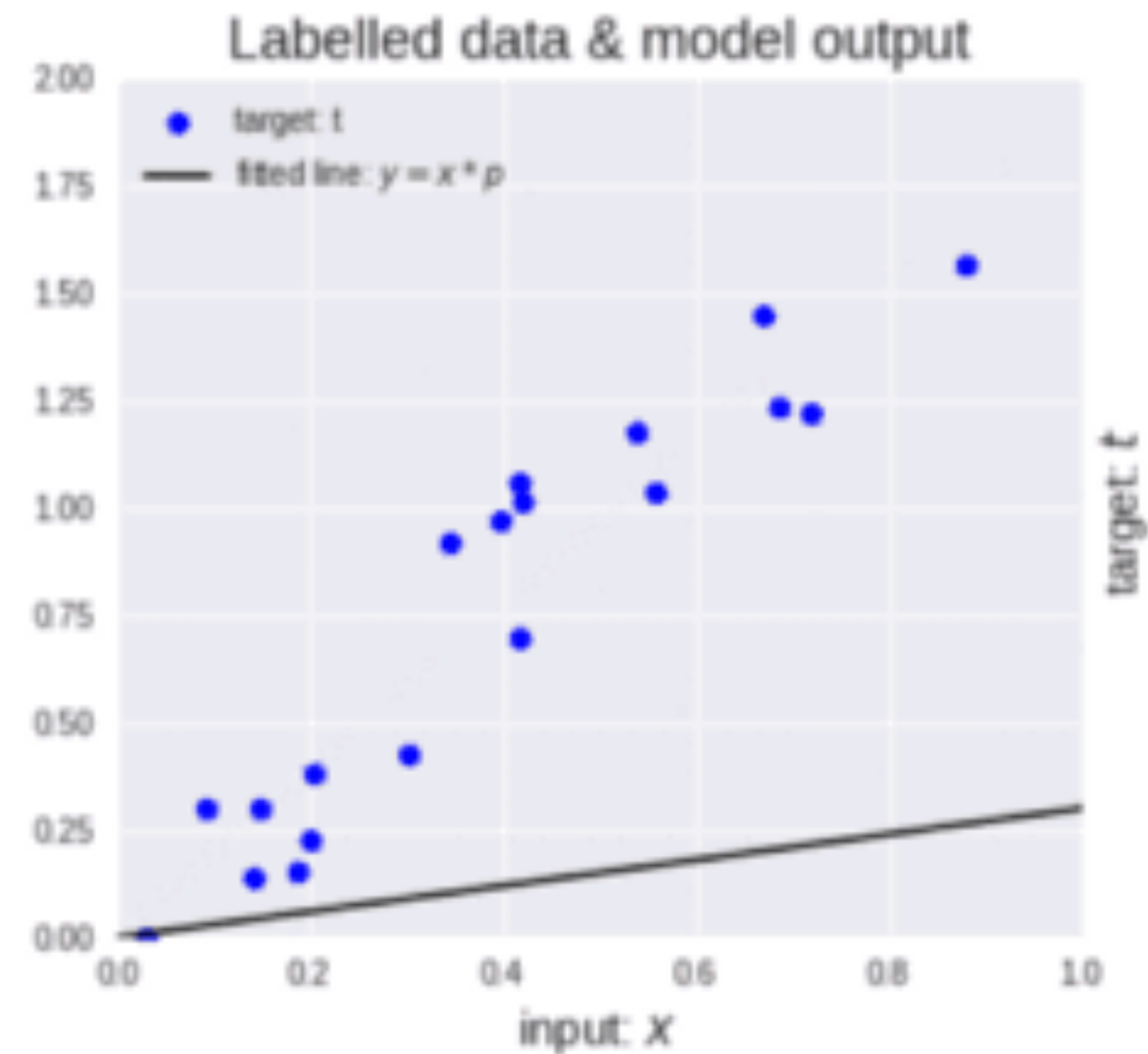
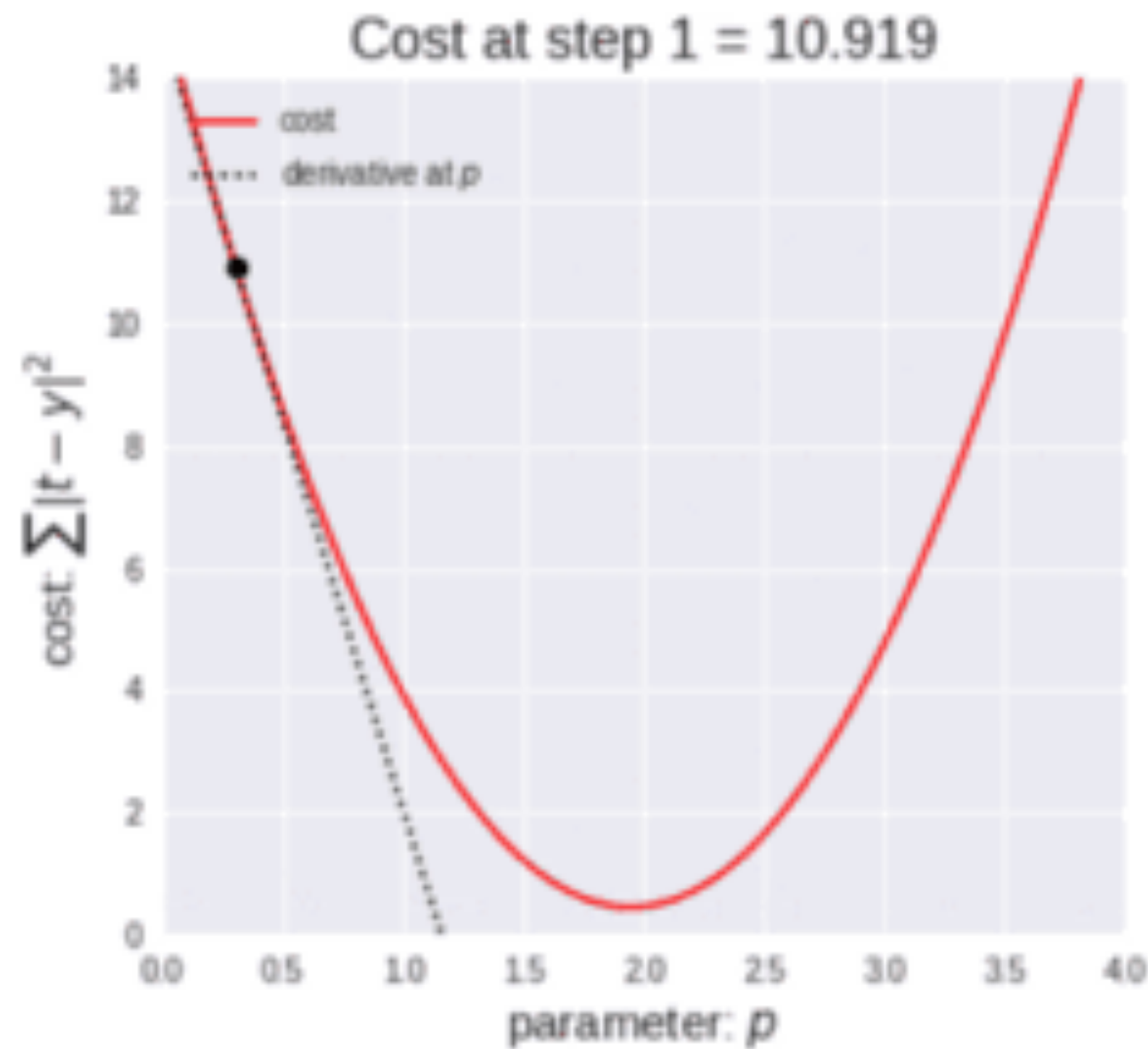
- Loss function $f(\vec{w}) = \|X^T \vec{w} - Y\|^2 = (X^T \vec{w} - Y)^T (X^T \vec{w} - Y)$
 - Set $f(\vec{w}) = \sum_{i=1}^n (x^i \cdot \vec{w} - y^i)^2$
 - Compute $\frac{\partial f}{\partial w_k} = \sum_{i=1}^n 2(\vec{w}^T x^i - y^i) x_k^i$
 - Compute $\nabla f = [\frac{\partial f}{\partial w_1}, \dots, \frac{\partial f}{\partial w_k}, \dots, \frac{\partial f}{\partial w_d}]$
 - Use gradient descent

Linear Regression with Gradient Descent



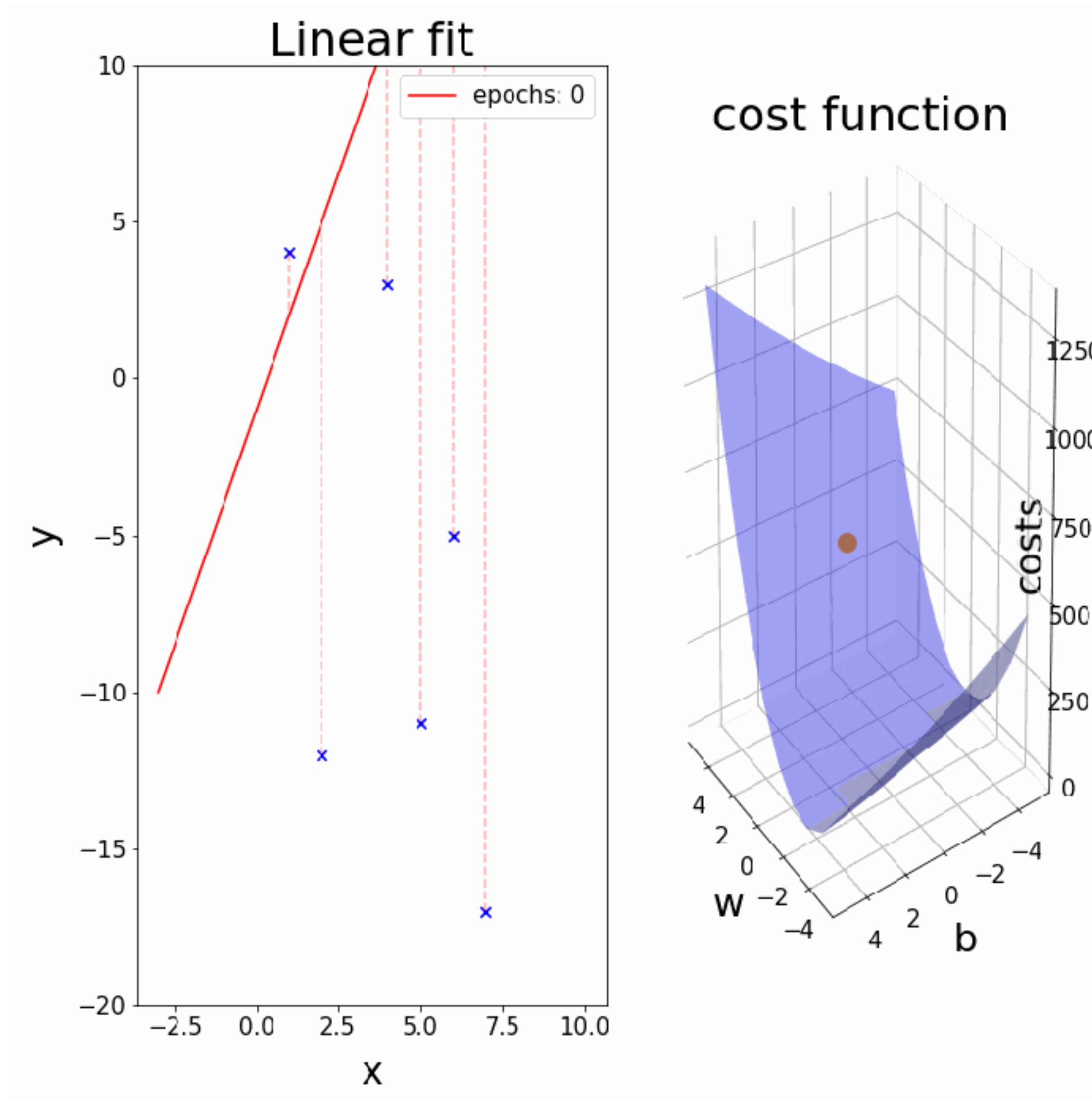
Credits: <https://www.kdnuggets.com/>

Linear Regression with Gradient Descent



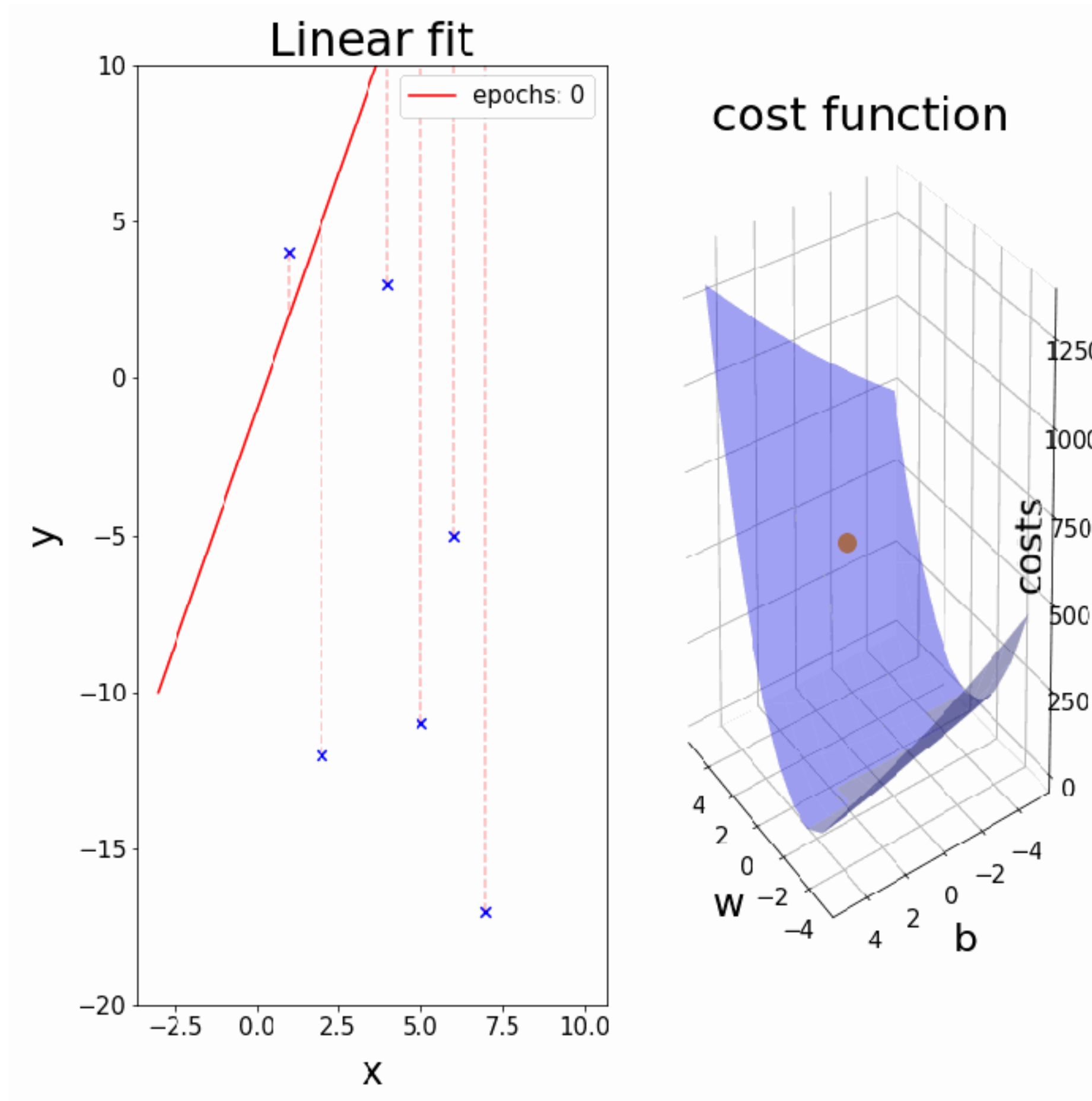
Credits: <https://www.kdnuggets.com/>

Linear Regression with Gradient Descent



Credits: Tobias Roeschl

Linear Regression with Gradient Descent

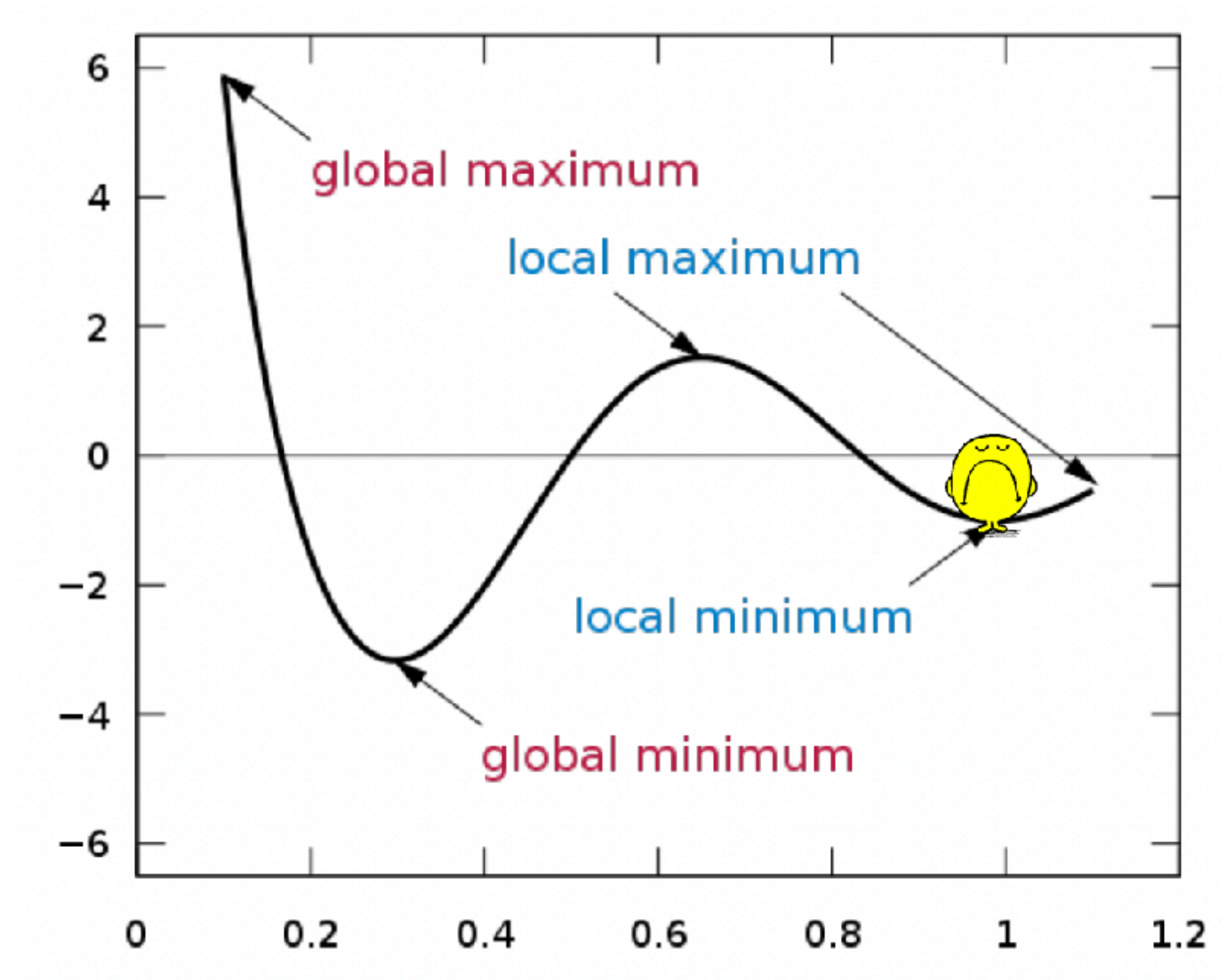


Credits: Tobias Roeschl

Gradient Descent Algorithm for Machine Learning

- Given: cost / loss/ objective function $f(\vec{\theta}, D)$. Where $\vec{\theta} \in \mathbb{R}^d$.
- Goal: find $\vec{\theta}^*$ such that $f(\vec{\theta}^*, D) = \min_{\vec{\theta}} f(\vec{\theta}, D)$.
- Gradient descent solution:
 - Start from initial guess $\vec{\theta}^0$ and learning rate α
 - Update $\vec{\theta}^{i+1} \leftarrow \vec{\theta}^i - \alpha \nabla f(\vec{\theta}^i, D)$
 - Repeat until change in θ is small, or maximum number of steps reached.

Key issue: Local minima

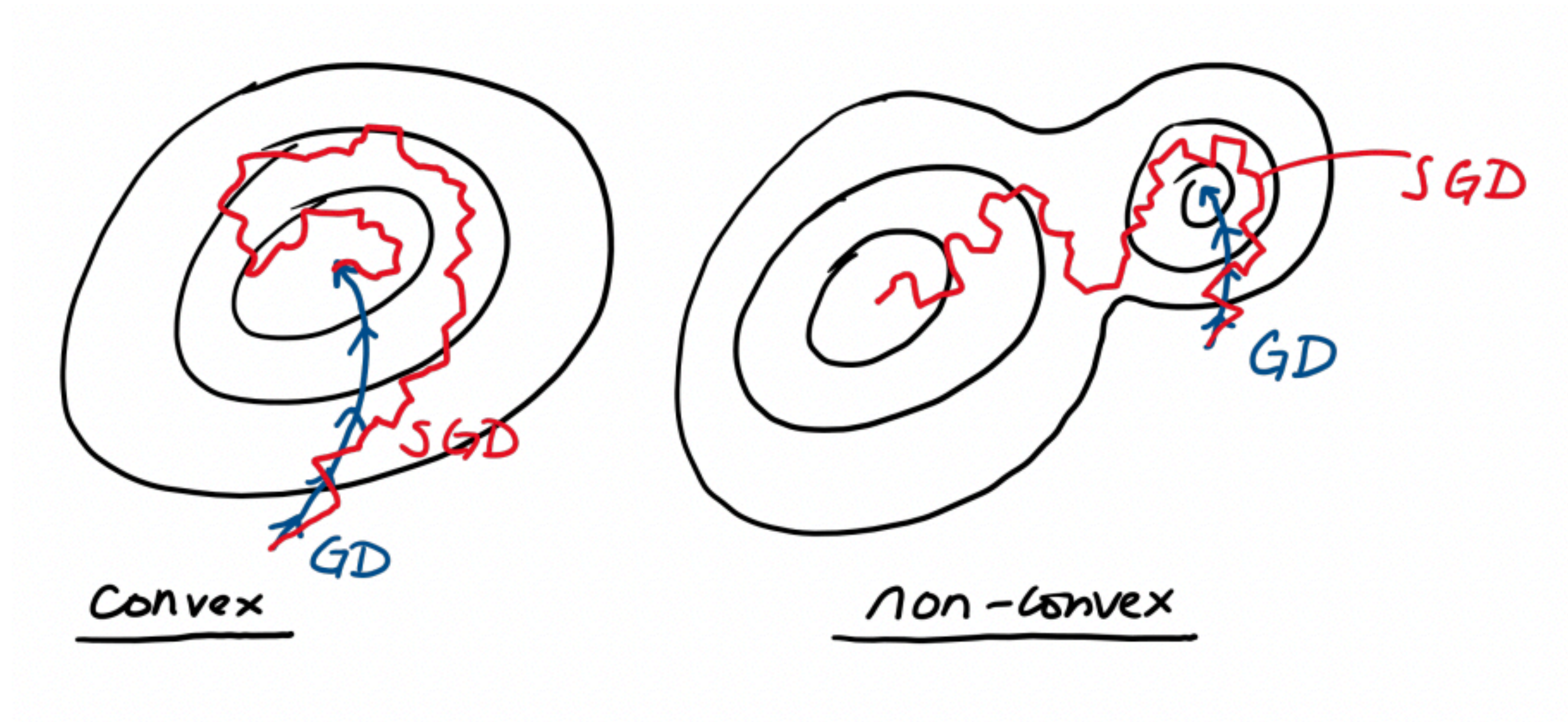


Credits: Michail Michailidis & Patrick Maiden

One solution: Stochastic Gradient Descent (SGD)

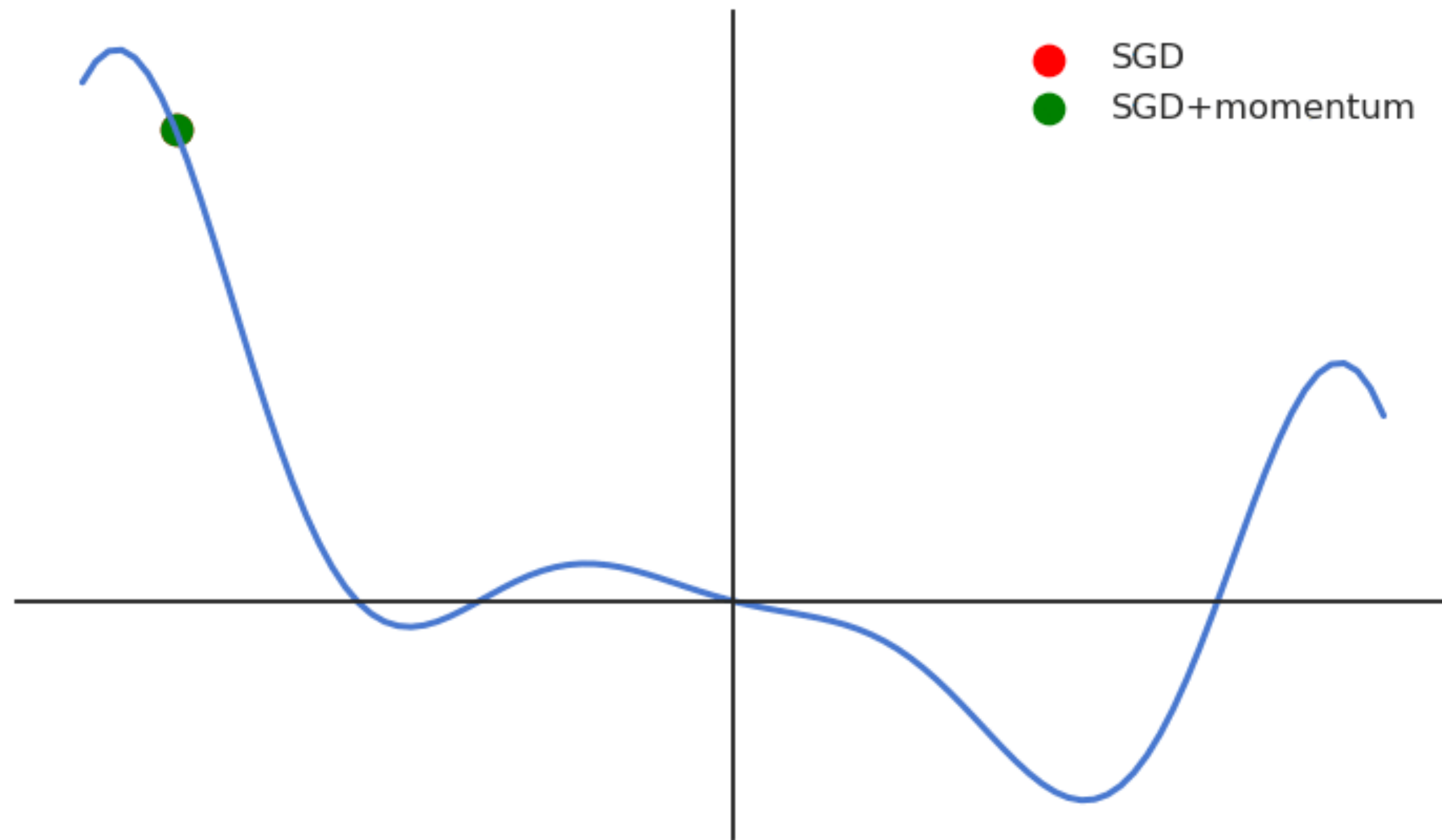
- Gradient descent solution:
 - Start from initial guess $\vec{\theta}^0$ and learning rate α
 - Update $\vec{\theta}^{i+1} \leftarrow \vec{\theta}^i - \alpha \nabla f(\vec{\theta}, D)$
- SGD:
 - Sample single or multiple datapoints $d \sim D$
 - Update $\vec{\theta}^{i+1} \leftarrow \vec{\theta}^i - \alpha \nabla f(\vec{\theta}, d)$

One solution: Stochastic Gradient Descent (SGD)



Credits: Stanley Chan

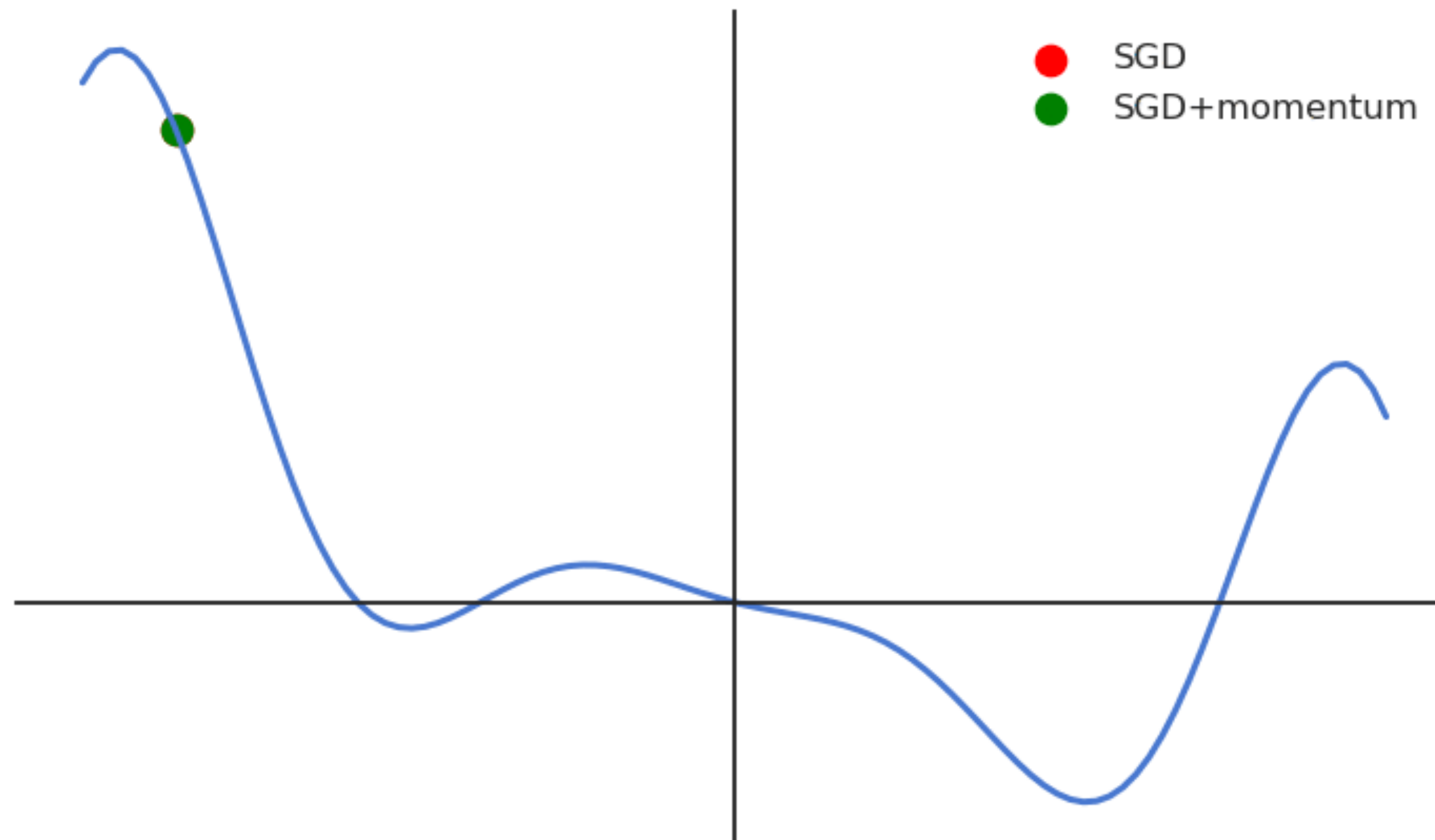
Another solution: Momentum



$$v^{t+1} = \beta v^t + \alpha \nabla f(\vec{w})$$

$$\vec{\theta}^{t+1} \leftarrow \vec{\theta}^t - v^{t+1}$$

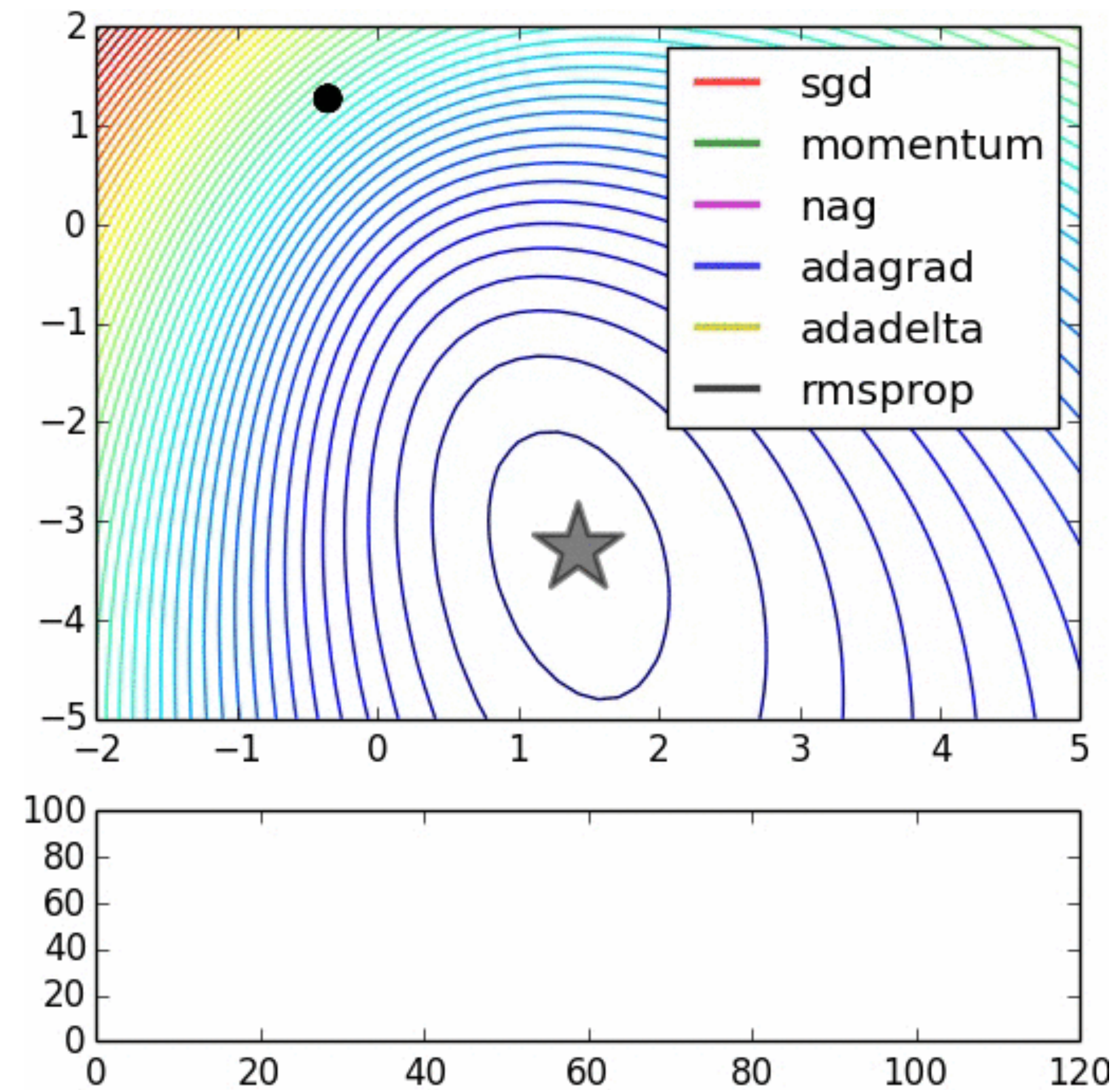
Another solution: Momentum



$$v^{t+1} = \beta v^t + \alpha \nabla f(\vec{w})$$

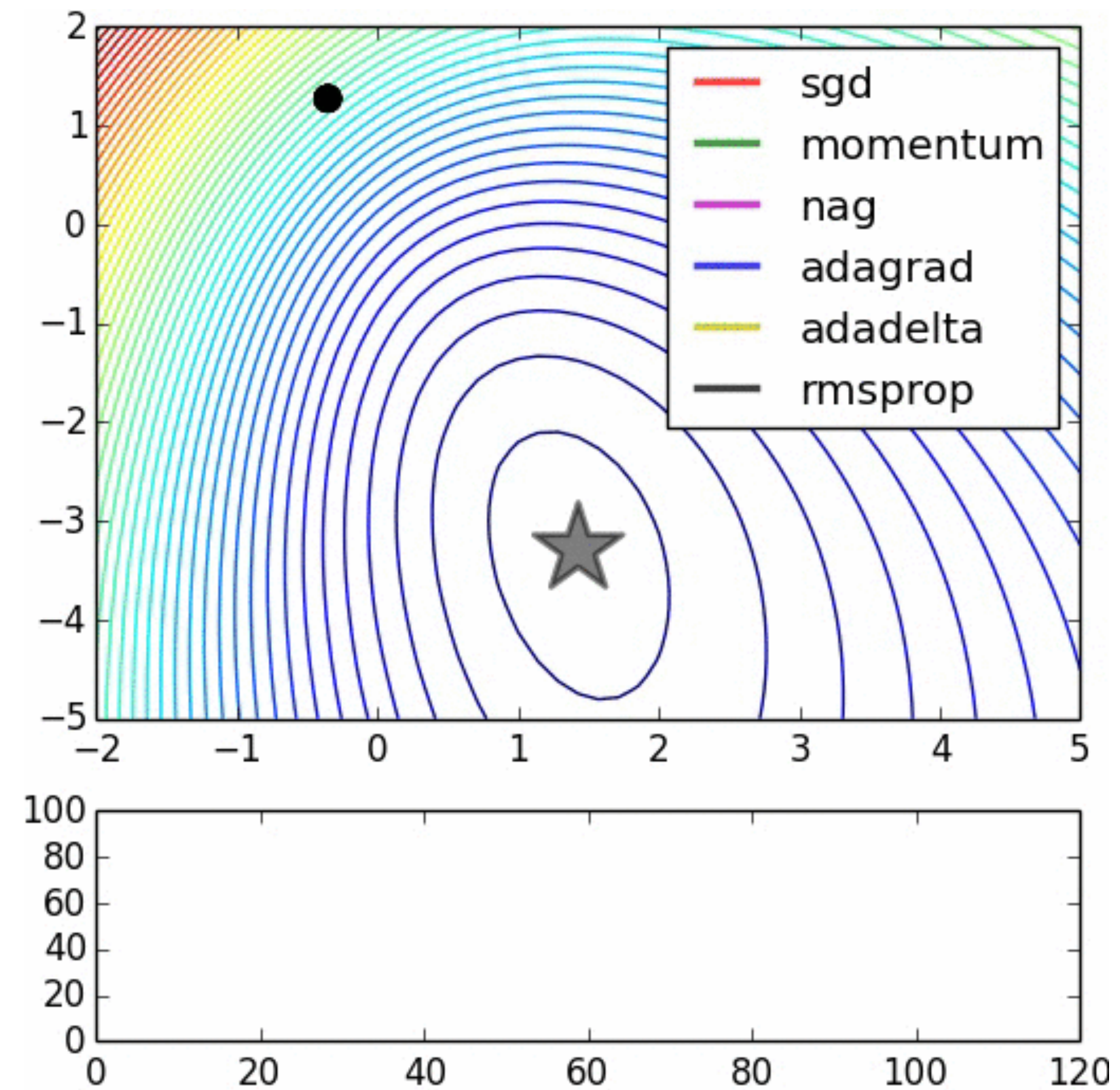
$$\vec{\theta}^{t+1} \leftarrow \vec{\theta}^t - v^{t+1}$$

More solutions have been proposed!



Credits: Alec Radford

More solutions have been proposed!



Credits: Alec Radford

Additional Reading

- Interactive tutorial: <https://uclaacm.github.io/gradient-descent-visualiser/>
- Book chapter: <https://www.cs.utah.edu/~jeffp/IDABook/T6-GD.pdf>
- SGD + variants: <https://runder.io/optimizing-gradient-descent>

Questions?