# 6.036: Introduction to Machine Learning

**Lecture start:** Tuesdays 9:35am

**Who's talking?** Prof. Tamara Broderick

**Questions?** Ask on Piazza: "lecture (week) 3" folder

**Materials:** slides, video will all be available on Canvas

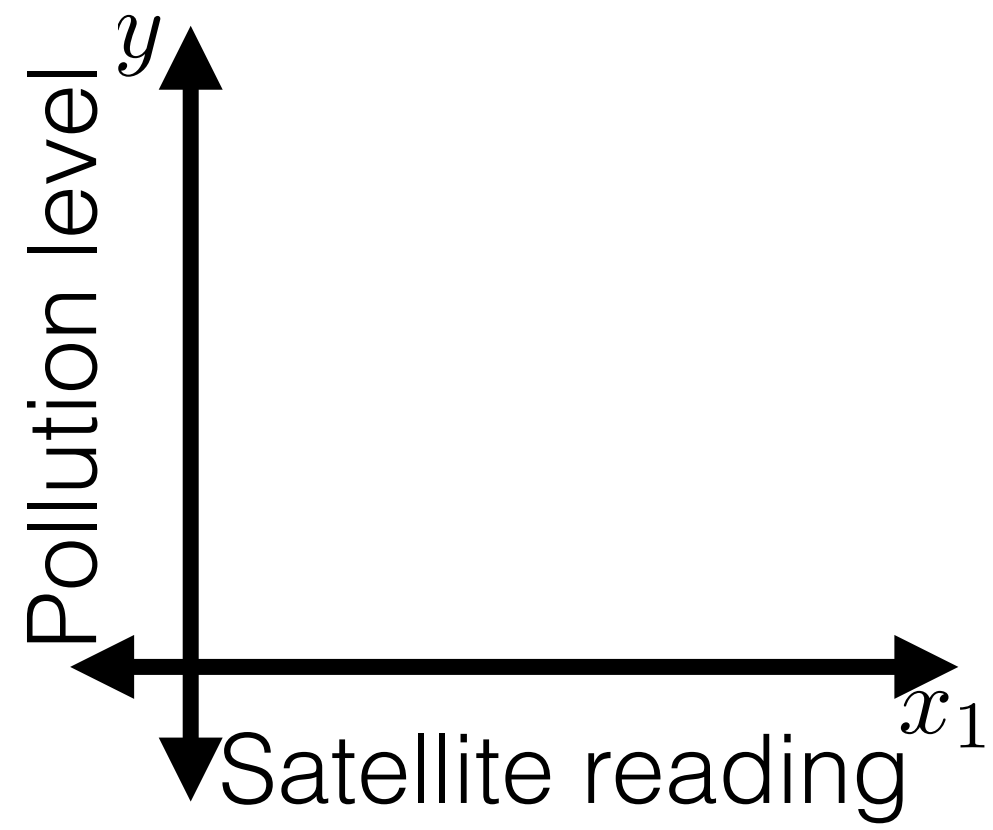**Live Zoom feed:** https://mit.zoom.us/j/94238622313

## Last Time

I. Machine learning setup

II. Linear regression

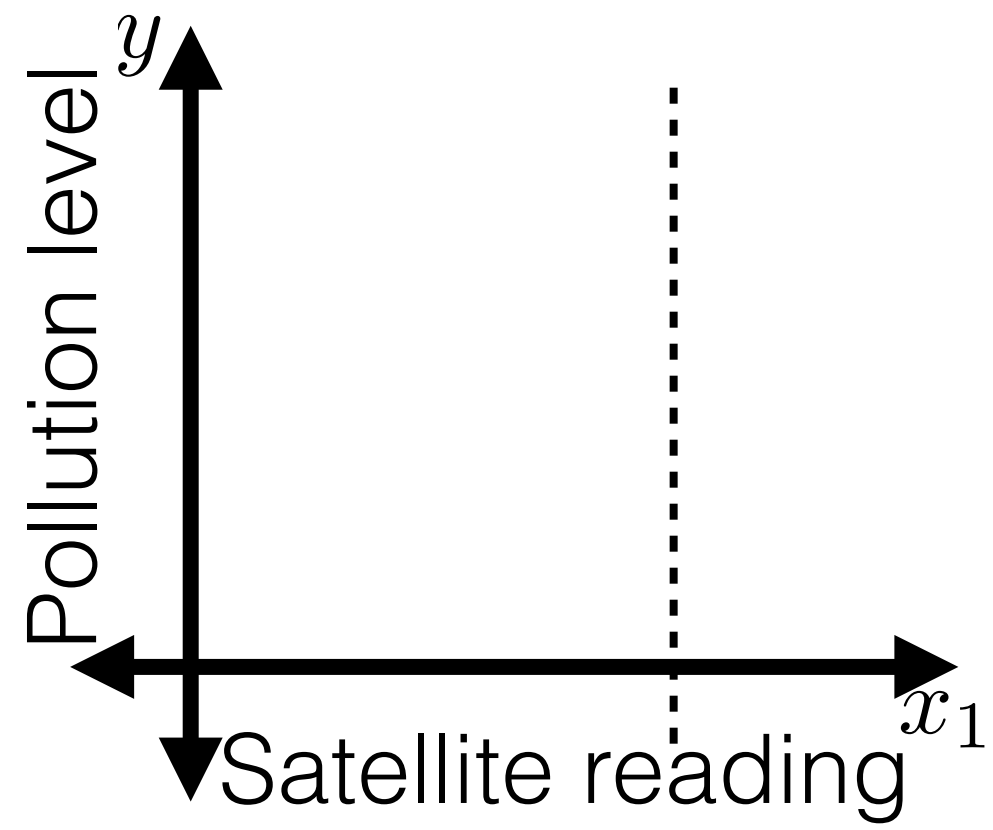III. Regularization

## Today's Plan

I. Gradient descent

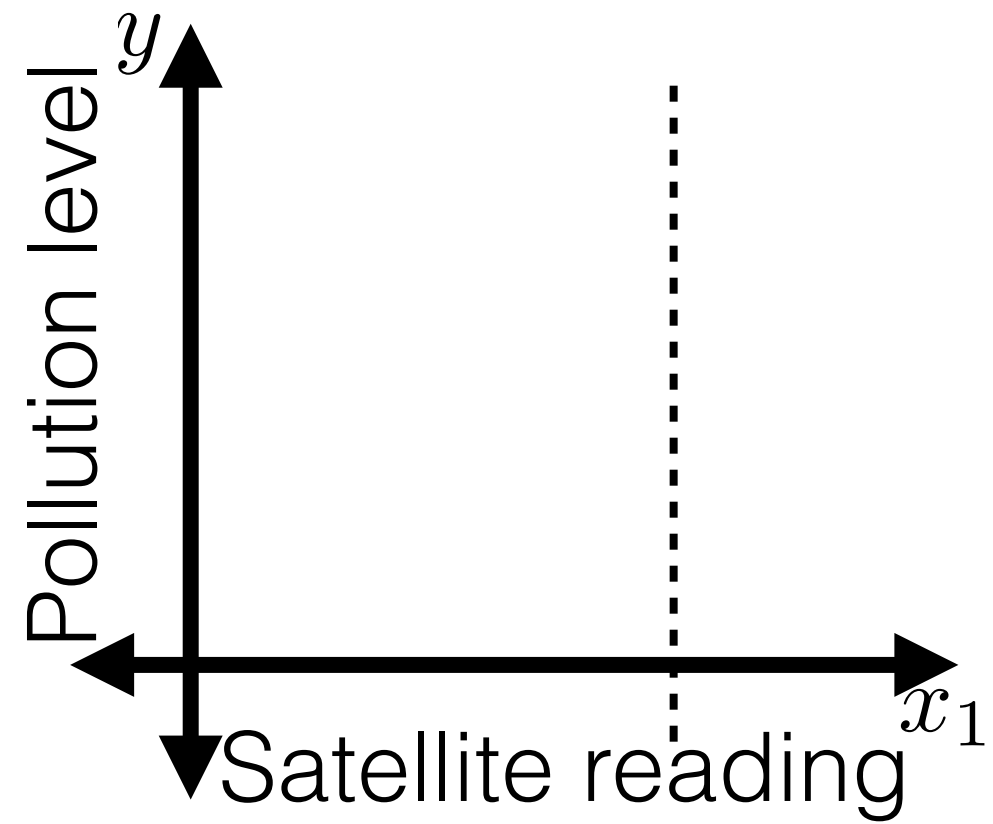II. Stochastic gradient descent (SGD)
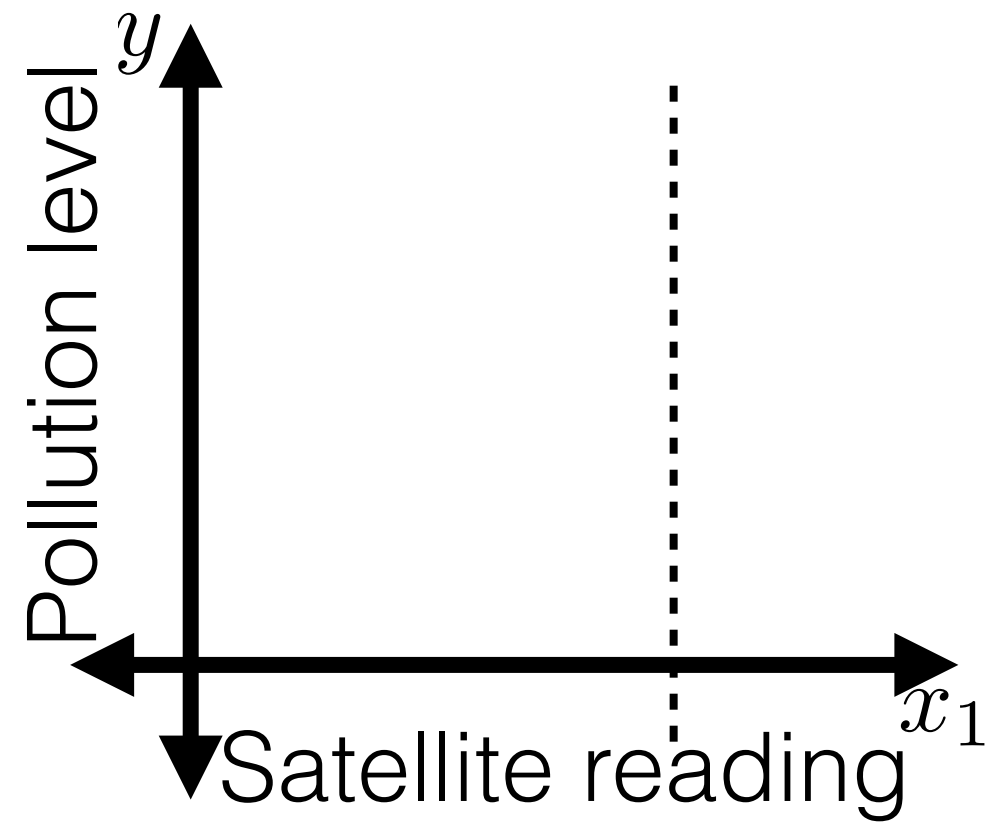
# Recall

# Recall

# Recall

# Recall

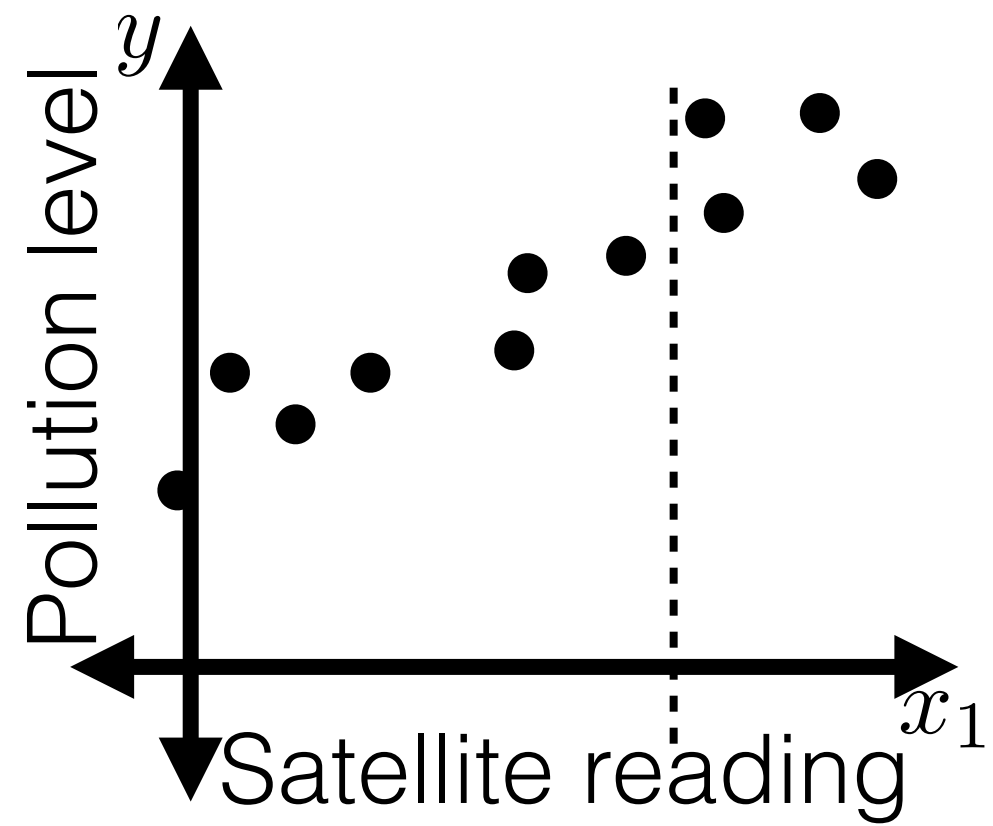- A general ML approach:

# Recall

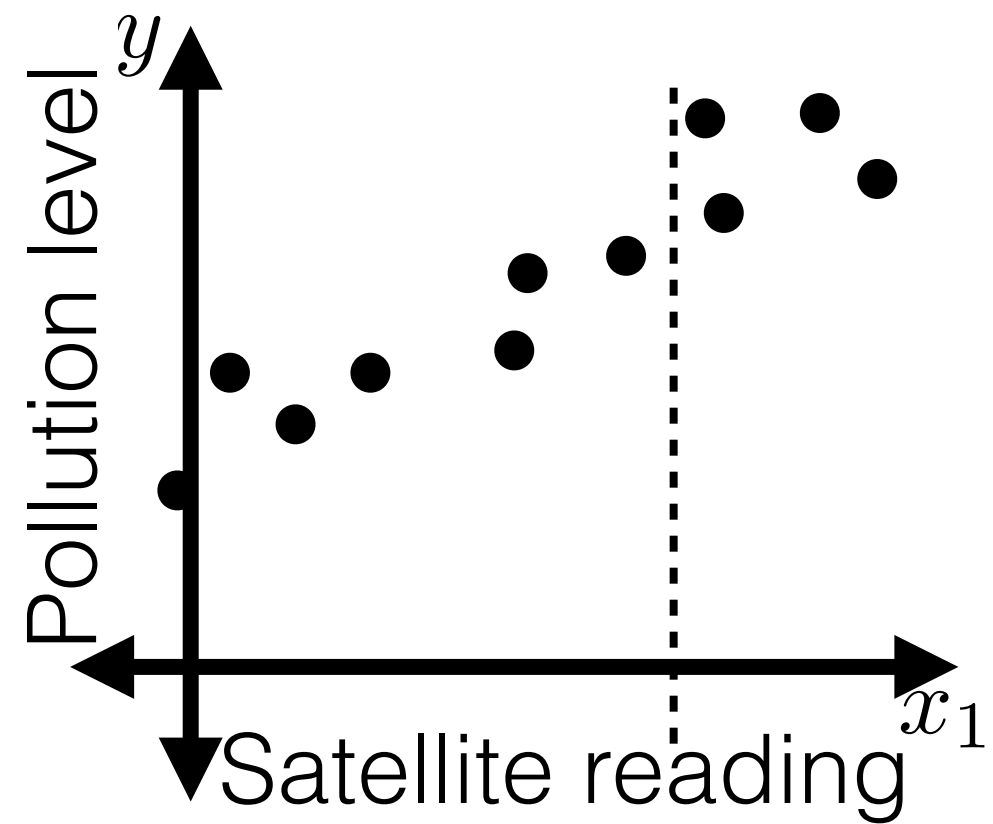- A general ML approach:
  - Collect data

# Recall

- A general ML approach:
  - Collect data

# Recall

- A general ML approach:
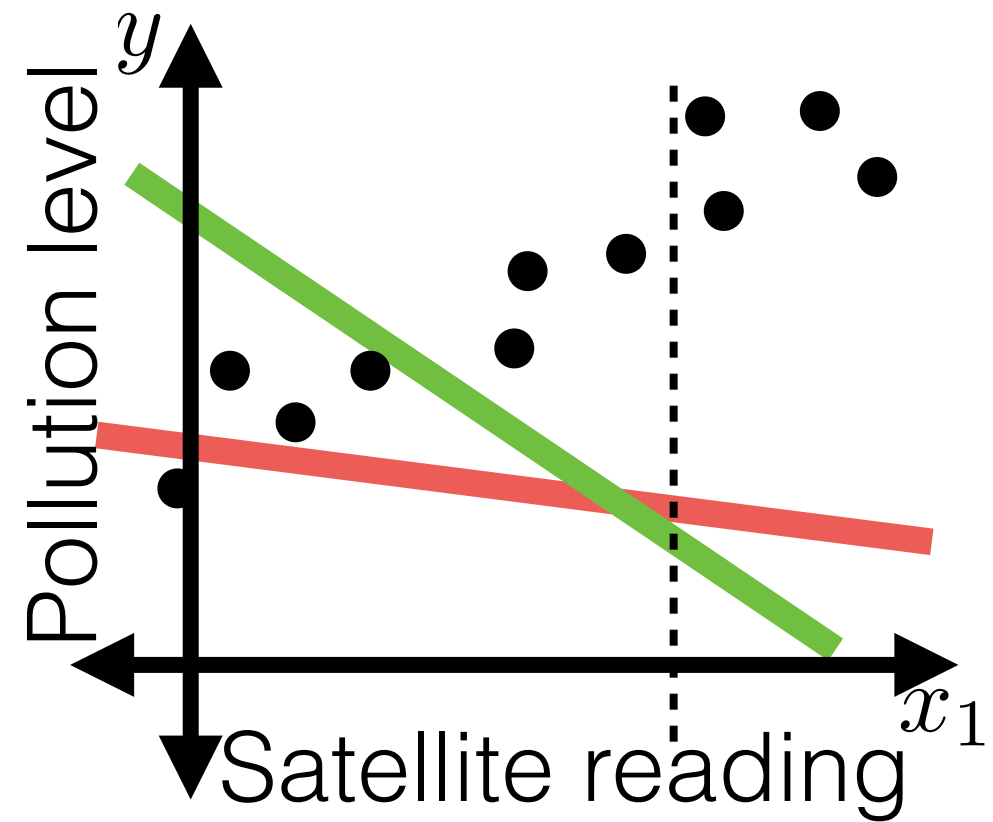  - Collect data
  - Choose hypothesis class
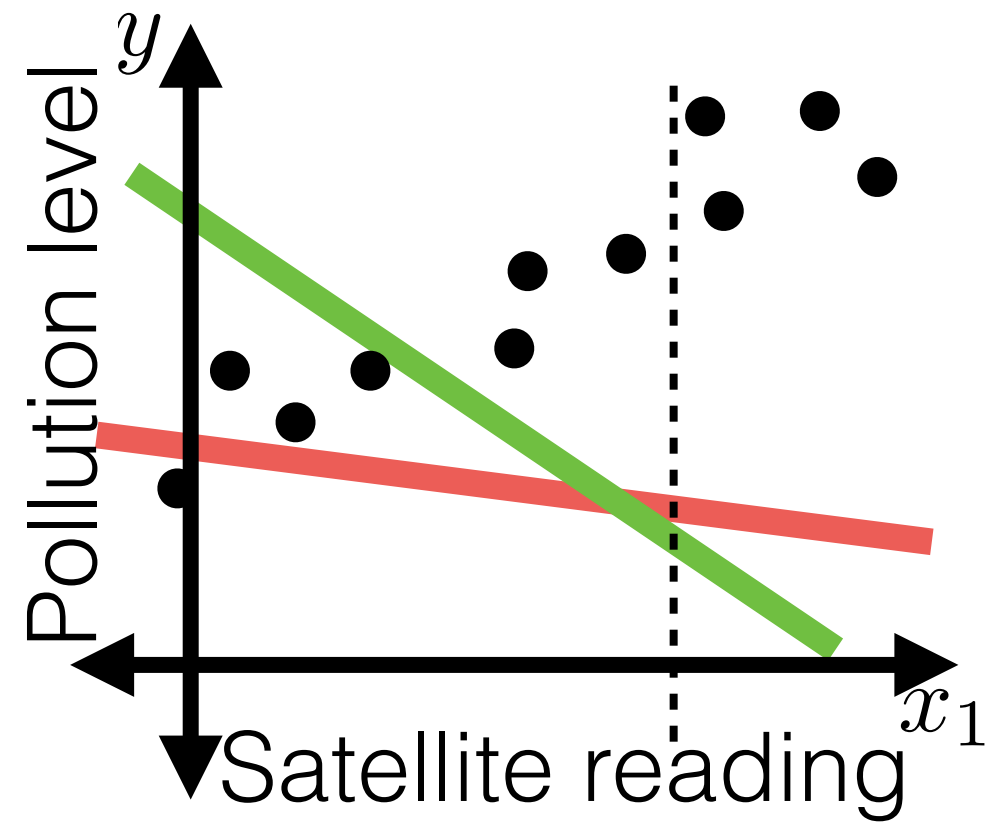
# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
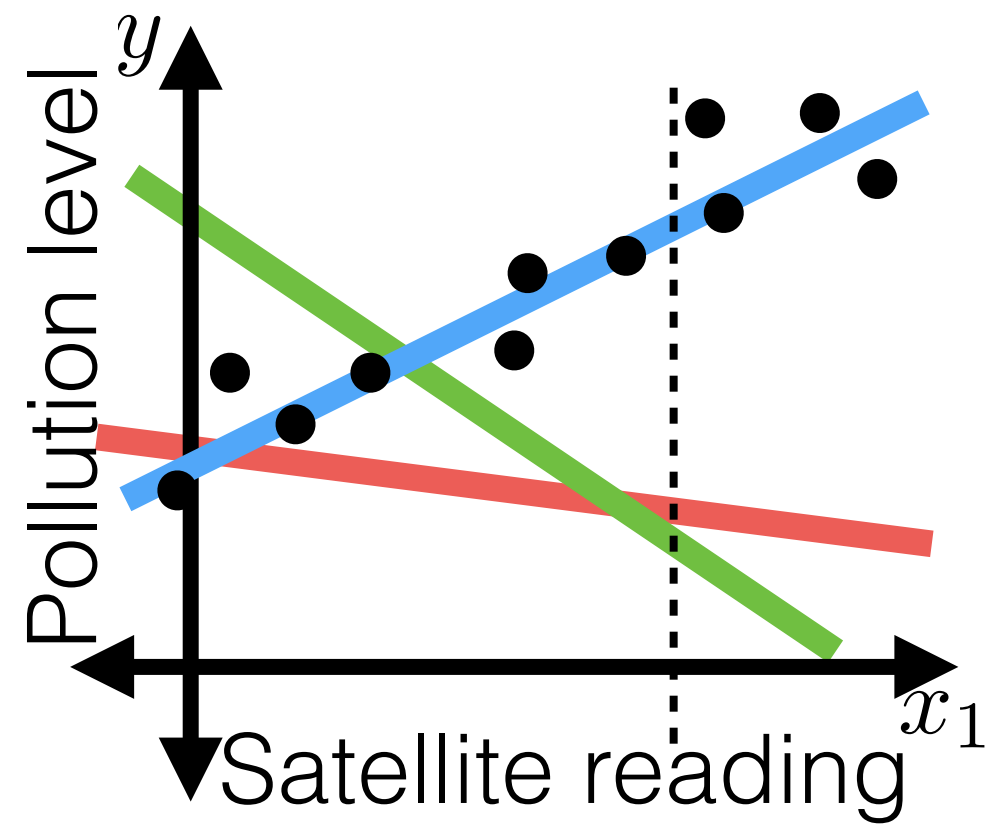  - Choose "good" hypothesis by minimizing training loss + regularizer

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer



$$\frac{1}{n}\sum_{i=1}^{n} L(h(x^{(i)};\Theta), y^{(i)}) + \lambda R(\Theta) \quad (\lambda > 0)$$

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer

$$\frac{1}{n}\sum_{i=1}^{n} L(h(x^{(i)};\Theta), y^{(i)}) + \lambda R(\Theta) \quad (\lambda > 0)$$

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
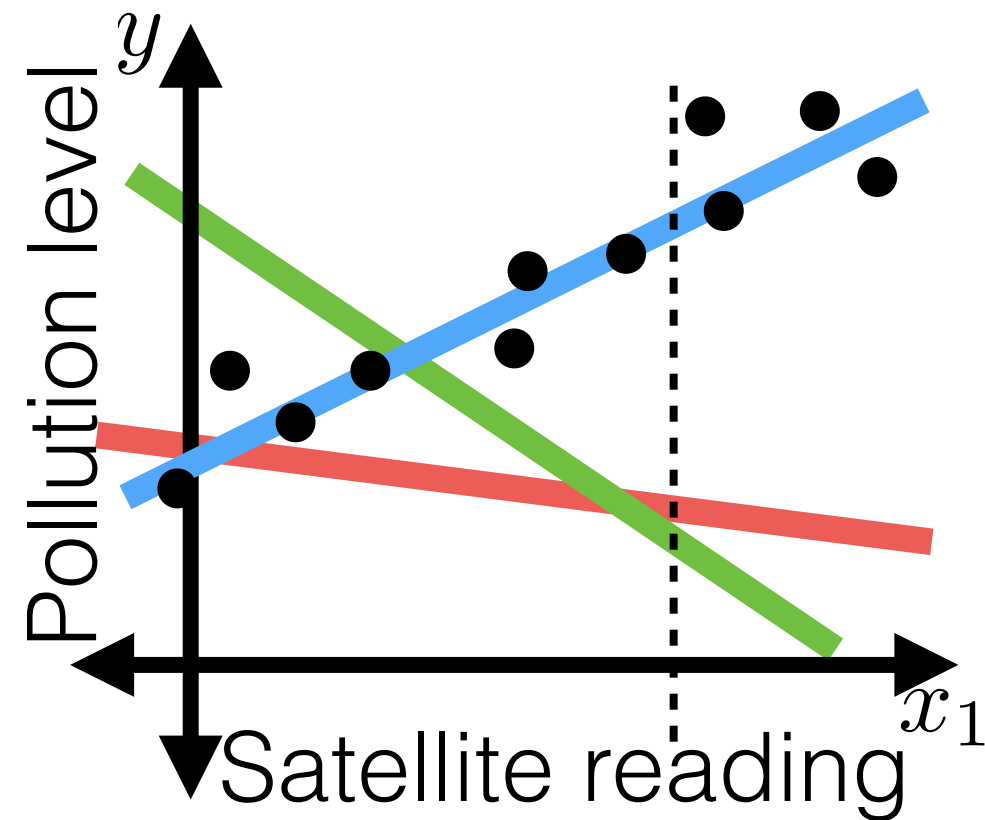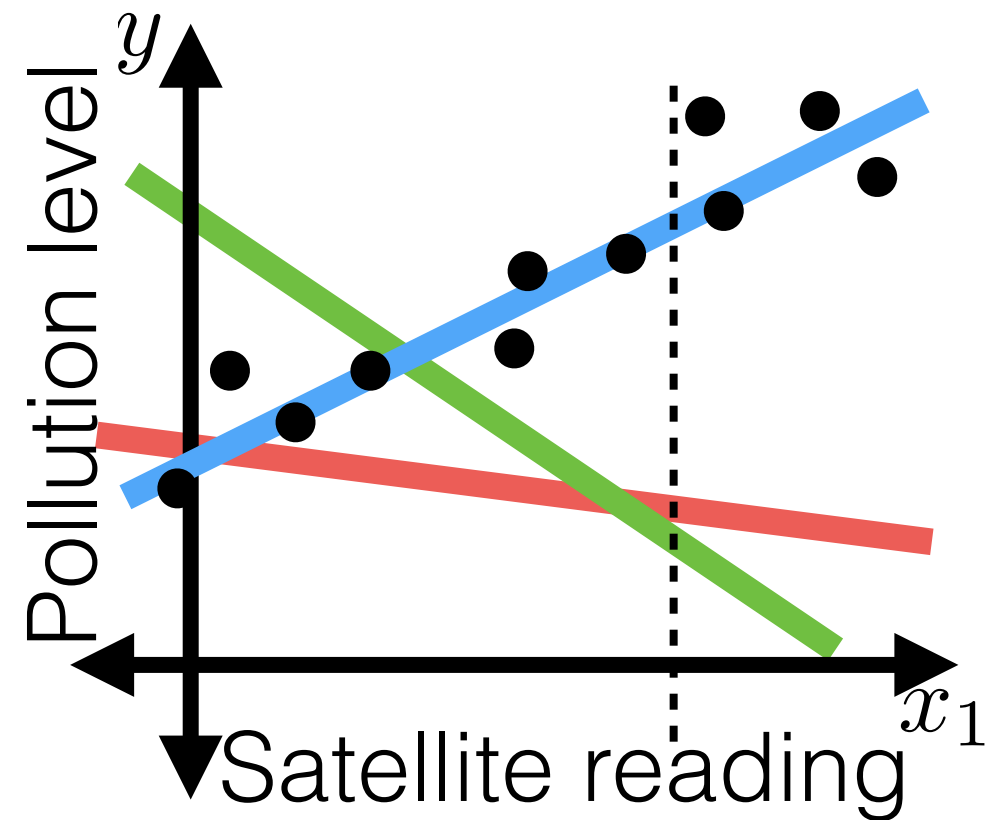  - Choose "good" hypothesis by minimizing training loss + regularizer



$$\frac{1}{n}\sum_{i=1}^{n} L(h(x^{(i)};\Theta),y^{(i)}) + \lambda R(\Theta) \qquad (\lambda > 0)$$
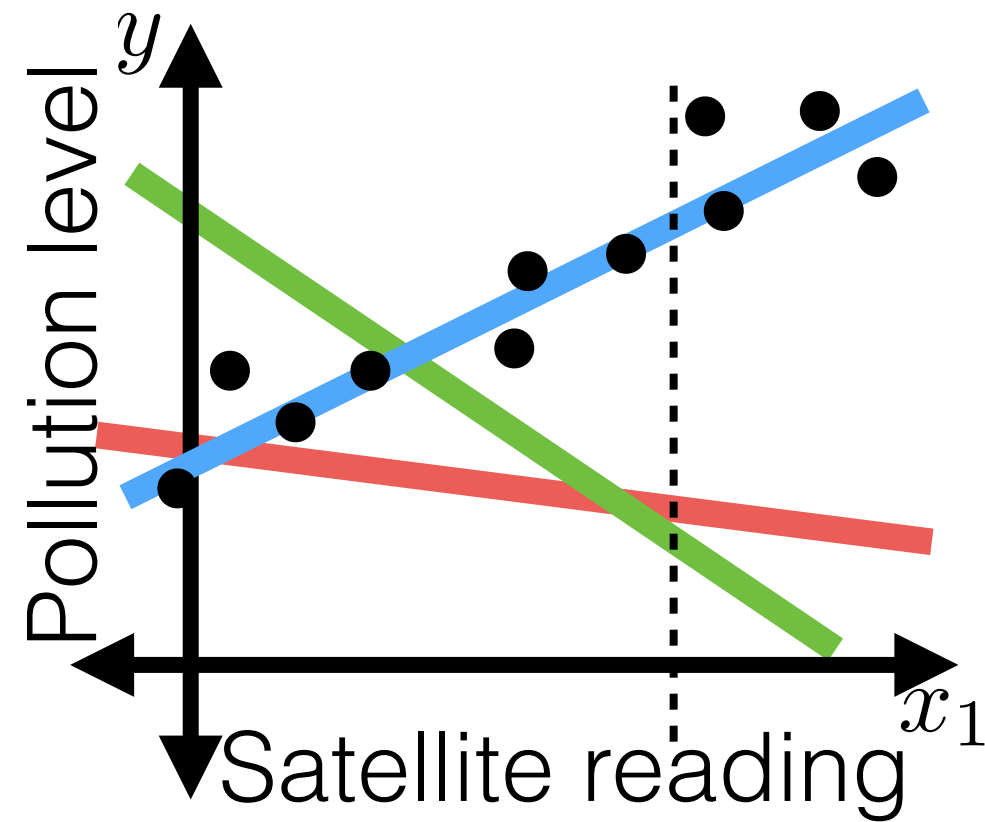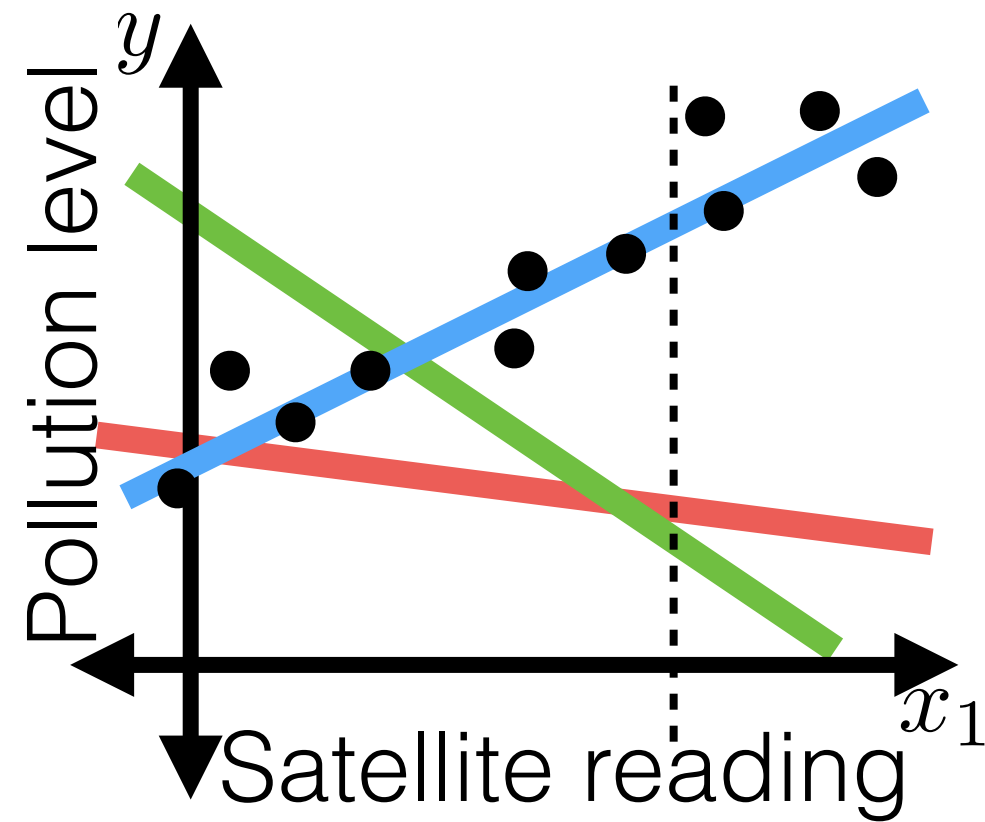
# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer



$$\frac{1}{n} \sum_{i=1}^{n} L(h(x^{(i)}; \Theta), y^{(i)}) + \lambda R(\Theta) \qquad (\lambda > 0)$$

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
- Example: ridge regression

$$\frac{1}{n}\sum_{i=1}^{n} L(h(x^{(i)}; \Theta), y^{(i)}) + \lambda R(\Theta) \quad (\lambda > 0)$$

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
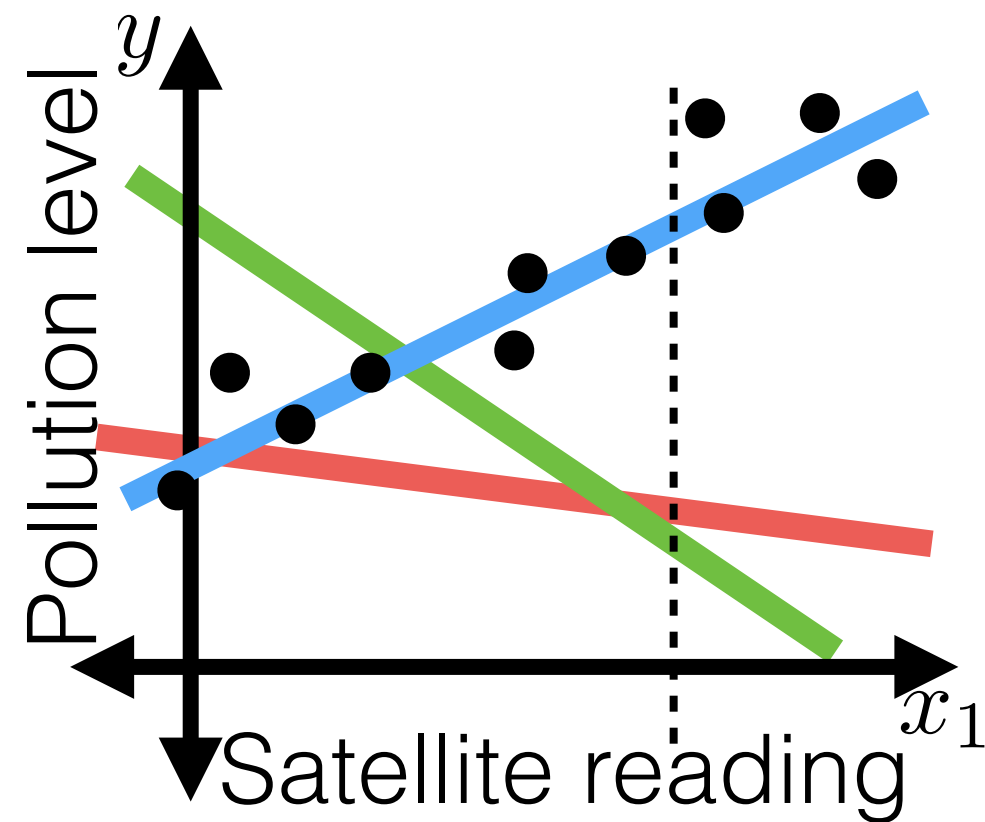- Example: ridge regression

$$\frac{1}{n} \sum_{i=1}^{n} L(h(x^{(i)}; \Theta), y^{(i)}) + \lambda R(\Theta) \quad (\lambda > 0)$$

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
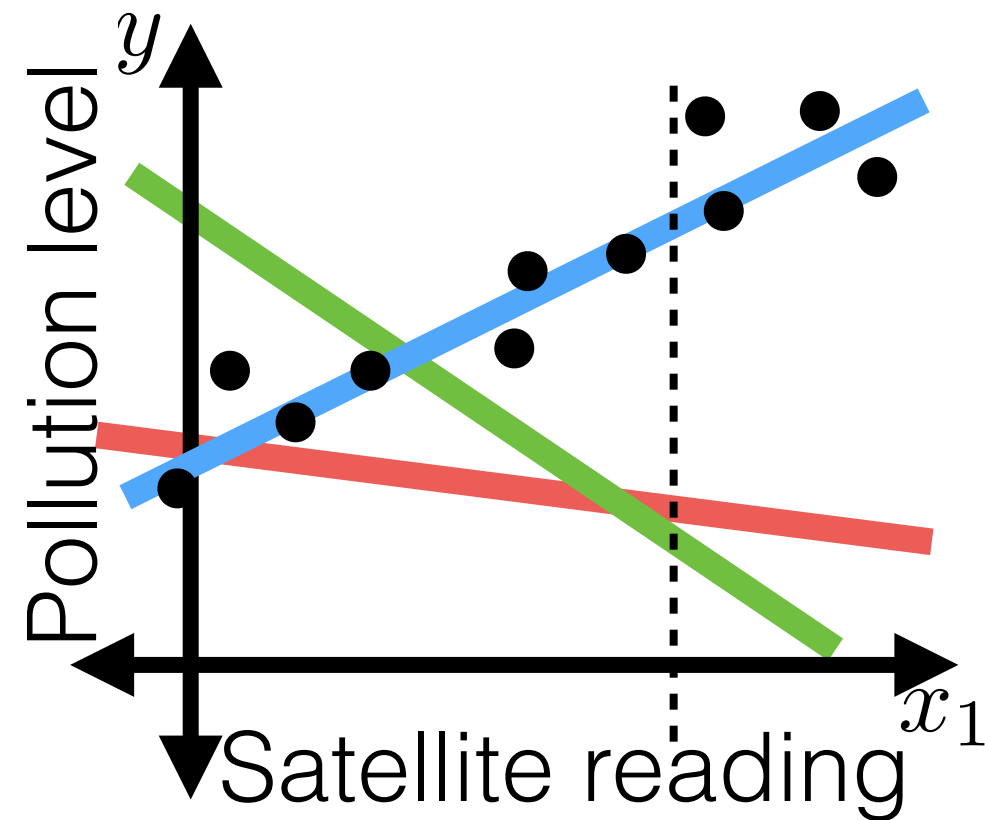- Example: ridge regression



linear regression hypothesis

$$\frac{1}{n} \sum_{i=1}^{n} L(h(x^{(i)}; \Theta), y^{(i)}) + \lambda R(\Theta) \quad (\lambda > 0)$$

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
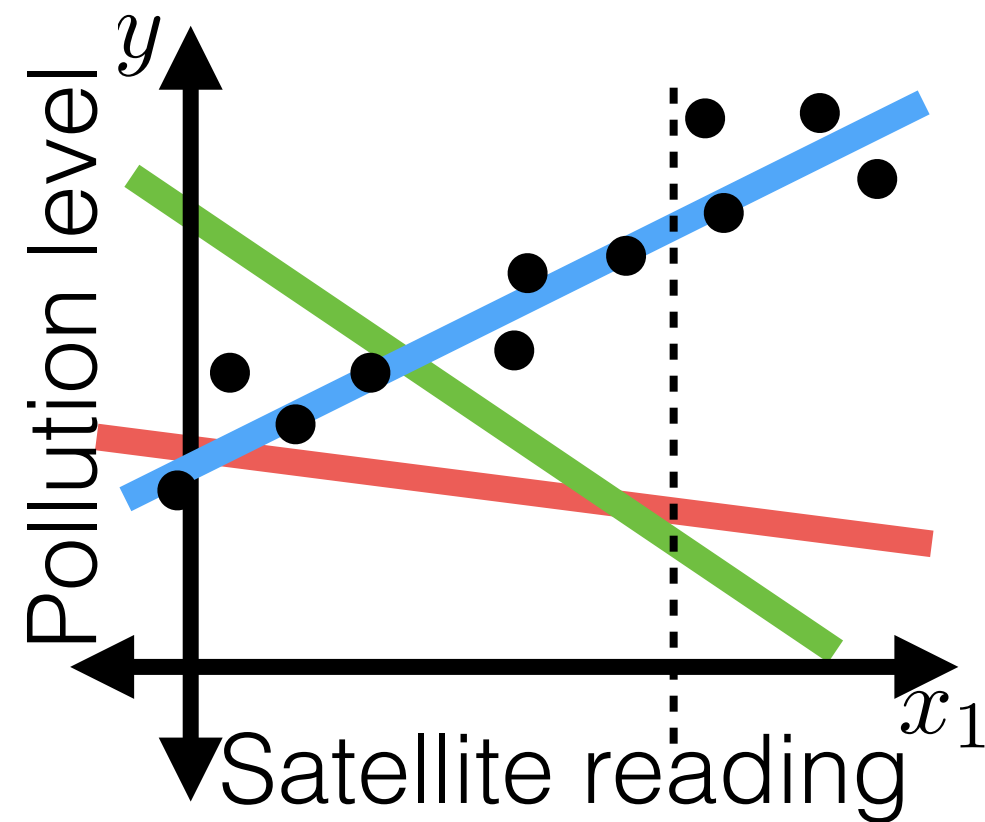- Example: ridge regression



linear regression hypothesis

$$\frac{1}{n} \sum_{i=1}^{n} L(h(x^{(i)}; \Theta), y^{(i)}) + \lambda R(\Theta) \quad (\lambda > 0)$$

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
- Example: ridge regression

$$\frac{1}{n} \sum_{i=1}^{n} L(\boxed{\theta^\top x^{(i)} + \theta_0}, y^{(i)}) + \lambda R(\Theta) \quad (\lambda > 0)$$
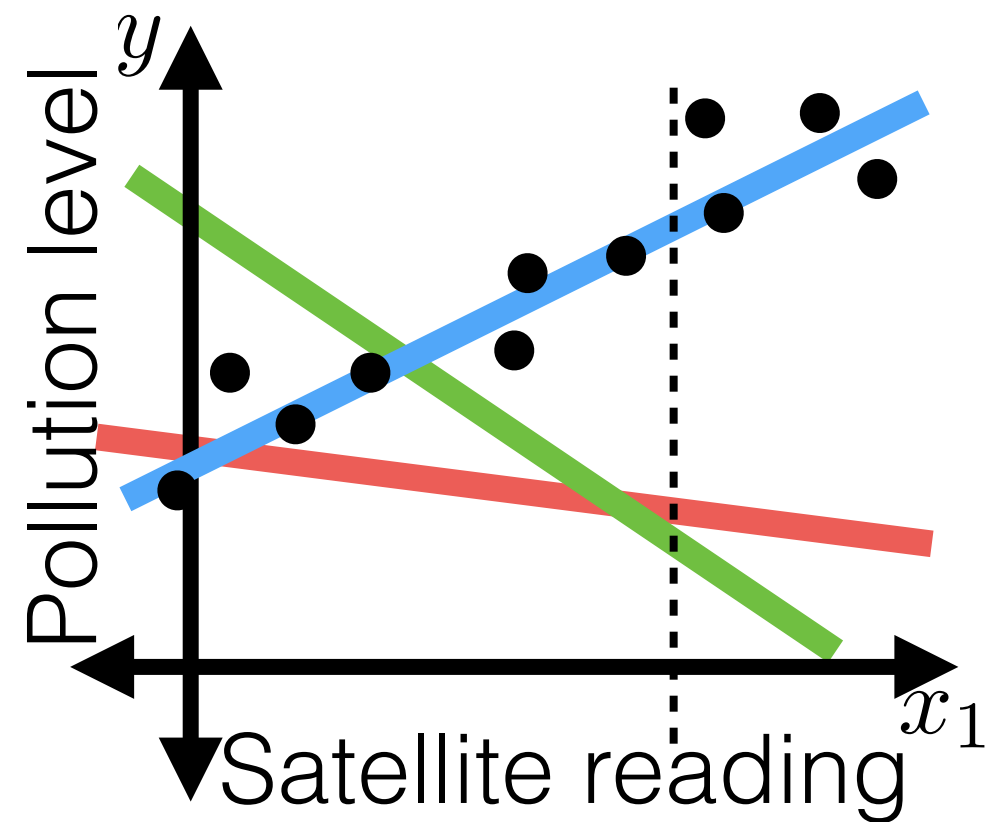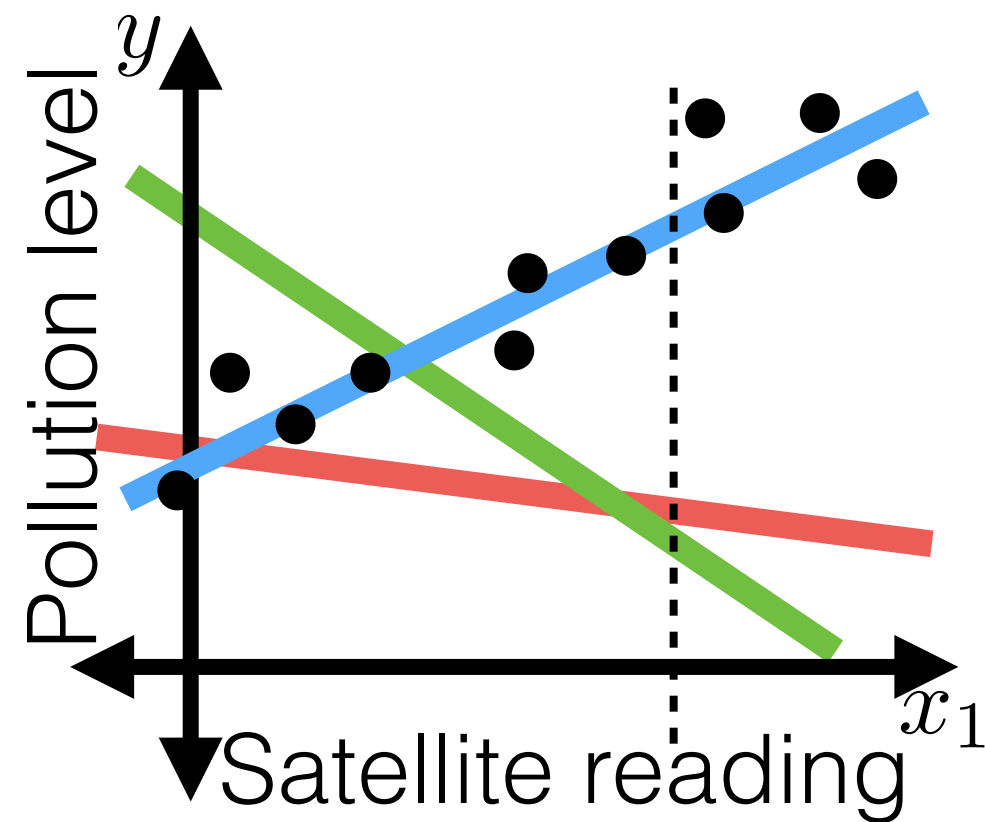


linear regression hypothesis

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
- Example: ridge regression

$$\frac{1}{n}\sum_{i=1}^{n} L(\theta^{\top} x^{(i)} + \theta_0, y^{(i)}) + \lambda R(\Theta) \quad (\lambda > 0)$$



linear regression hypothesis

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
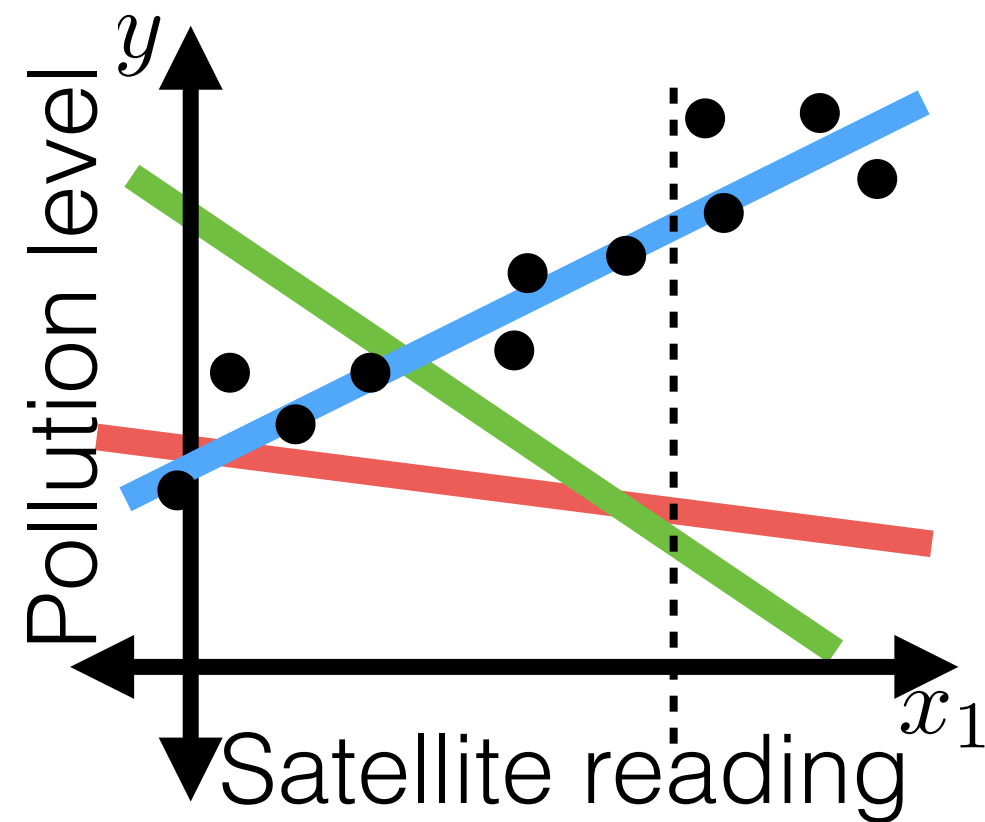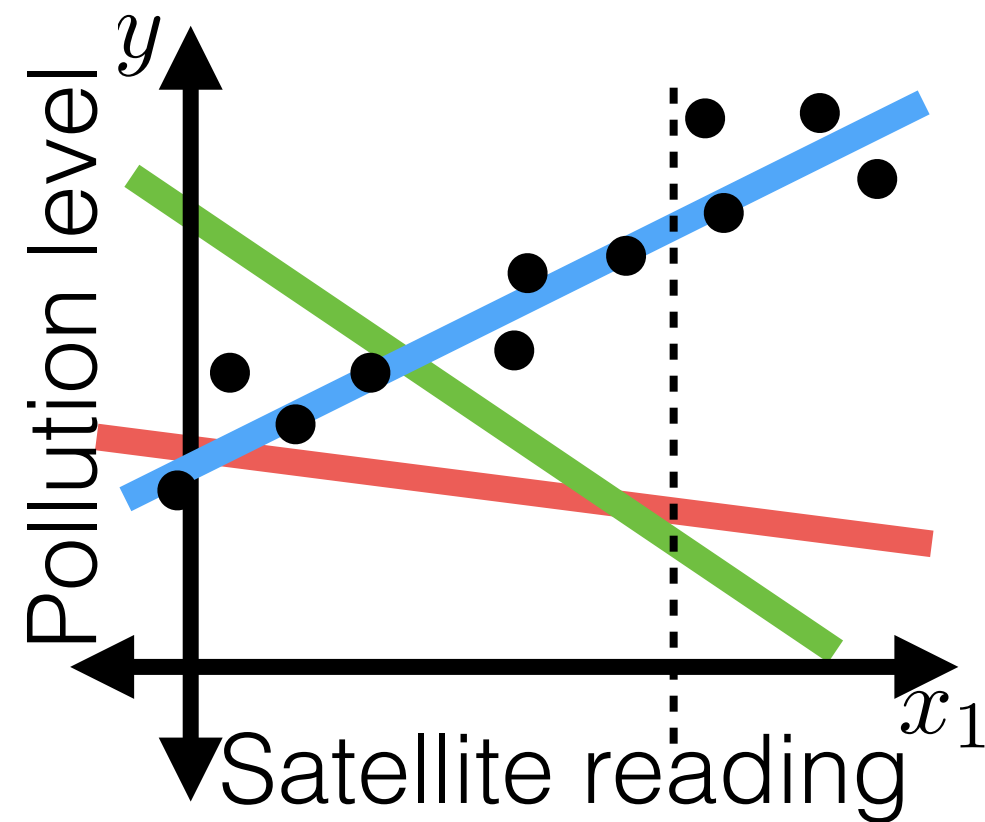- Example: ridge regression



linear regression hypothesis

$$\frac{1}{n}\sum_{i=1}^{n} L(\theta^\top x^{(i)} + \theta_0, y^{(i)}) + \lambda R(\Theta) \quad (\lambda > 0)$$

2

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
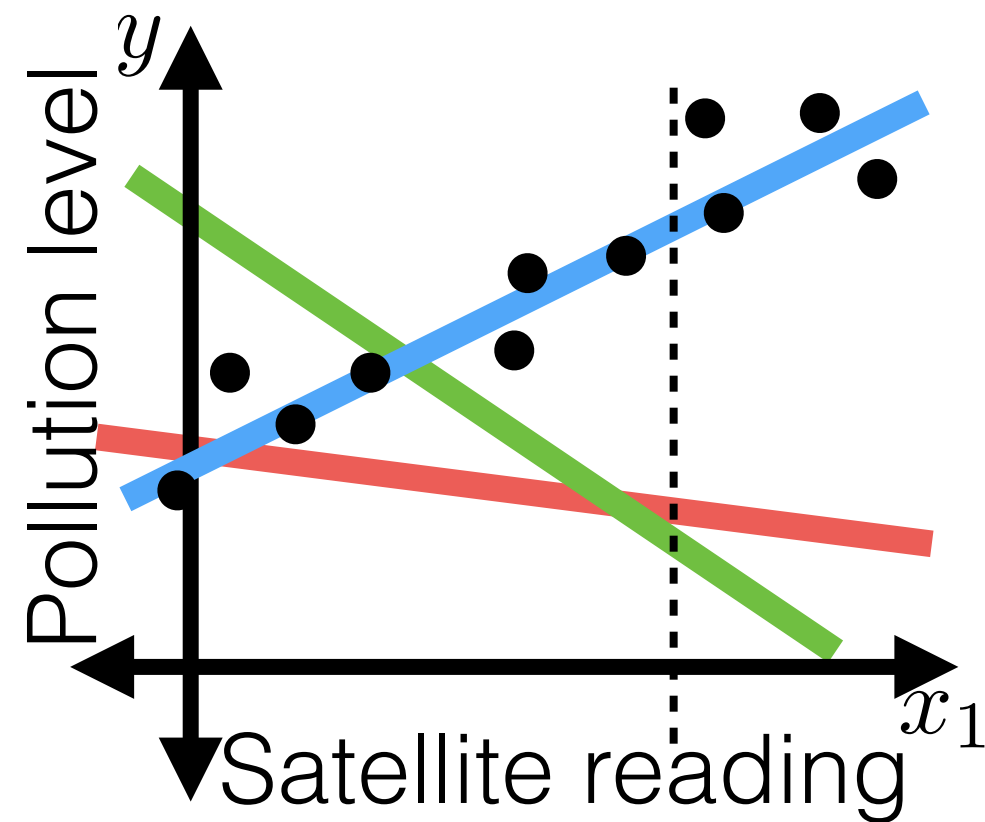- Example: ridge regression

linear regression hypothesis
squared loss $L(g,a) = (g-a)^2$

$$\frac{1}{n}\sum_{i=1}^{n} L(\theta^\top x^{(i)} + \theta_0, y^{(i)}) + \lambda R(\Theta) \quad (\lambda > 0)$$


Pollution level vs Satellite reading ($x_1$, $y$)

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
- Example: ridge regression



linear regression hypothesis
squared loss $L(g,a) = (g\text{-}a)^2$

$$\frac{1}{n} \sum_{i=1}^{n} L(\theta^\top x^{(i)} + \theta_0, y^{(i)}) + \lambda R(\Theta) \quad (\lambda > 0)$$

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
- Example: ridge regression



linear regression hypothesis
squared loss $L(g,a) = (g-a)^2$

$$\frac{1}{n}\sum_{i=1}^{n}(\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda R(\Theta)\ (\lambda > 0)$$

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
- Example: ridge regression

$$\frac{1}{n}\sum_{i=1}^{n}(\theta^{\top}x^{(i)}+\theta_0-y^{(i)})^2+\lambda R(\Theta)\ (\lambda>0)$$
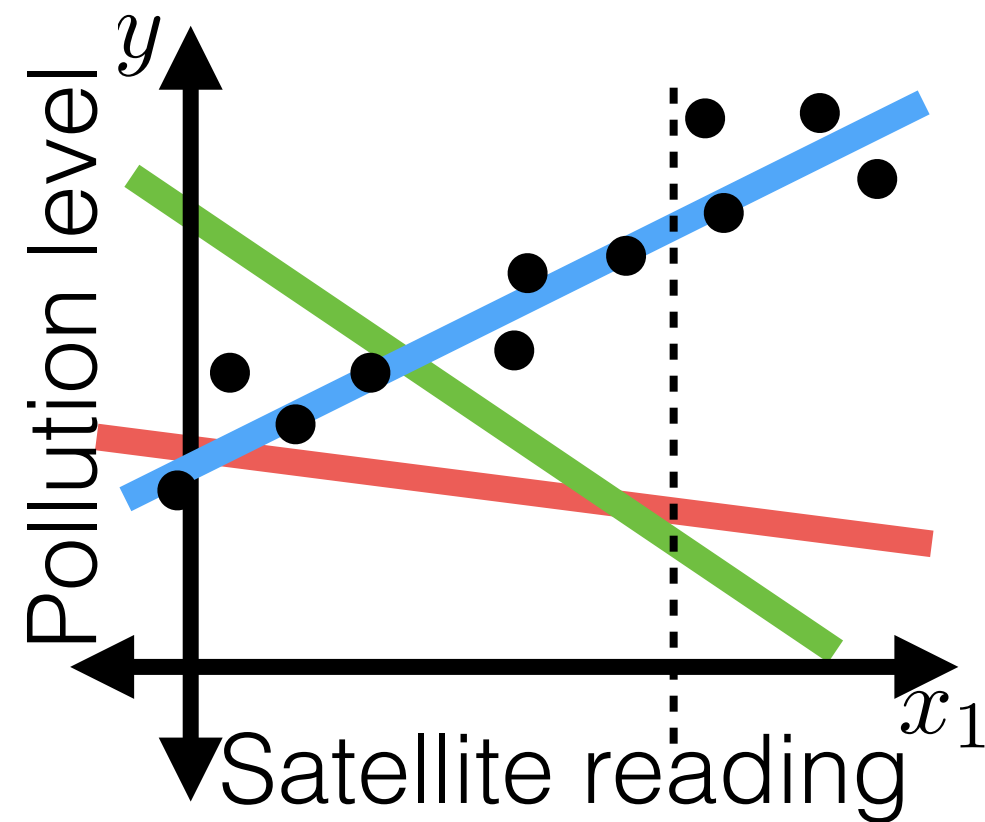
linear regression hypothesis
squared loss $L(g,a) = (g\text{-}a)^2$

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
- Example: ridge regression

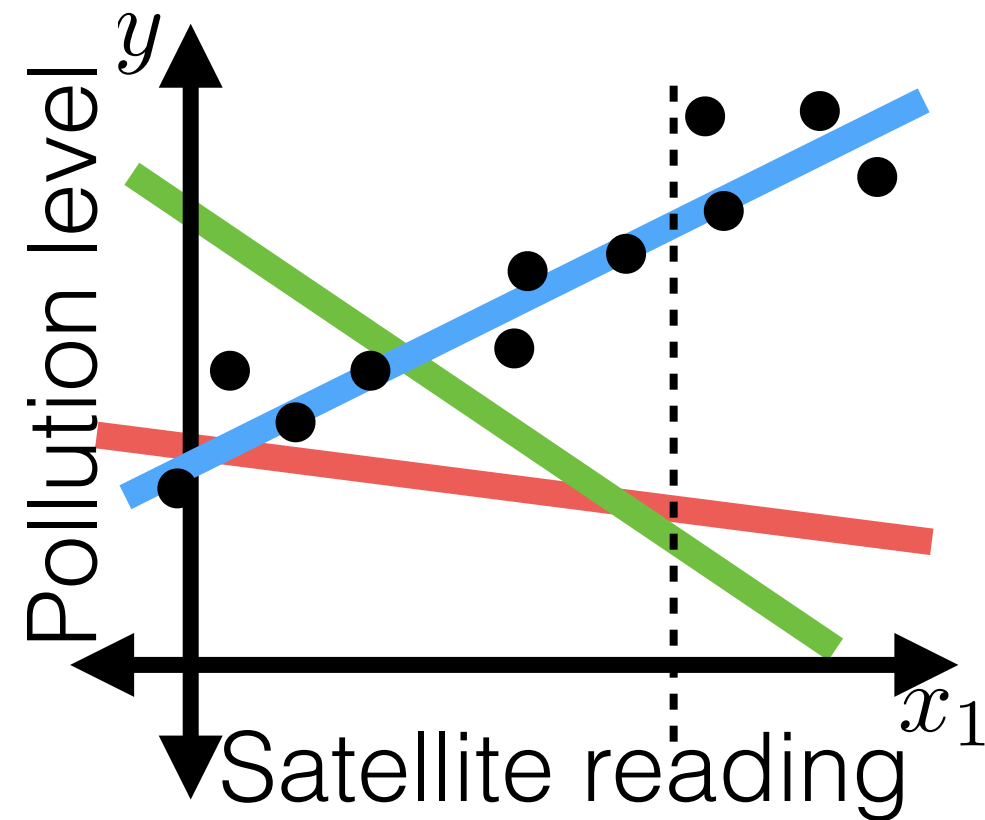$$\frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda R(\Theta) \; (\lambda > 0)$$



linear regression hypothesis
squared loss $L(g,a) = (g\text{-}a)^2$
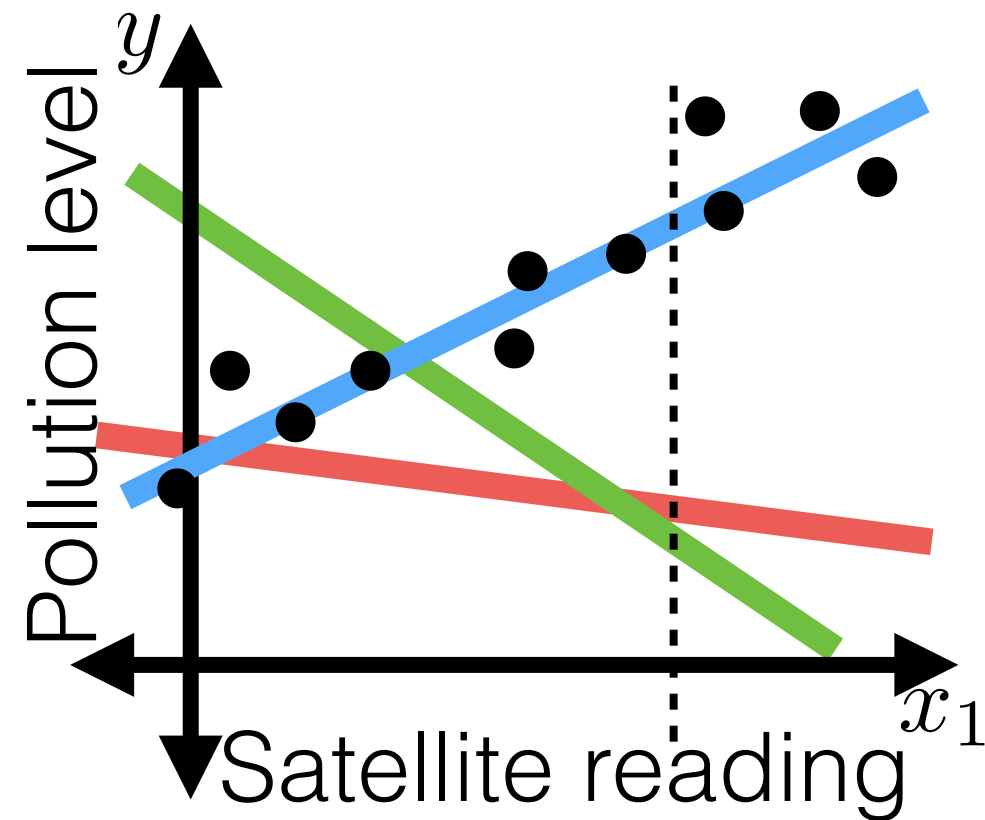squared-norm as regularizer

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
- Example: ridge regression

$$\frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda R(\Theta) \ (\lambda > 0)$$



linear regression hypothesis
squared loss $L(g,a) = (g-a)^2$
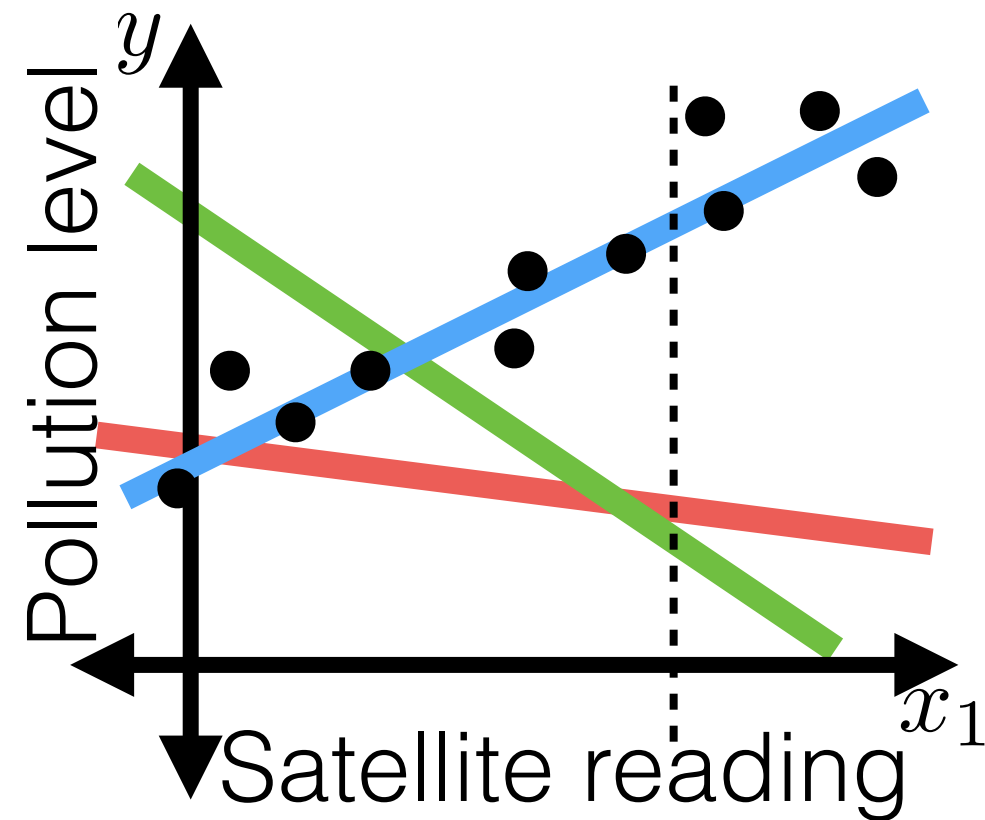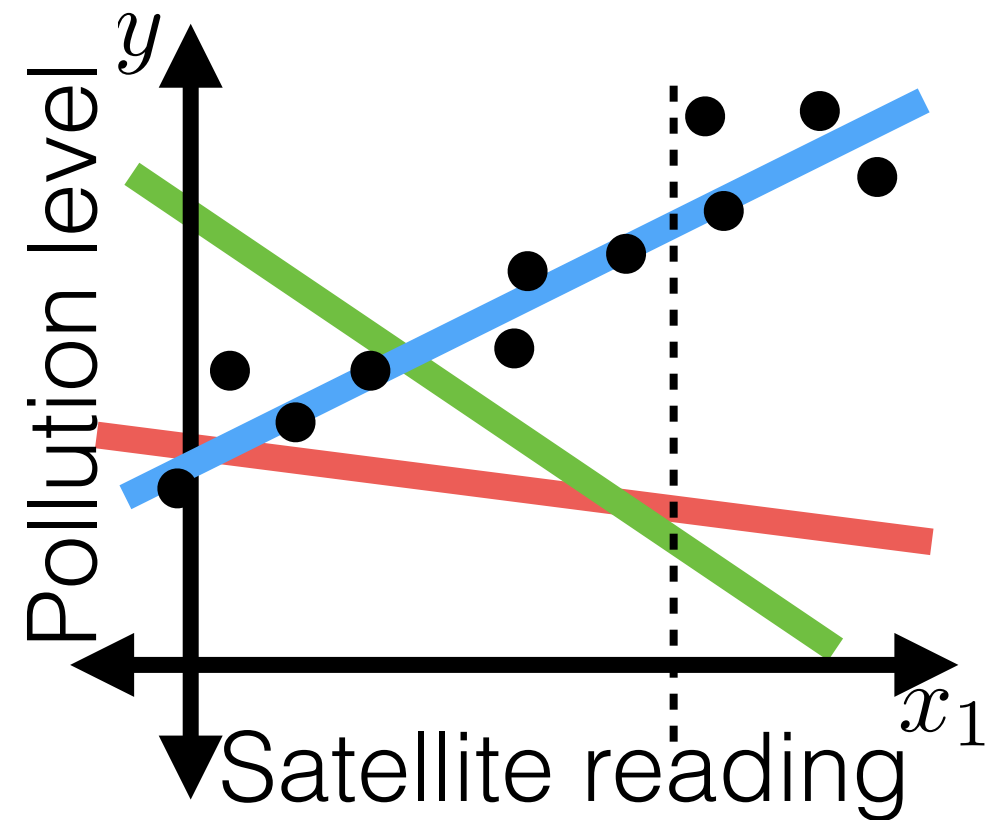
squared-norm as regularizer

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
- Example: ridge regression

$$\frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$



Pollution level
Satellite reading

linear regression hypothesis
squared loss $L(g,a) = (g\text{-}a)^2$
squared-norm as regularizer

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
- Example: ridge regression

$$\frac{1}{n}\sum_{i=1}^{n}(\theta^{\top}x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda\|\theta\|^2 \quad (\lambda > 0)$$



linear regression hypothesis
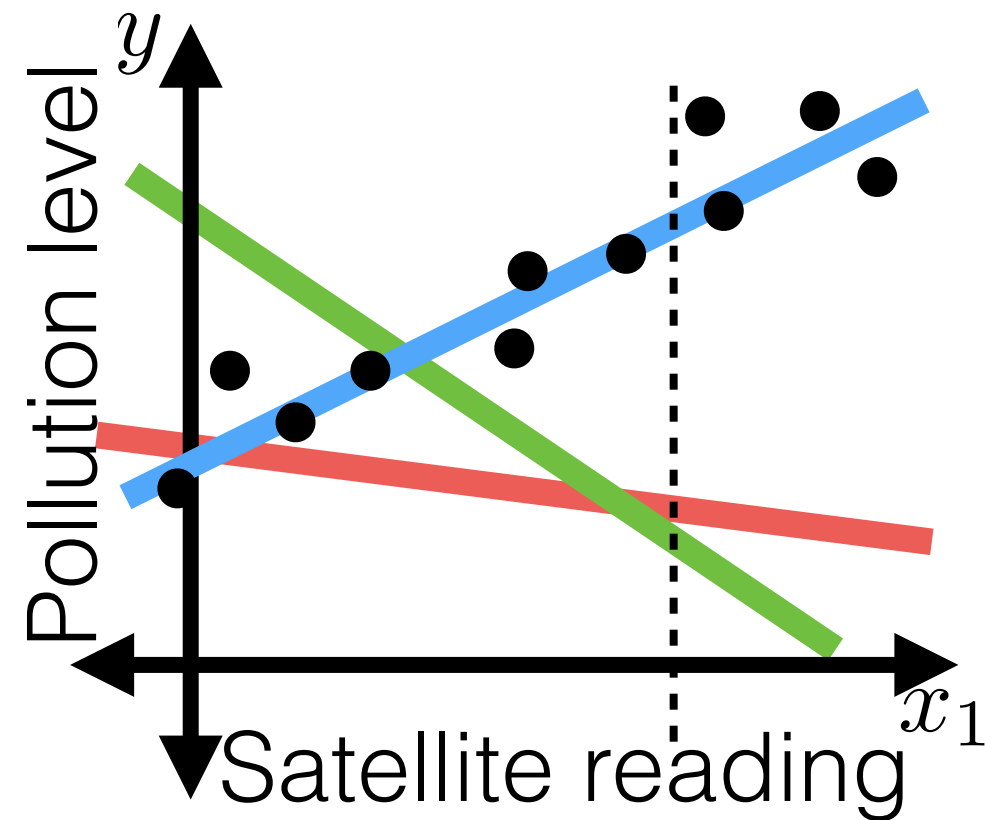squared loss $L(g,a) = (g-a)^2$
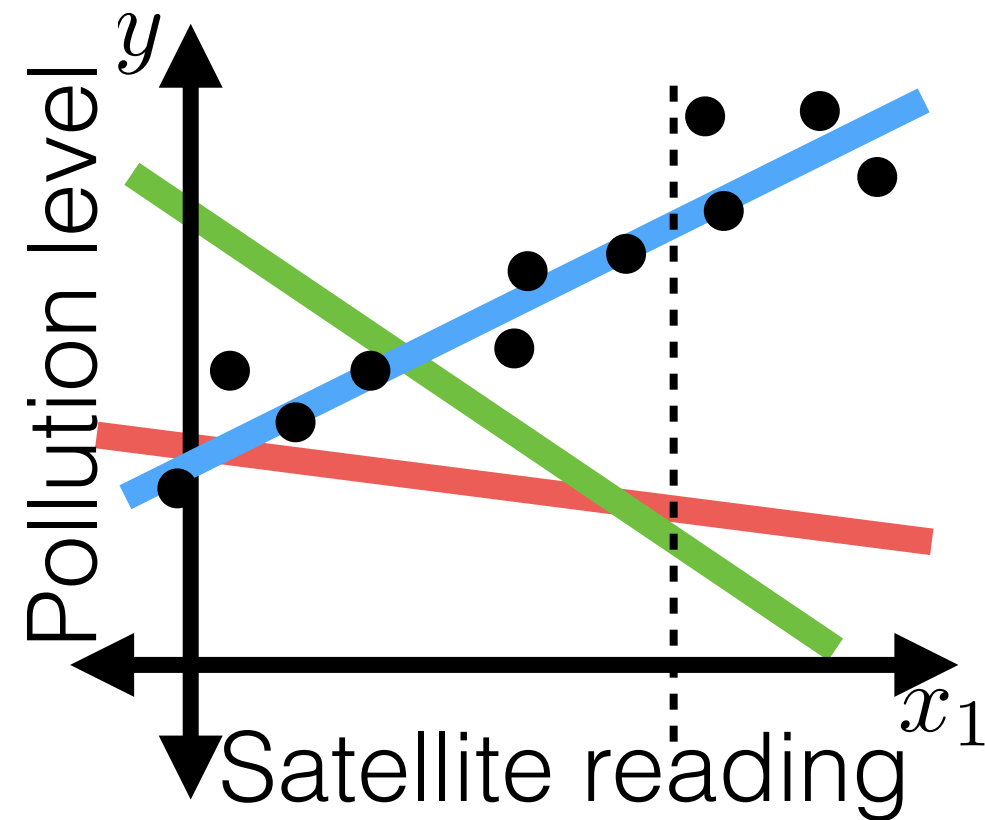squared-norm as regularizer

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
- Example: ridge regression

$$J_{\mathrm{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$



linear regression hypothesis
squared loss $L(g,a) = (g\text{-}a)^2$
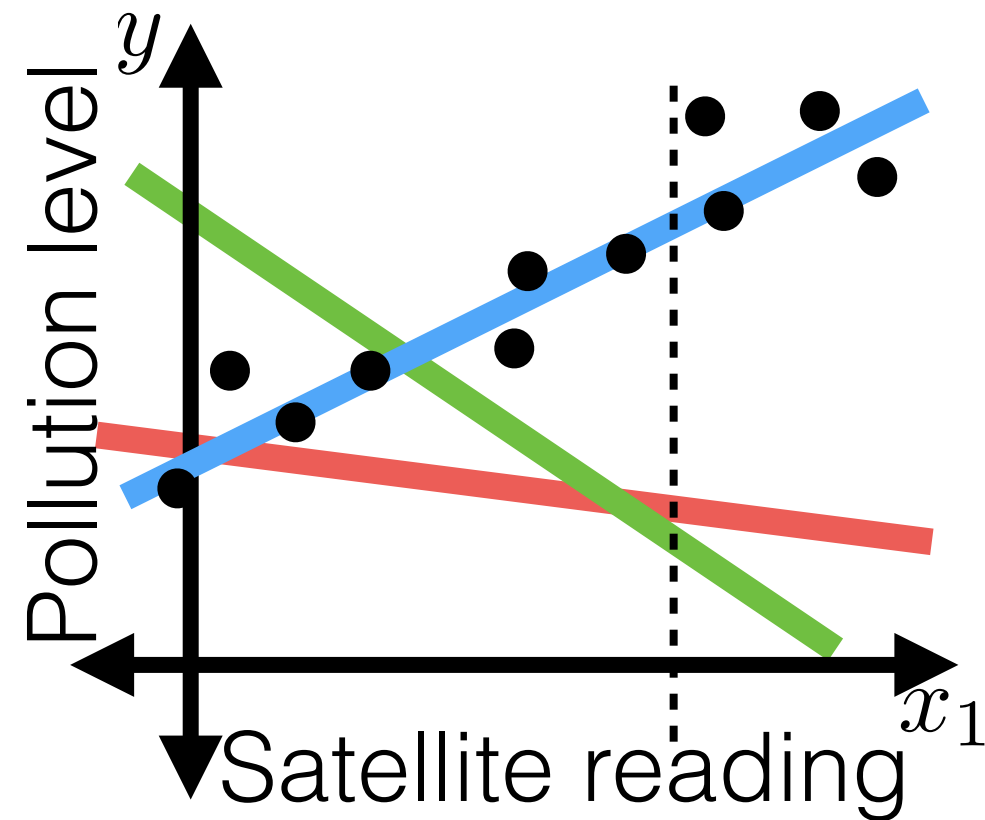squared-norm as regularizer

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer
- Example: ridge regression

$$J_{\mathrm{ridge}}(\theta, \theta_0) = \frac{1}{n}\sum_{i=1}^{n}(\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda\|\theta\|^2 \quad (\lambda > 0)$$



linear regression hypothesis
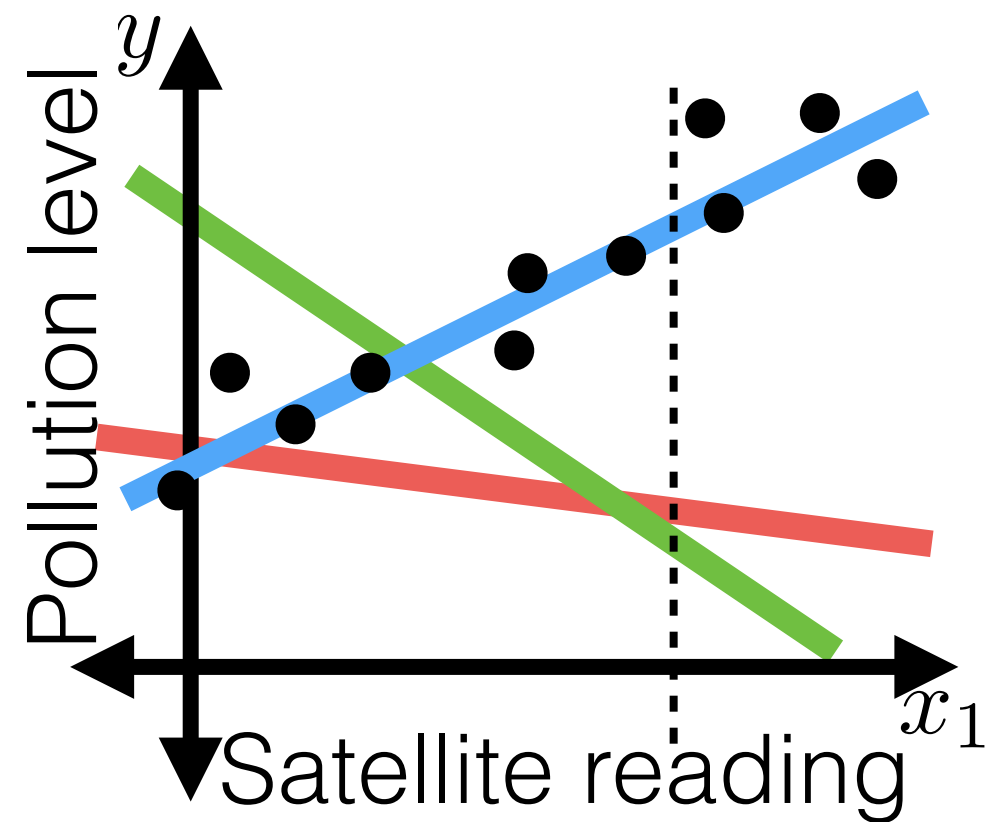squared loss $L(g,a) = (g\text{-}a)^2$
squared-norm as regularizer

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer



linear regression hypothesis
squared loss $L(g,a) = (g-a)^2$

squared-norm as regularizer

- Example: ridge regression

$$J_{\mathrm{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

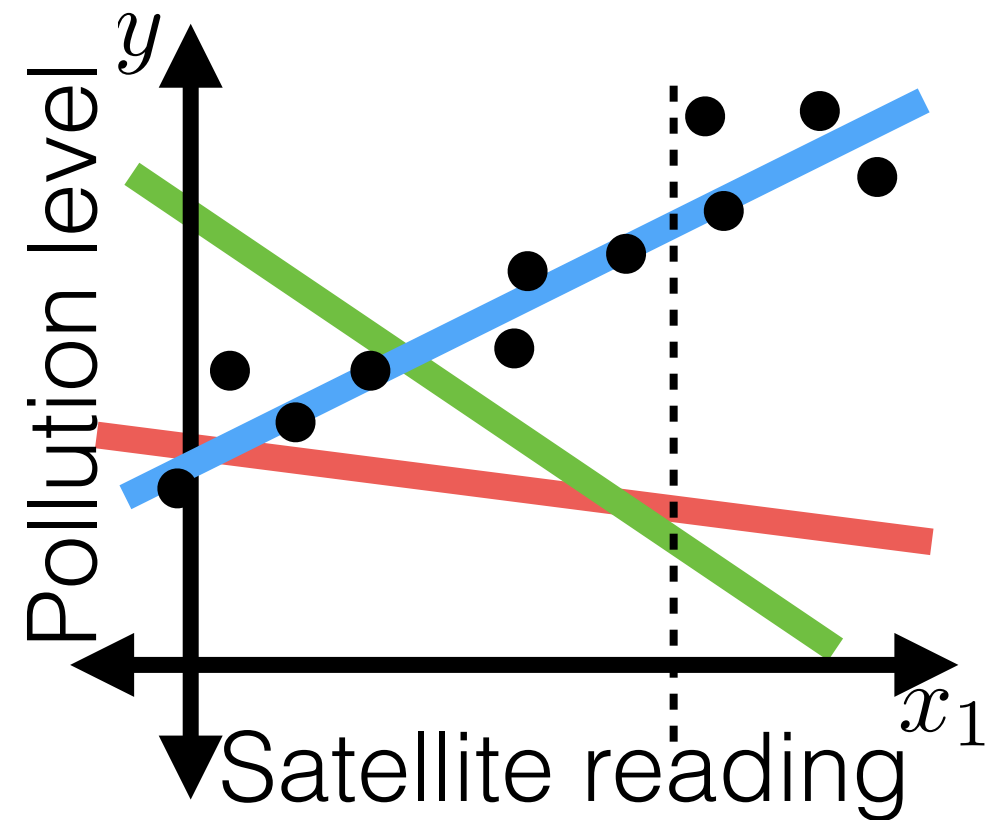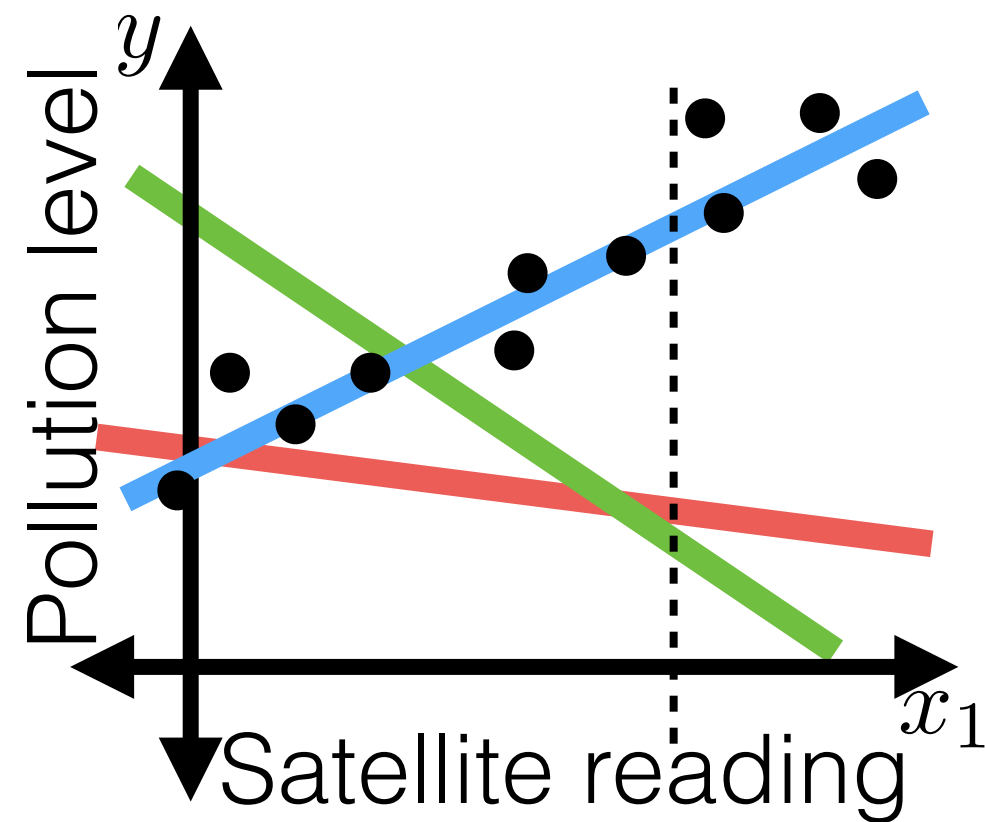- "All models are wrong, but some are useful" -George Box

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer



linear regression hypothesis
squared loss $L(g,a) = (g\text{-}a)^2$
squared-norm as regularizer

- Example: ridge regression

$$J_{\mathrm{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

- "All models are wrong, but some are useful" -George Box
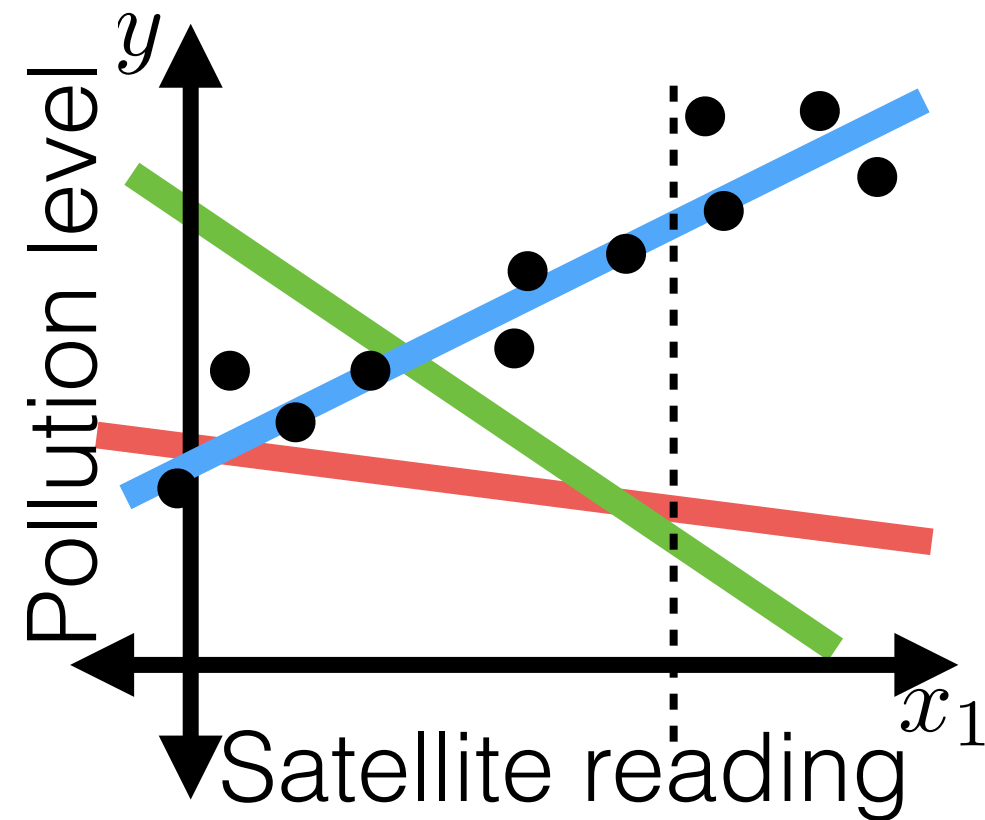- Limitations of a closed-form solution for objective minimizer

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer



linear regression hypothesis
squared loss $L(g,a) = (g-a)^2$

- Example: ridge regression

$$J_{\mathrm{ridge}}(\theta, \theta_0) = \frac{1}{n}\sum_{i=1}^{n}(\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda\|\theta\|^2 \quad (\lambda > 0)$$

squared-norm as regularizer

- "All models are wrong, but some are useful" -George Box
- Limitations of a closed-form solution for objective minimizer
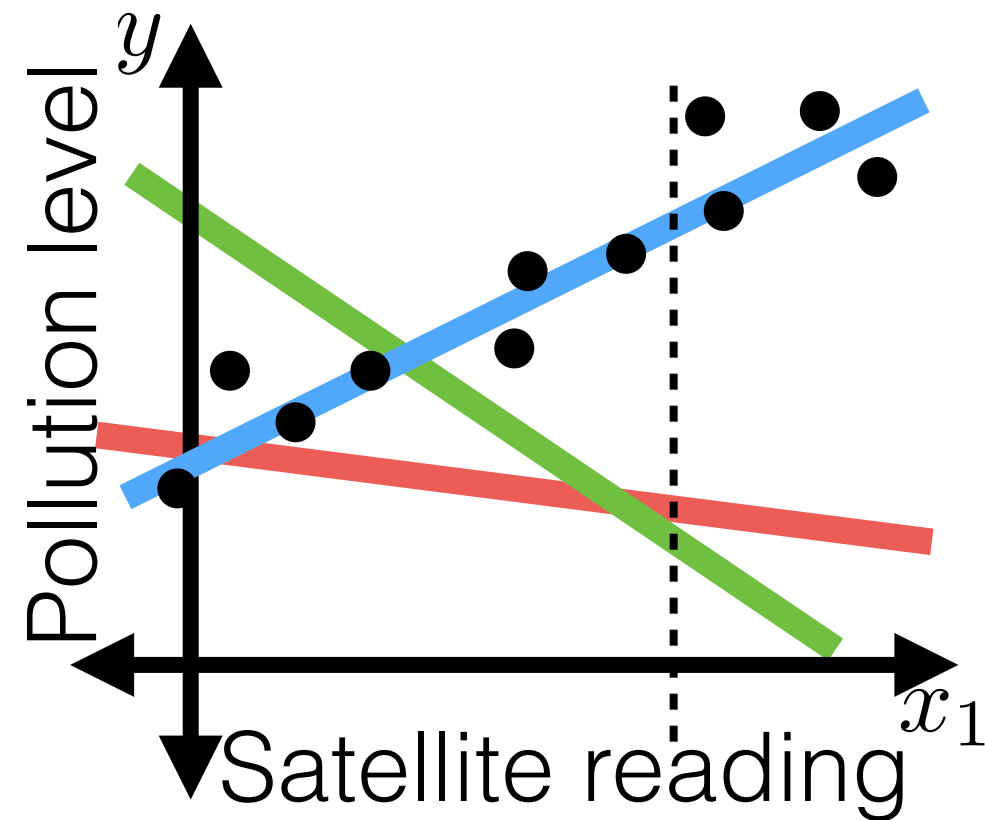  - Other hypotheses or loss or regularizer

2

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer



linear regression hypothesis
squared loss $L(g,a) = (g-a)^2$

squared-norm as regularizer

- Example: ridge regression

$$J_{\mathrm{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

- "All models are wrong, but some are useful" -George Box
- Limitations of a closed-form solution for objective minimizer
  - Other hypotheses or loss or regularizer: maybe no closed-form solution, or difficult

2

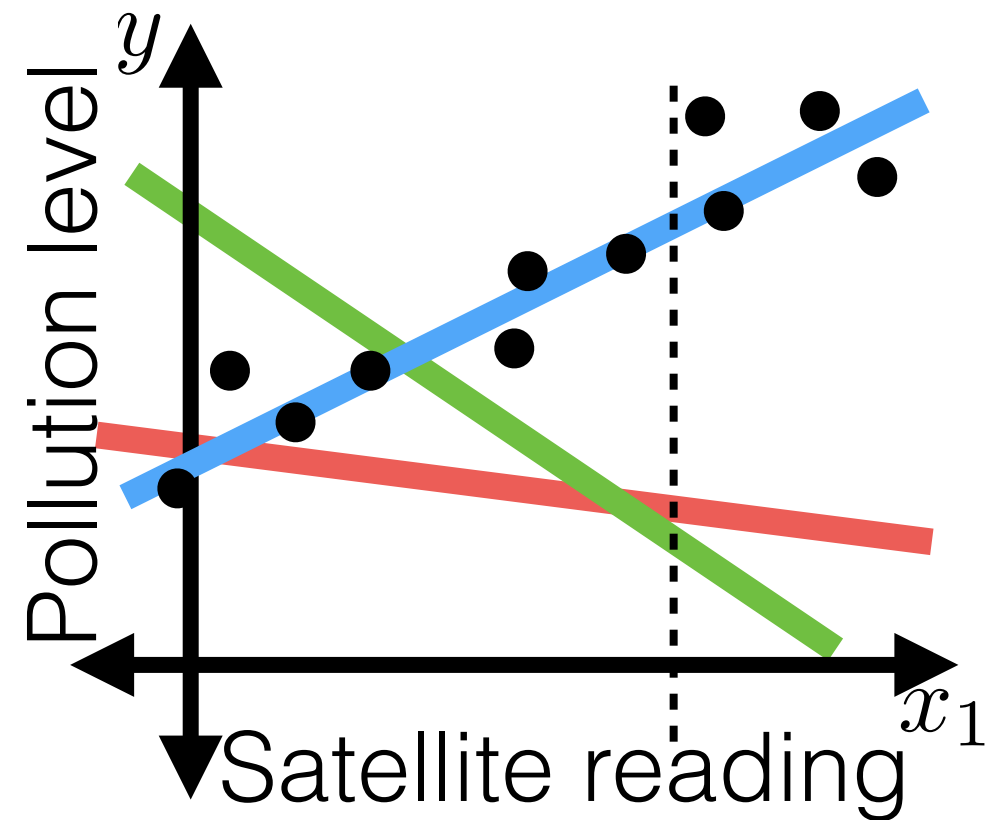# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer



linear regression hypothesis
squared loss $L(g,a) = (g-a)^2$

- Example: ridge regression

$$J_{\mathrm{ridge}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda \|\theta\|^2 \quad (\lambda > 0)$$

squared-norm as regularizer

- "All models are wrong, but some are useful" -George Box
- Limitations of a closed-form solution for objective minimizer
  - Other hypotheses or loss or regularizer: maybe no closed-form solution, or difficult

e.g.

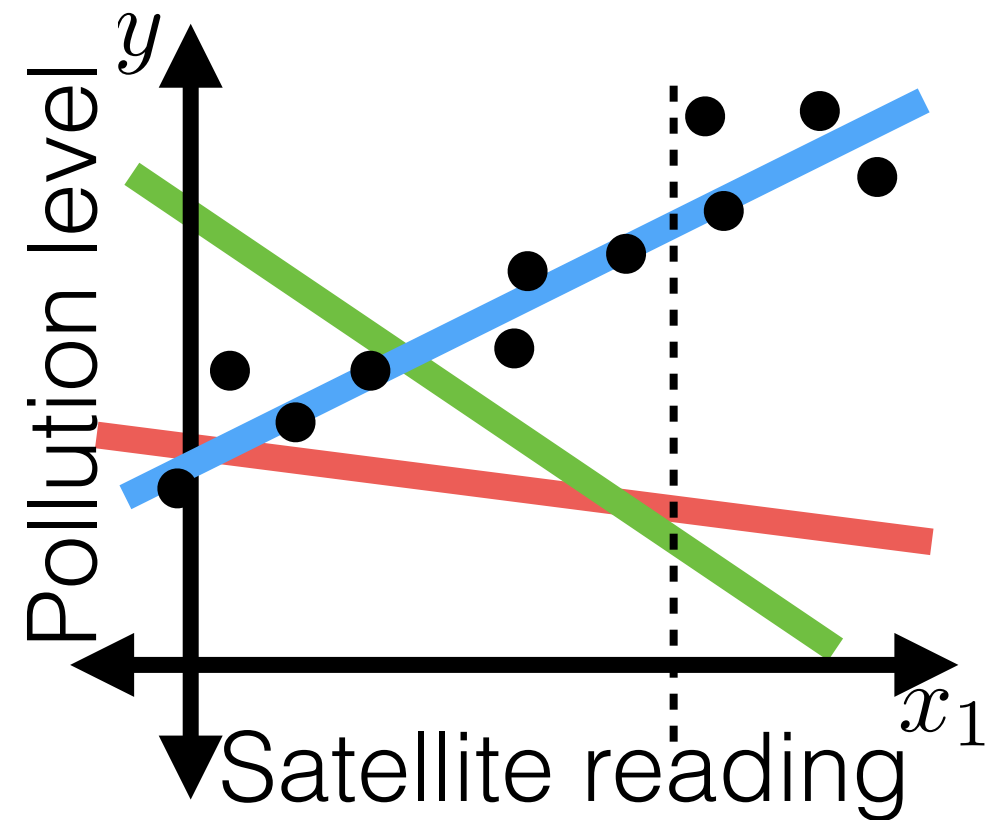# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer



linear regression hypothesis
squared loss $L(g,a) = (g-a)^2$

- Example: ridge regression

$$J_{\mathrm{ridge}}(\theta, \theta_0) = \frac{1}{n}\sum_{i=1}^{n}(\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda\|\theta\|^2 \quad (\lambda > 0)$$

squared-norm as regularizer

- "All models are wrong, but some are useful" -George Box

- Limitations of a closed-form solution for objective minimizer
  - Other hypotheses or loss or regularizer: maybe no closed-form solution, or difficult

e.g. $L(g,a) =$

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer



linear regression hypothesis
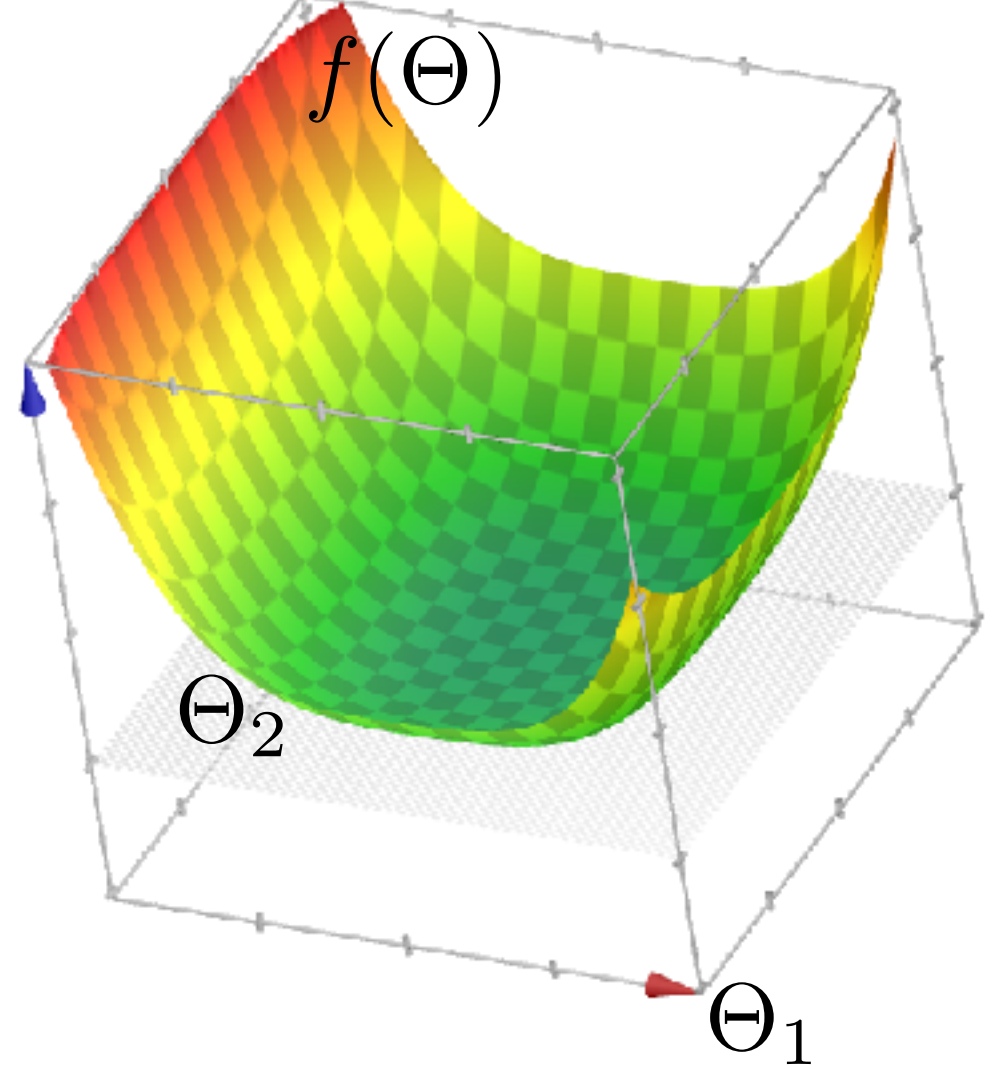squared loss $L(g,a) = (g-a)^2$

- Example: ridge regression

$$J_{\mathrm{ridge}}(\theta, \theta_0) = \frac{1}{n}\sum_{i=1}^{n}(\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda\|\theta\|^2 \quad (\lambda > 0)$$
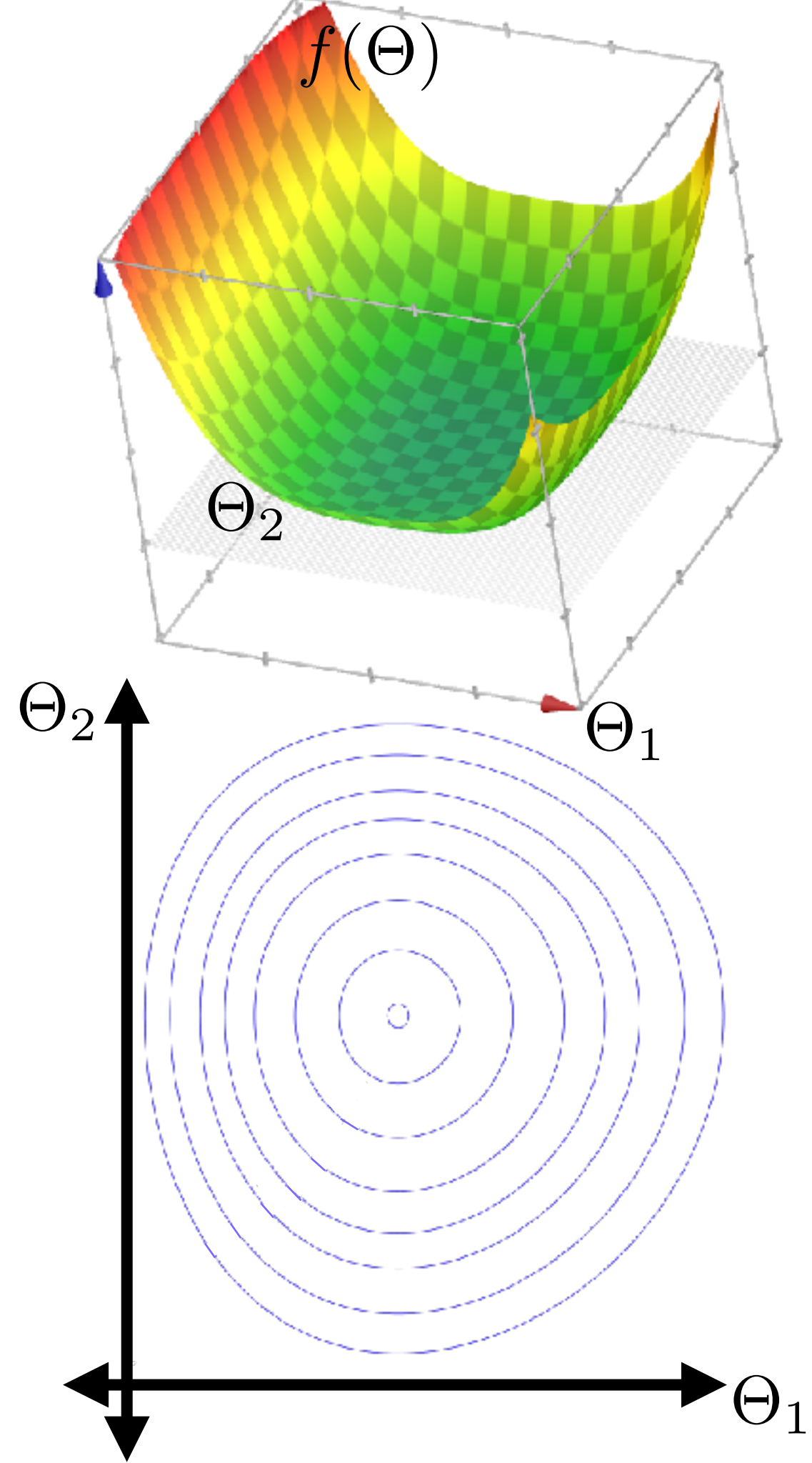
squared-norm as regularizer

- "All models are wrong, but some are useful" -George Box
- Limitations of a closed-form solution for objective minimizer
  - Other hypotheses or loss or regularizer: maybe no closed-form solution, or difficult

e.g. $L(g,a) =$

$$\begin{cases} (g-a)^2 \text{ if } g > a \\ 5(g-a)^2 \text{ if } g \leq a \end{cases}$$

2

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer



linear regression hypothesis
squared loss $L(g,a) = (g\text{-}a)^2$

- Example: ridge regression

$$J_{\mathrm{ridge}}(\theta, \theta_0) = \frac{1}{n}\sum_{i=1}^{n}(\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda\|\theta\|^2 \quad (\lambda > 0)$$

squared-norm as regularizer

- "All models are wrong, but some are useful" -George Box
- Limitations of a closed-form solution for objective minimizer
  - Other hypotheses or loss or regularizer: maybe no closed-form solution, or difficult

  e.g. $L(g,a) =$
  $$\begin{cases} (g-a)^2 \text{ if } g > a \\ 5(g-a)^2 \text{ if } g \leq a \end{cases}$$

- Can be too slow to run, even in ridge regression

2

# Recall

- A general ML approach:
  - Collect data
  - Choose hypothesis class
  - Choose "good" hypothesis by minimizing training loss + regularizer



linear regression hypothesis
squared loss $L(g,a) = (g-a)^2$
squared-norm as regularizer

- Example: ridge regression
  e.g.

$$f(\Theta) = J_{\mathrm{ridge}}(\theta, \theta_0) = \frac{1}{n}\sum_{i=1}^{n}(\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2 + \lambda\|\theta\|^2 \quad (\lambda > 0)$$

- "All models are wrong, but some are useful" -George Box

- Limitations of a closed-form solution for objective minimizer
  - Other hypotheses or loss or regularizer: maybe no closed-form solution, or difficult

  e.g. $L(g,a) =$
  $$\begin{cases} (g-a)^2 \text{ if } g > a \\ 5(g-a)^2 \text{ if } g \le a \end{cases}$$

2 - Can be too slow to run, even in ridge regression

# Gradient descent

# Gradient descent



$f(\Theta)$

$\Theta_2$

$\Theta_1$

# Gradient descent



$f(\Theta)$

$\Theta_2$

$\Theta_1$

$\Theta_2$

$\Theta_1$

3

# Gradient descent

$f(\Theta)$

$\Theta_2$

$\Theta_1$

$\Theta_2$

$\Theta_1$

3

# Gradient descent



$f(\Theta)$

$\Theta_2$

$\Theta_1$

$\Theta_2$

$\Theta_1$

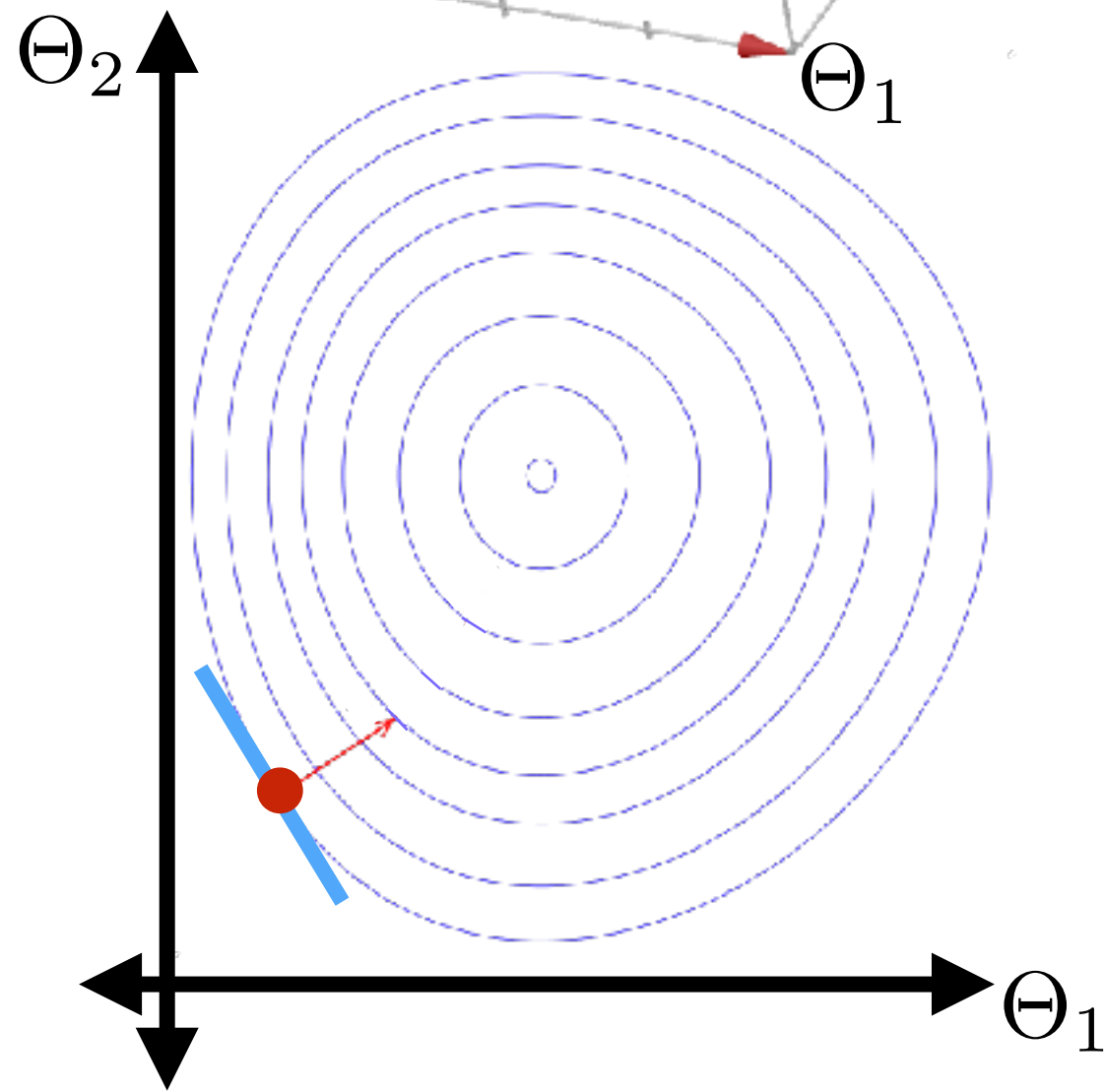# Gradient descent

$f(\Theta)$

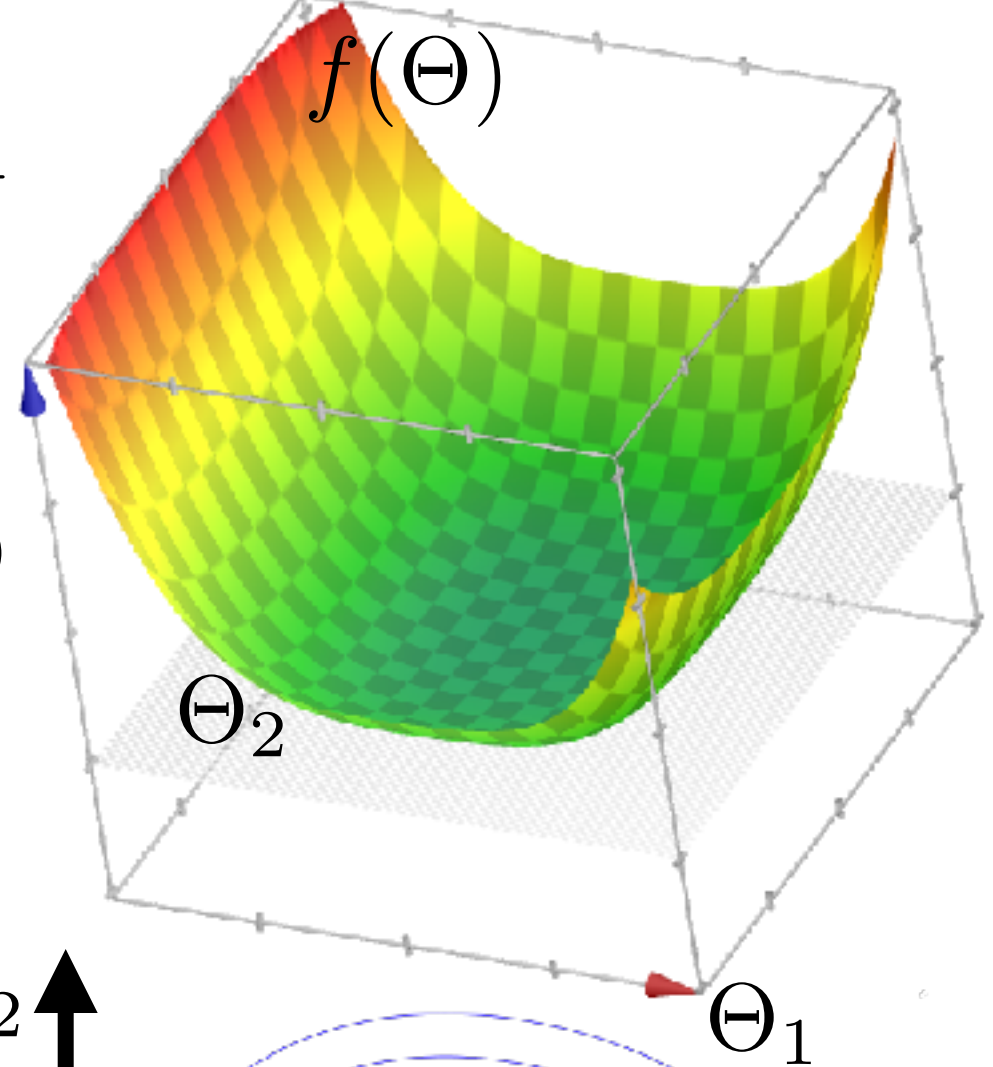$\Theta_2$

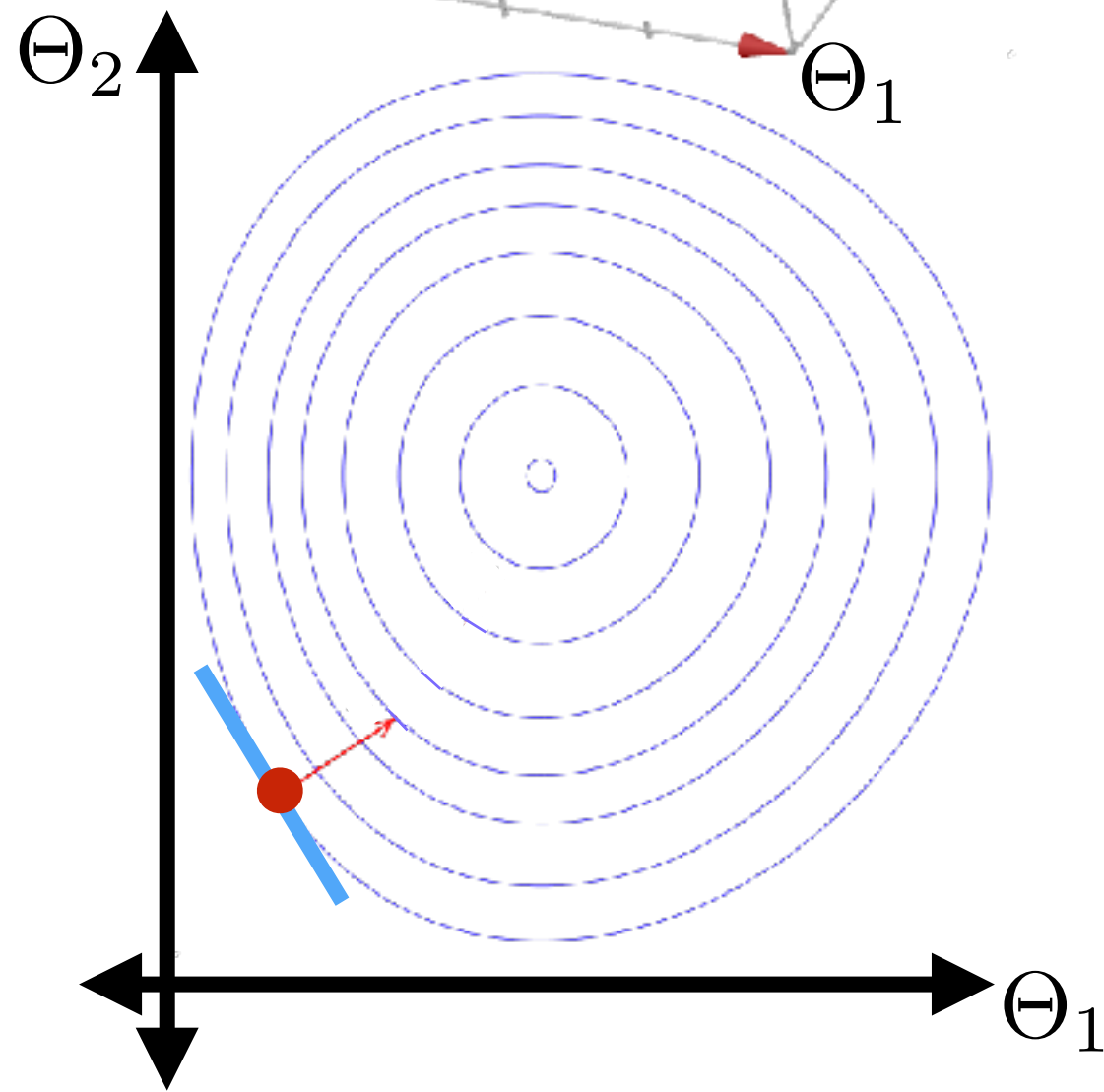$\Theta_1$

$\Theta_2$

$\Theta_1$

3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
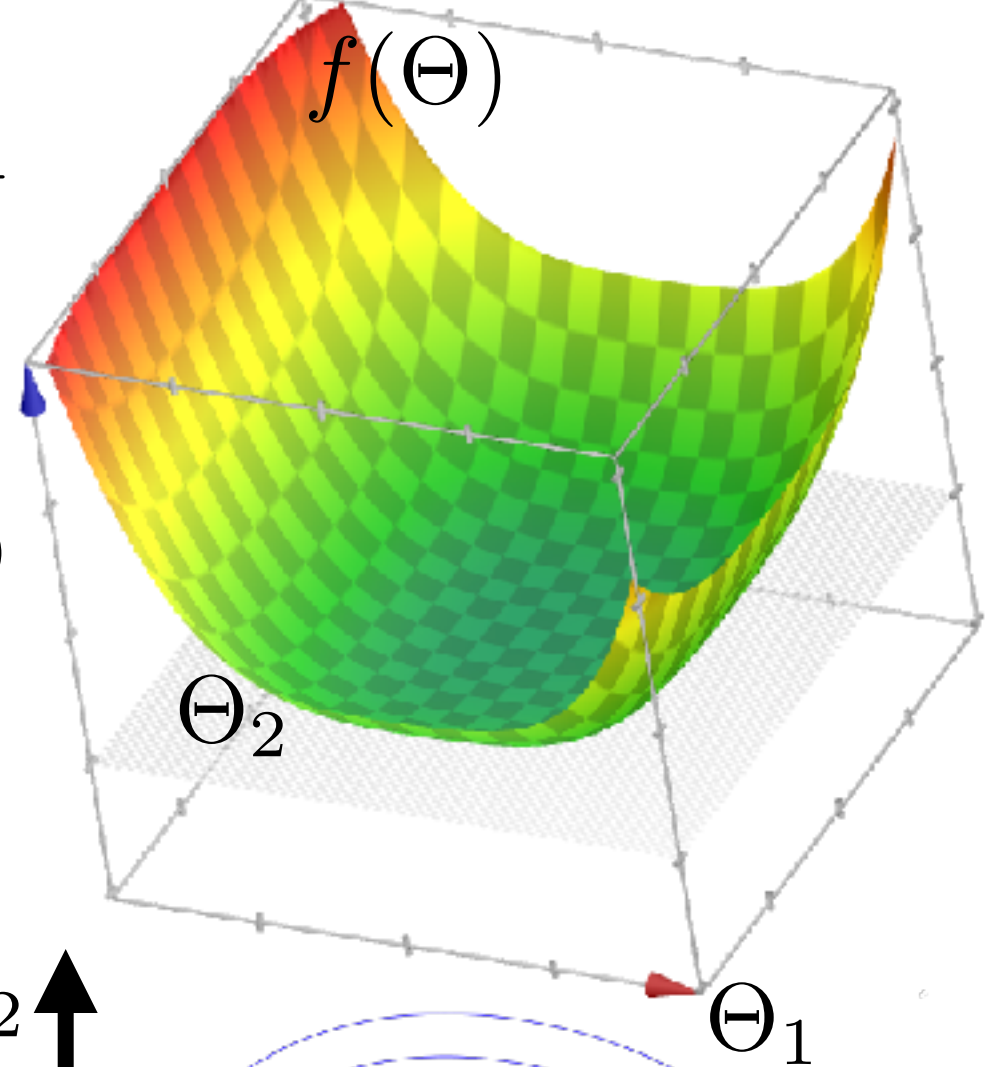  - with $\Theta \in \mathbb{R}^m$



3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

# Gradient descent

- Gradient $\nabla_{\Theta} f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^{\top}$
  - with $\Theta \in \mathbb{R}^m$

$f(\Theta)$

$\Theta_2$

$\Theta_1$
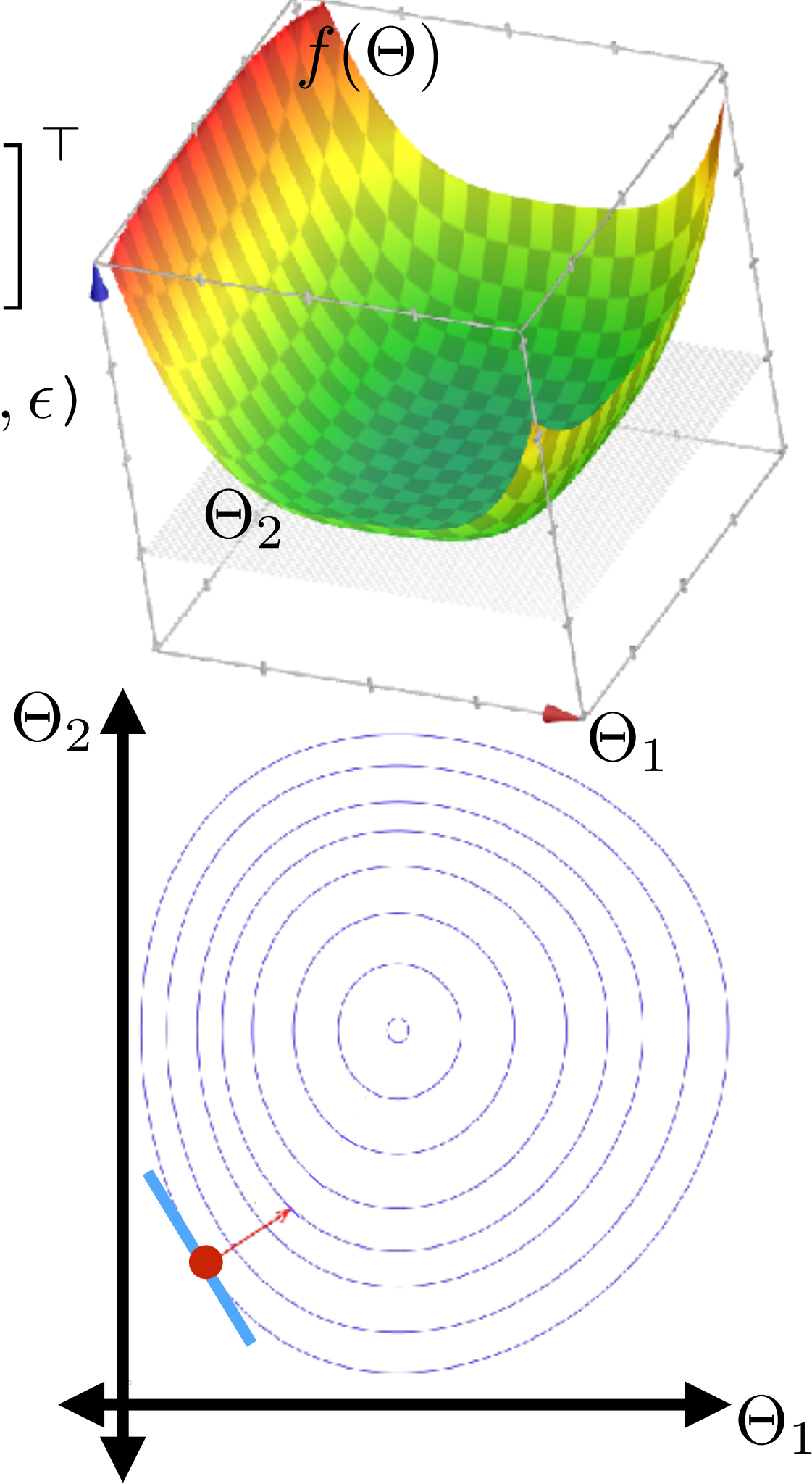
$\Theta_2$

$\Theta_1$
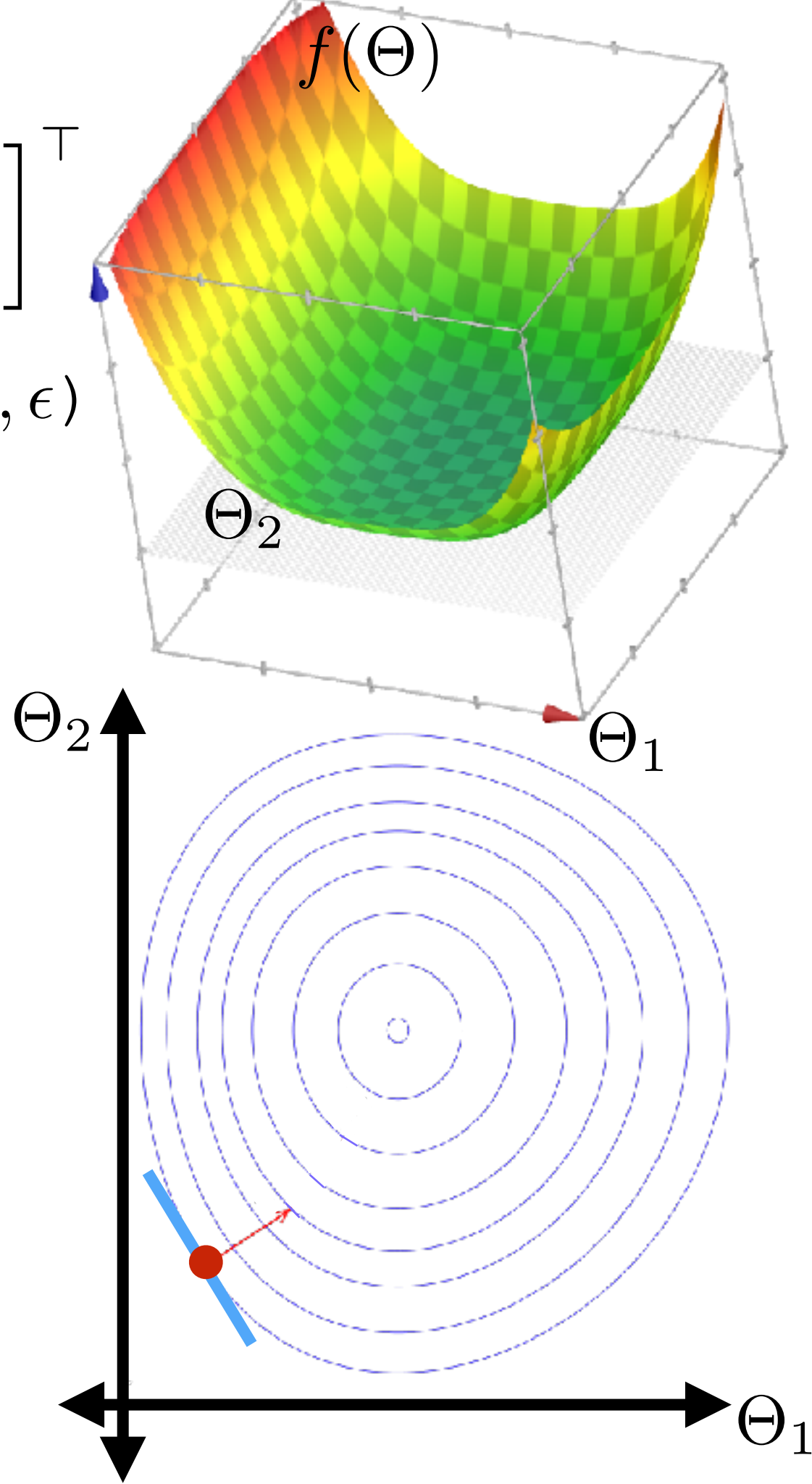
3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

# Gradient descent

- Gradient $\nabla_\Theta f = \left[\dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m}\right]^\top$
  - with $\Theta \in \mathbb{R}^m$
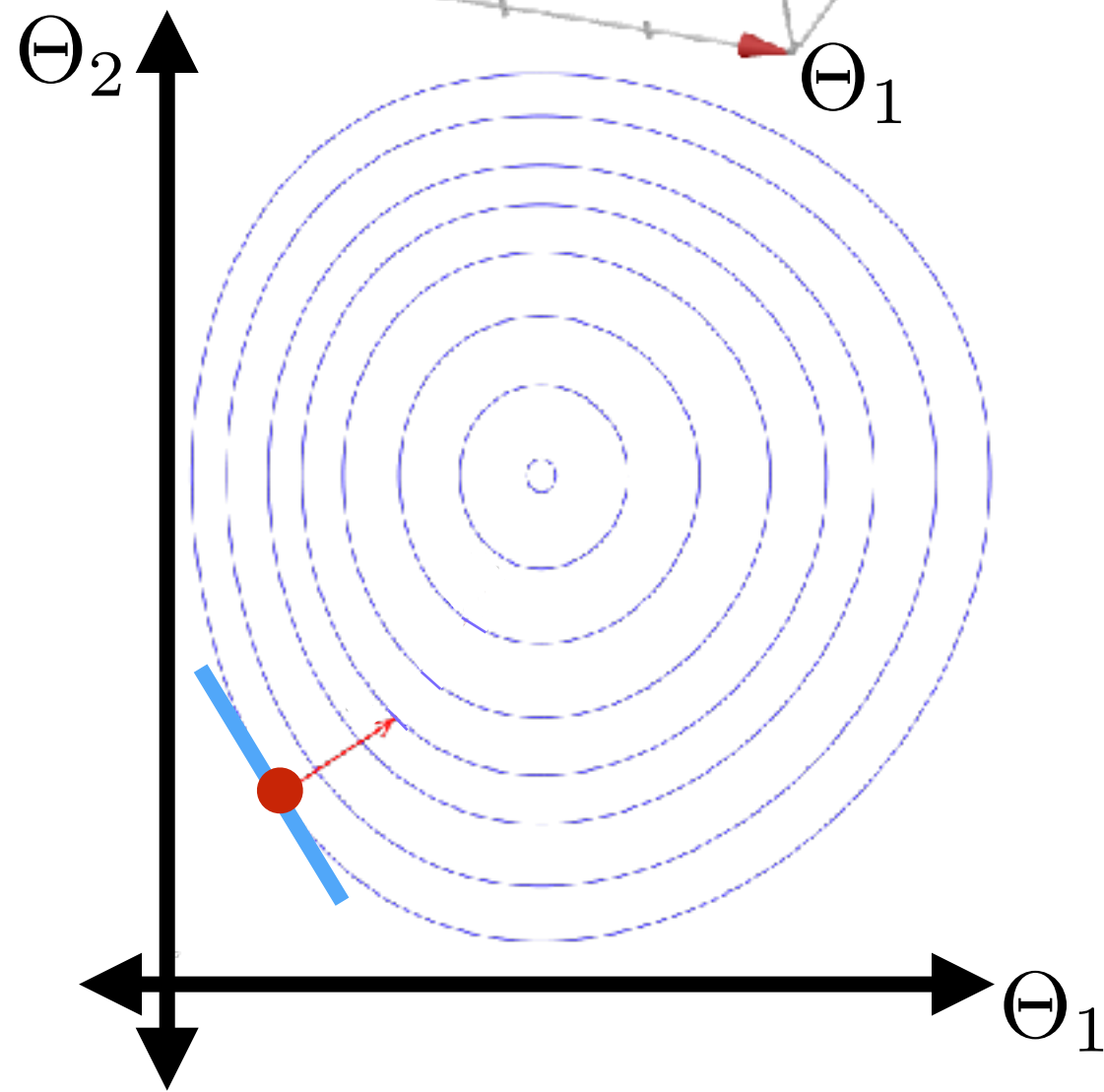
Gradient-Descent $(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$



3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
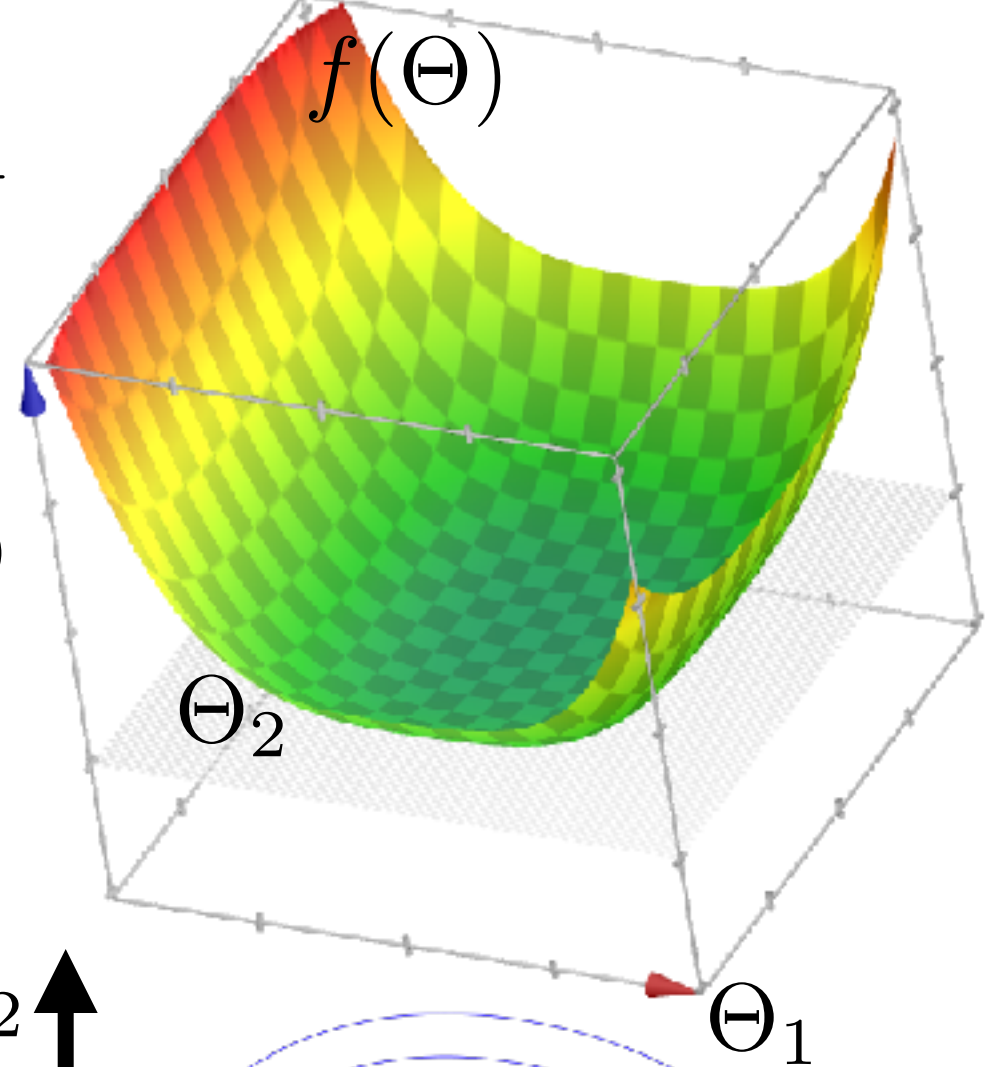  - with $\Theta \in \mathbb{R}^m$

Gradient-Descent $(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

Initialize $\Theta^{(0)} = \Theta_{\text{init}}$



3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

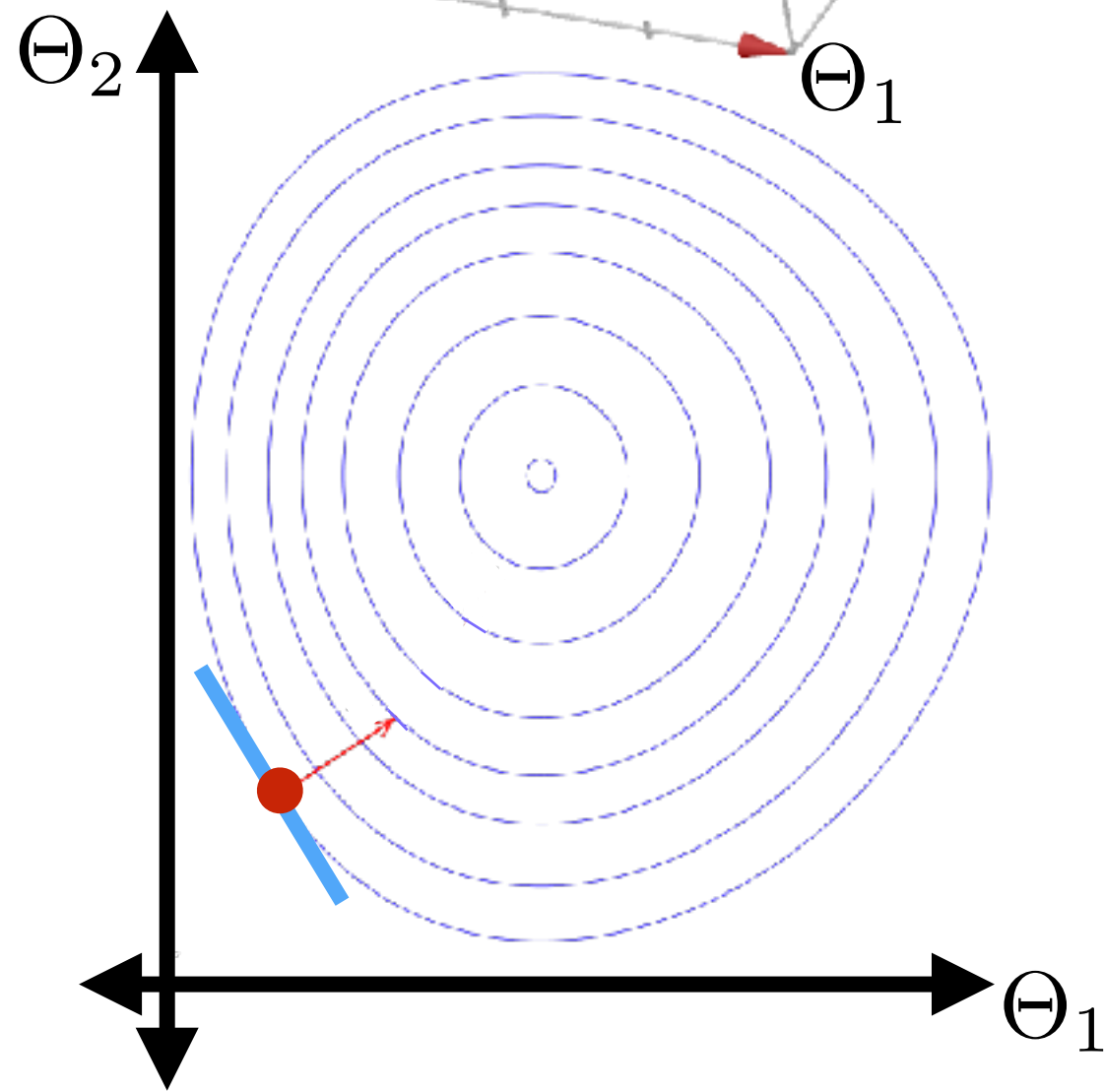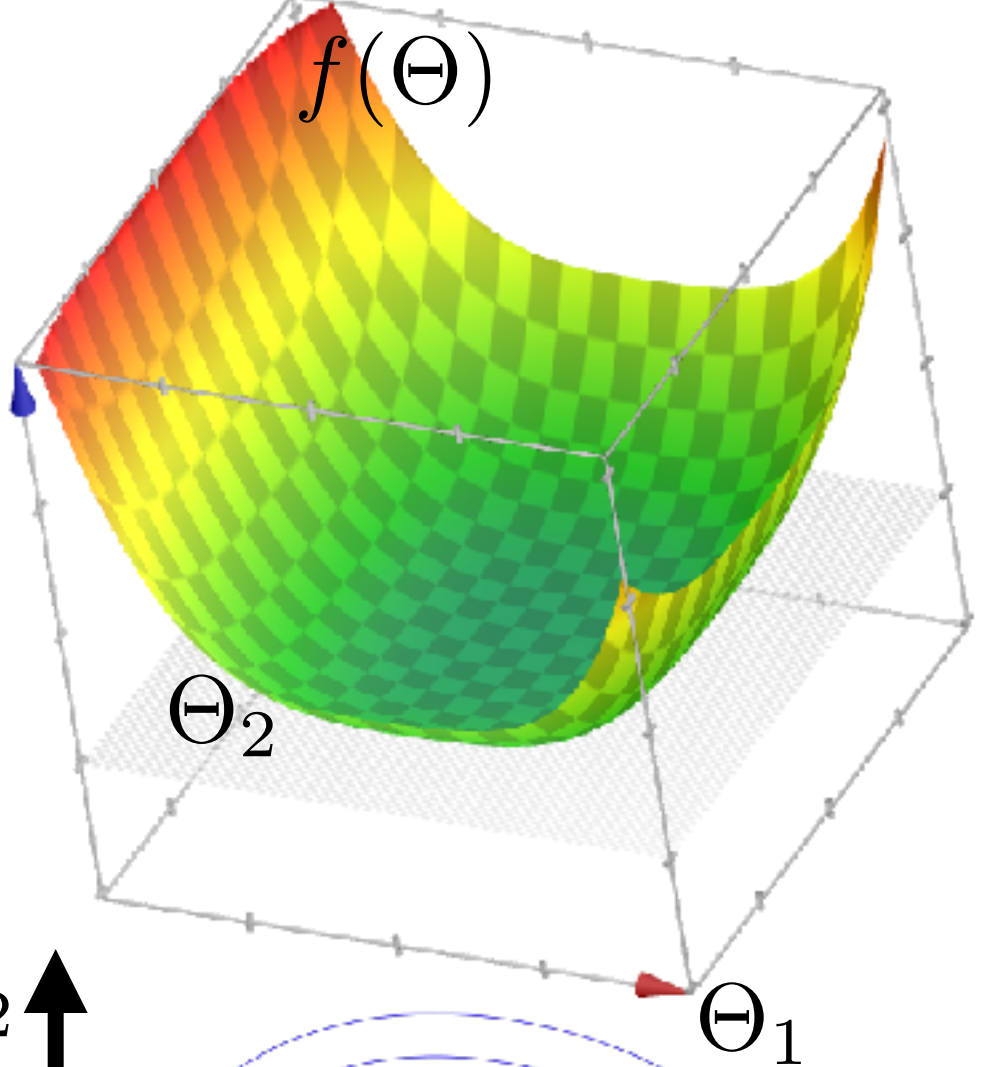Gradient-Descent $(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

Initialize $\Theta^{(0)} = \Theta_{\text{init}}$

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

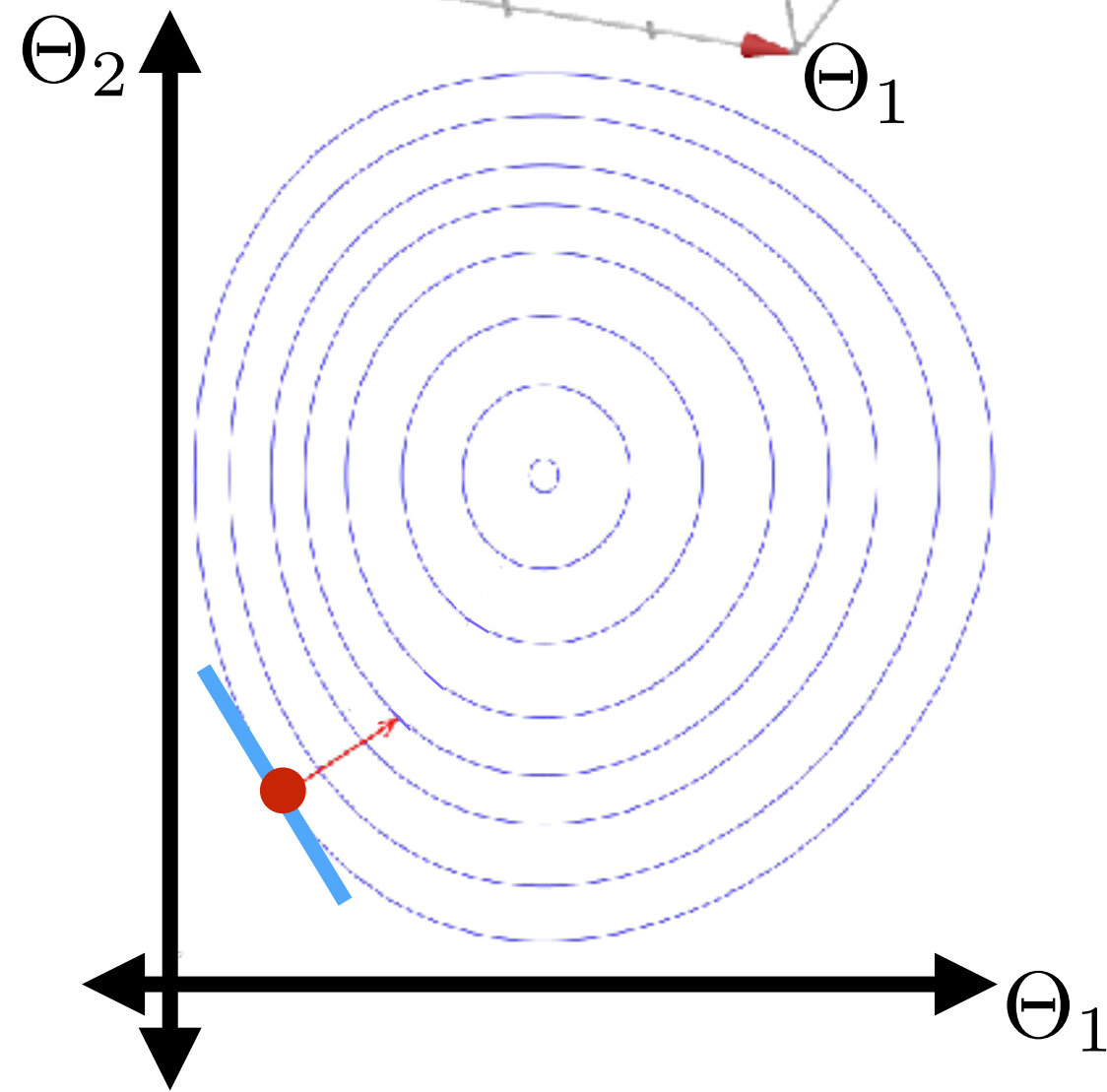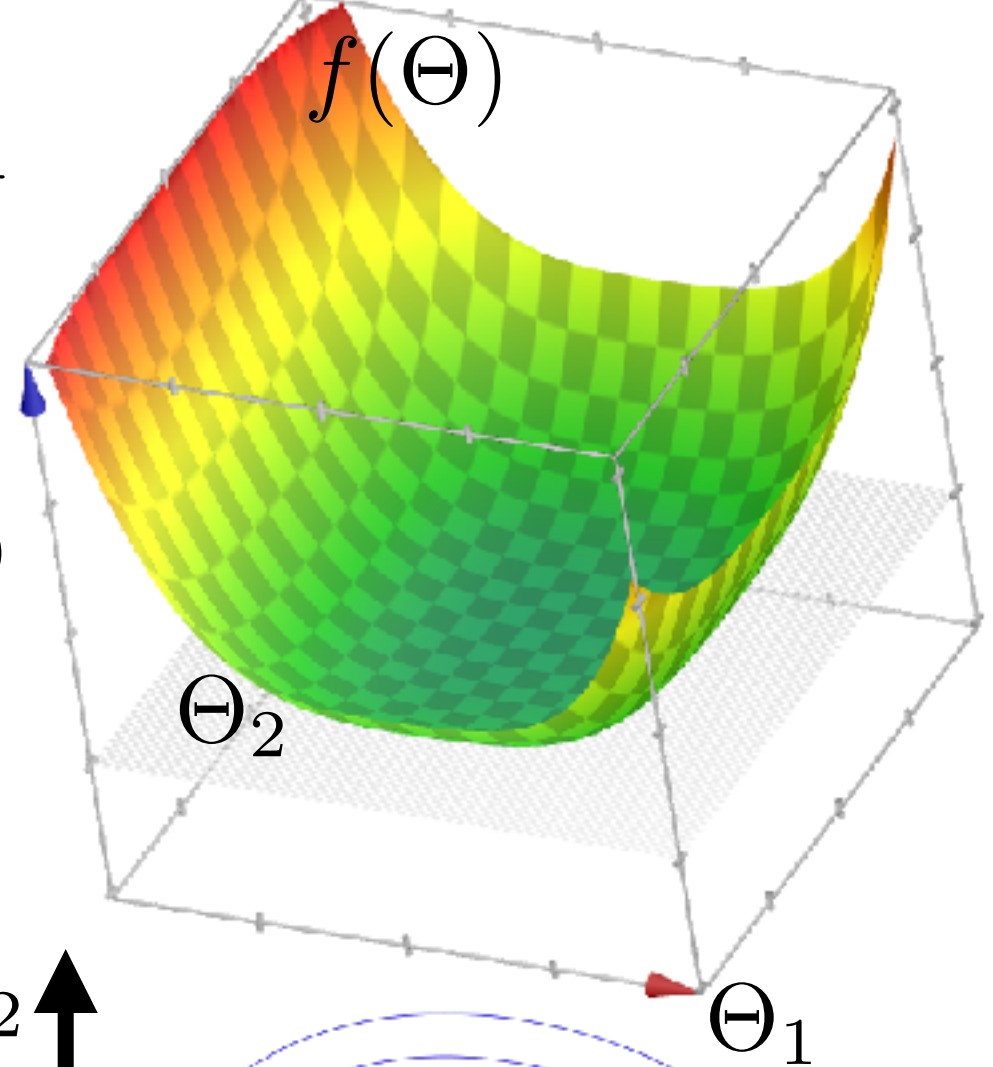Gradient-Descent $(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

Initialize $\Theta^{(0)} = \Theta_{\text{init}}$

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

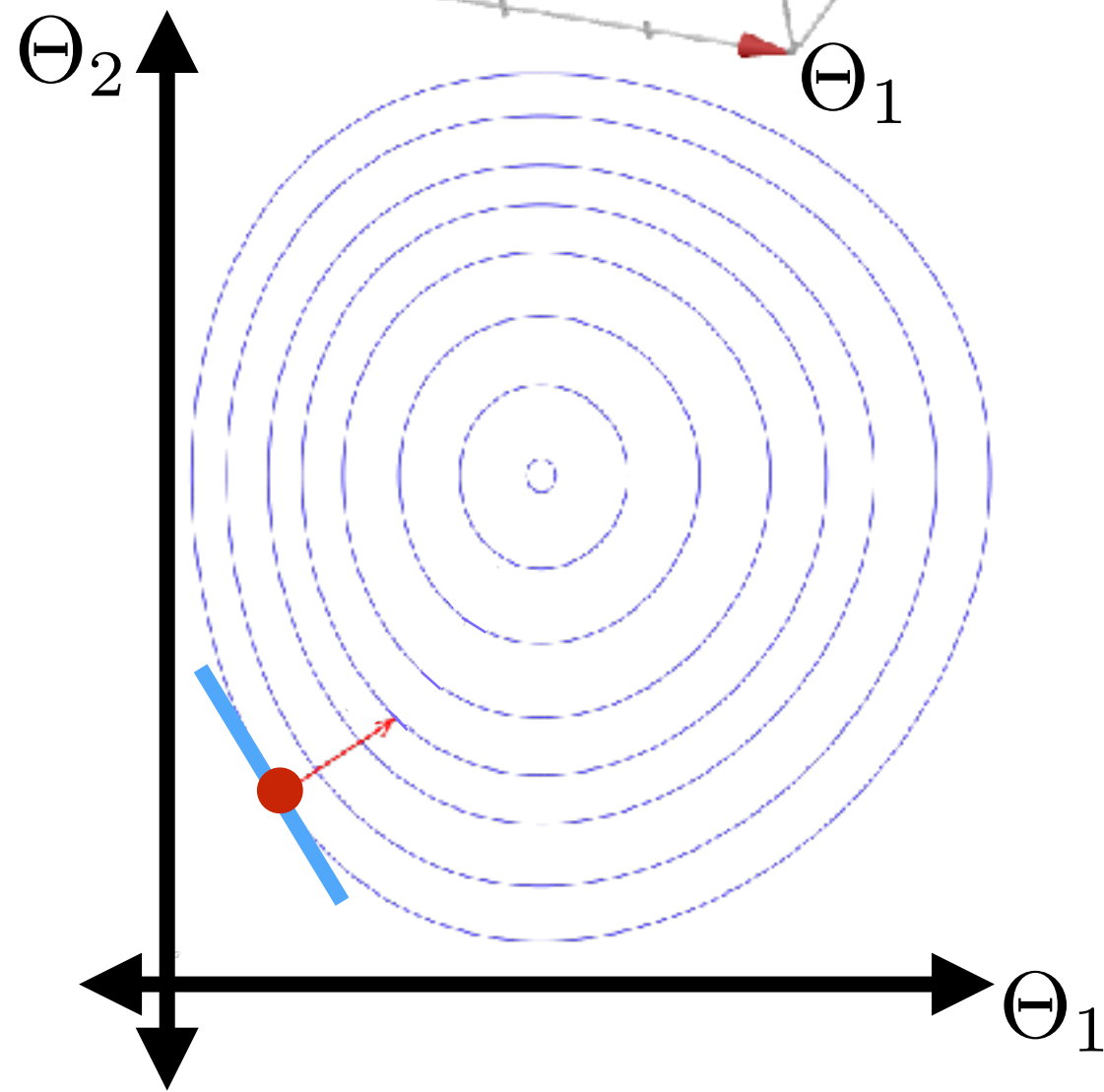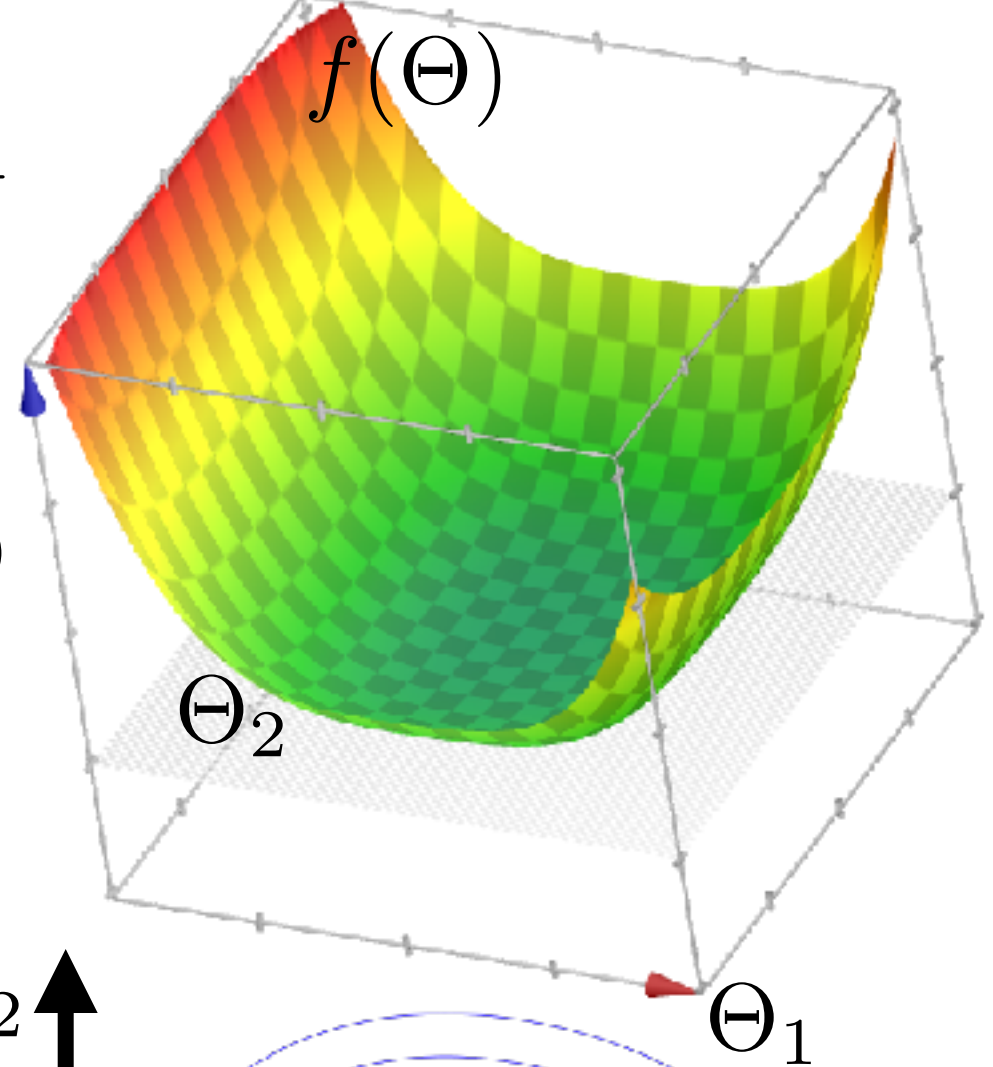Gradient-Descent $(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

Initialize $\Theta^{(0)} = \Theta_{\text{init}}$

Initialize t = 0

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$

  - with $\Theta \in \mathbb{R}^m$

```
Gradient-Descent (Θ_init, η, f, ∇_Θ f, ε)
  Initialize Θ^(0) = Θ_init
  Initialize t = 0
  repeat
```



3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

```
Gradient-Descent (Θ_init, η, f, ∇_Θ f, ε)
  Initialize Θ^(0) = Θ_init
  Initialize t = 0
  repeat
    t = t + 1
```



3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\text{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$
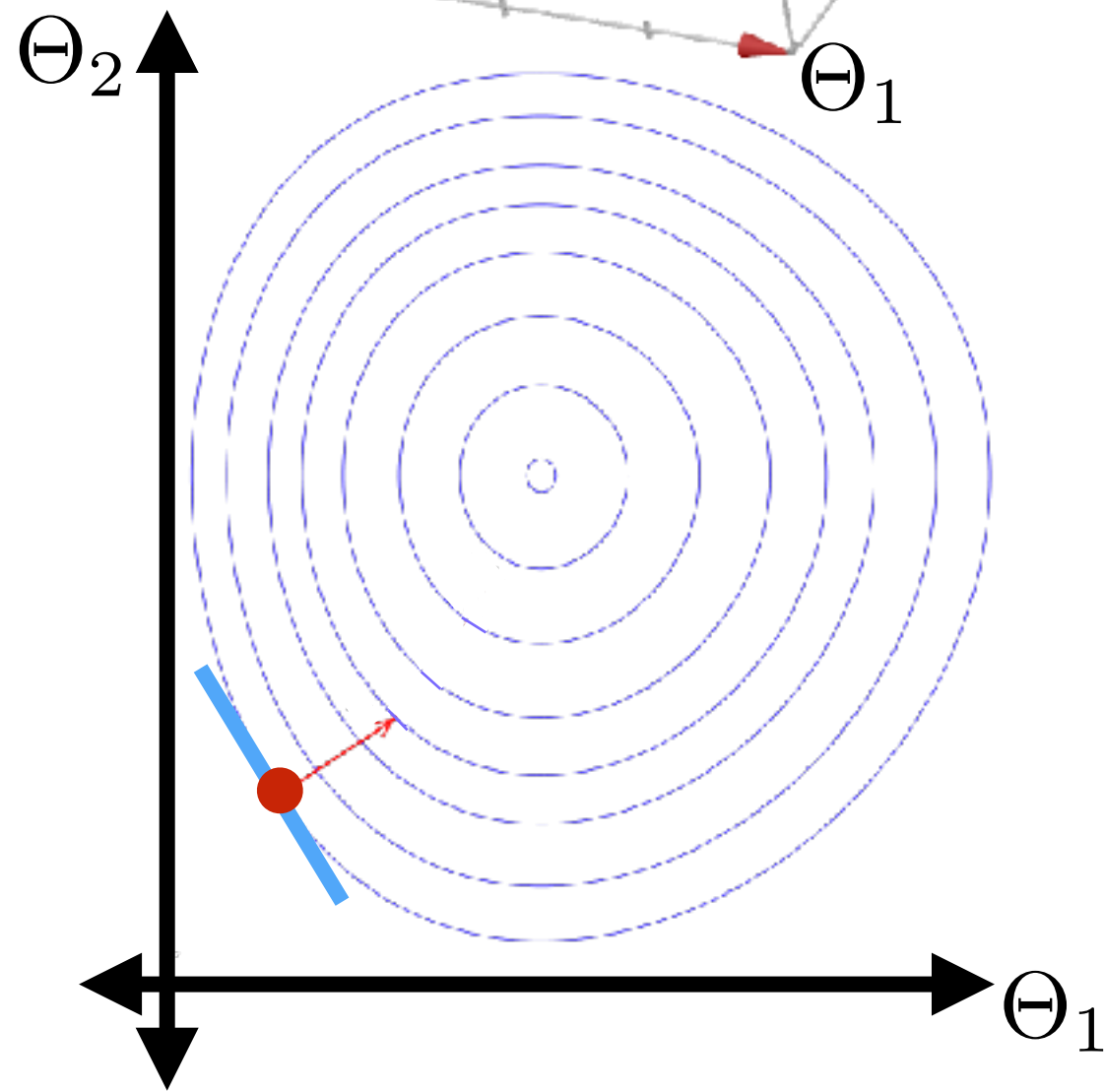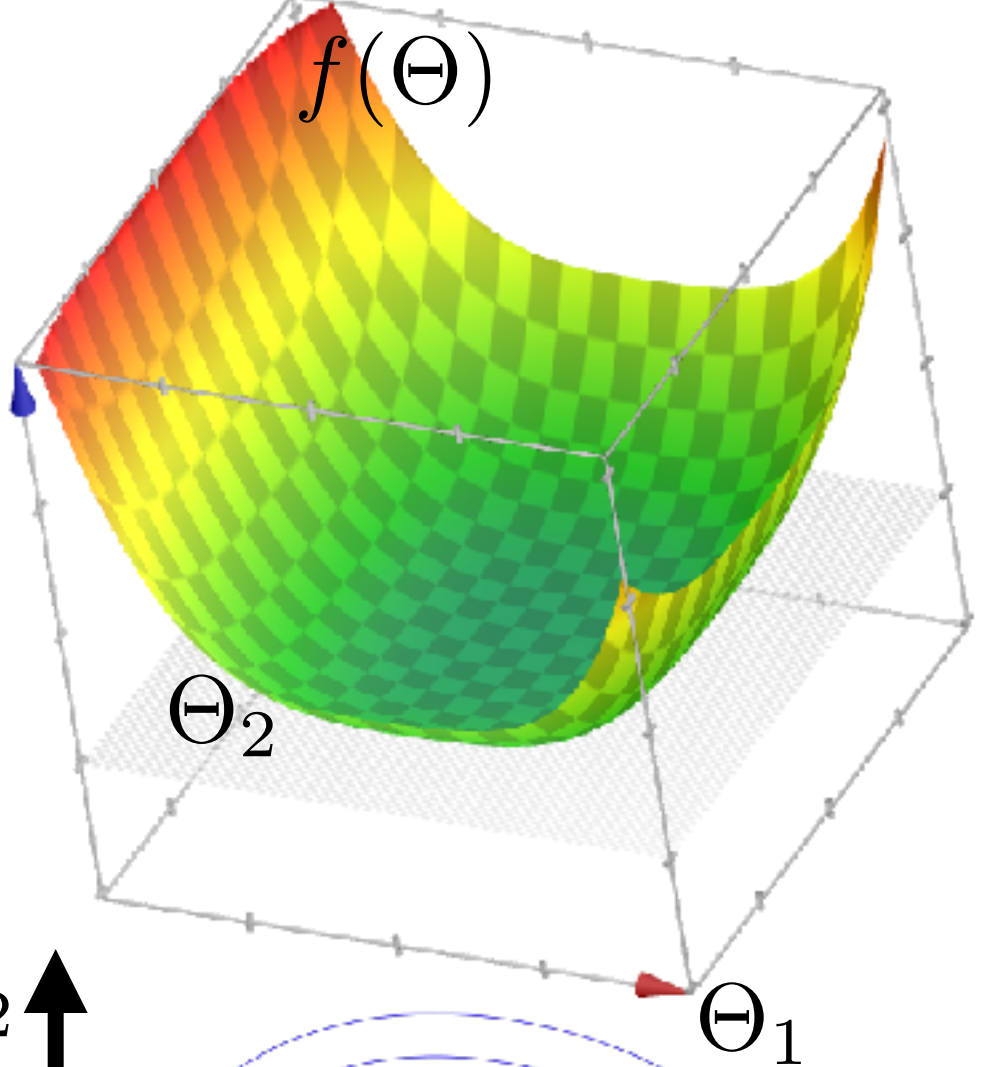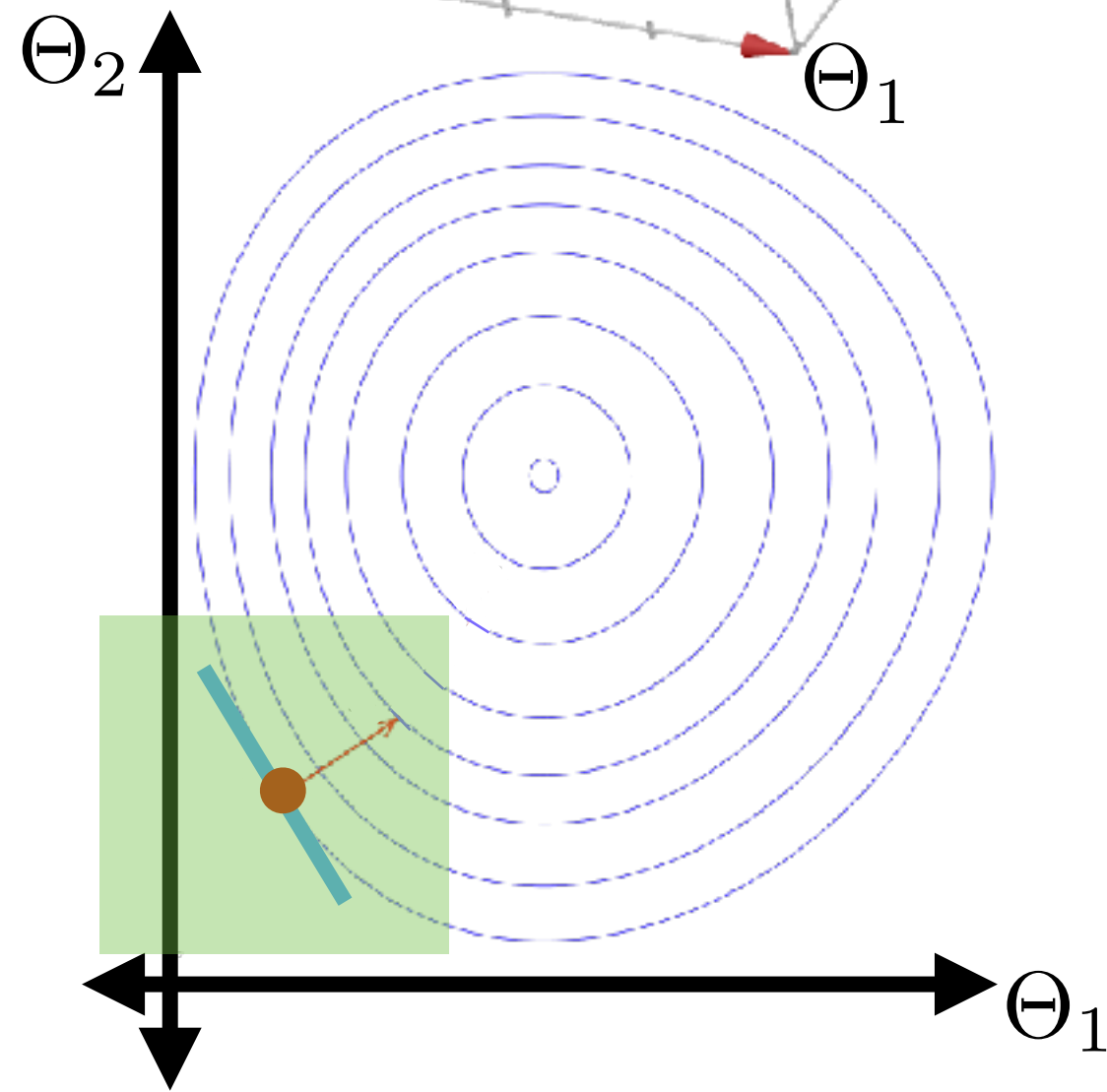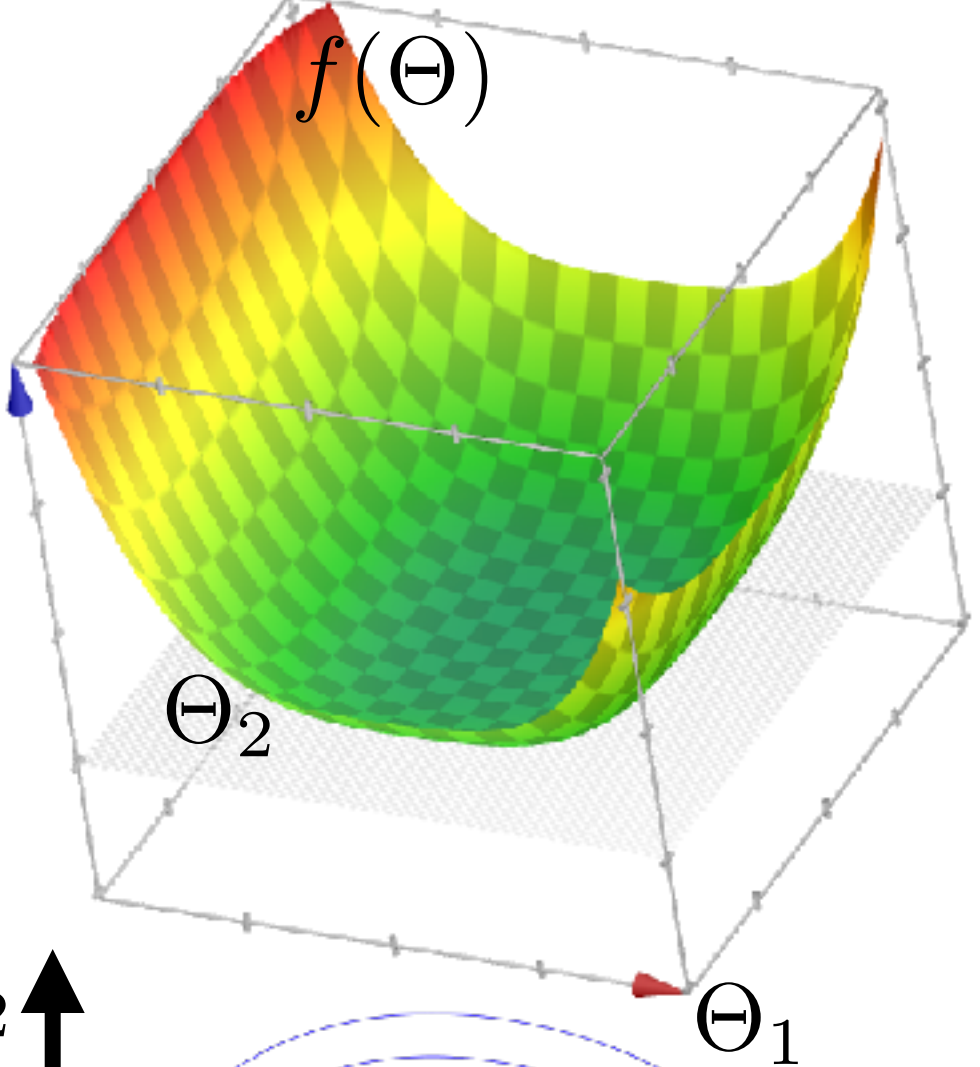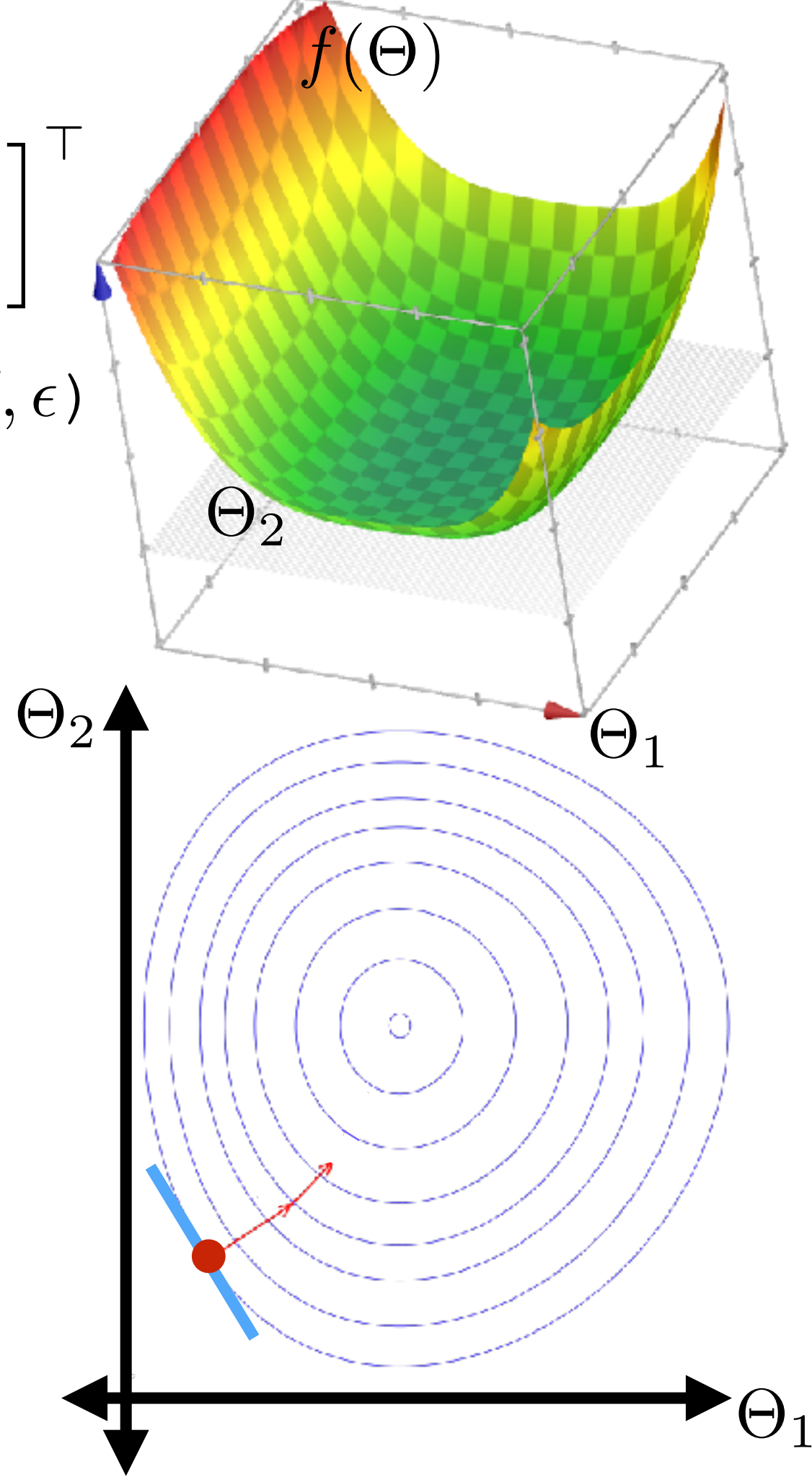
$\quad \text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\quad \text{Initialize } \texttt{t = 0}$

**repeat**

$\quad\quad \texttt{t = t + 1}$

$\quad\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

```
Gradient-Descent (Θ_init, η, f, ∇_Θ f, ε)
  Initialize Θ^(0) = Θ_init
  Initialize t = 0
```

**repeat**
```
  t = t + 1
```
$$\Theta^{(t)} = \boxed{\Theta^{(t-1)}} - \eta \nabla_\Theta f(\Theta^{(t-1)})$$



3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[\dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m}\right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\mathtt{Gradient\text{-}Descent}\,(\Theta_{\mathrm{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

$\mathtt{Initialize}\ \Theta^{(0)} = \Theta_{\mathrm{init}}$

$\mathtt{Initialize\ t\ =\ 0}$

**repeat**

   $\mathtt{t\ =\ t\ +\ 1}$

   $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$



3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\text{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$
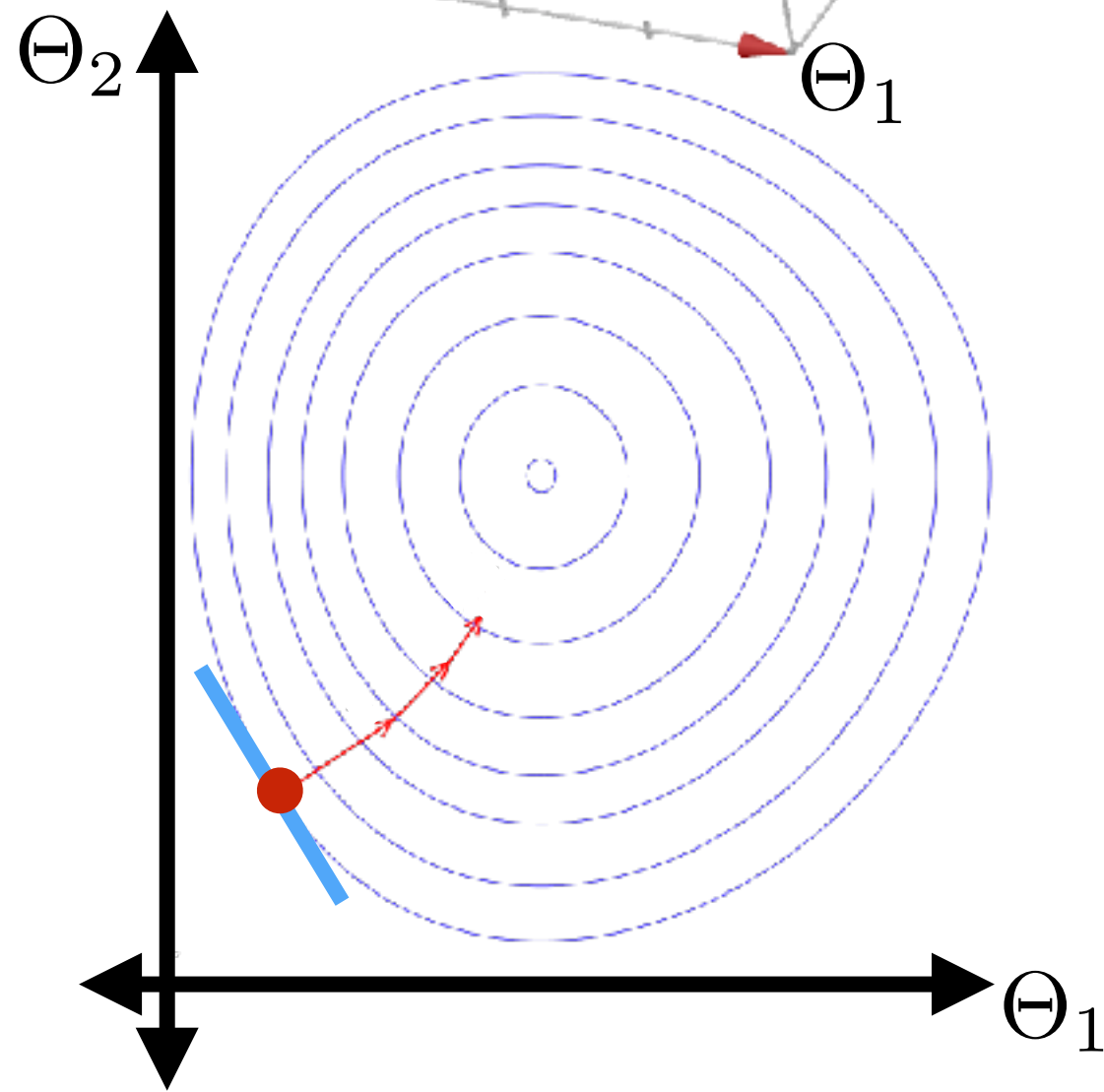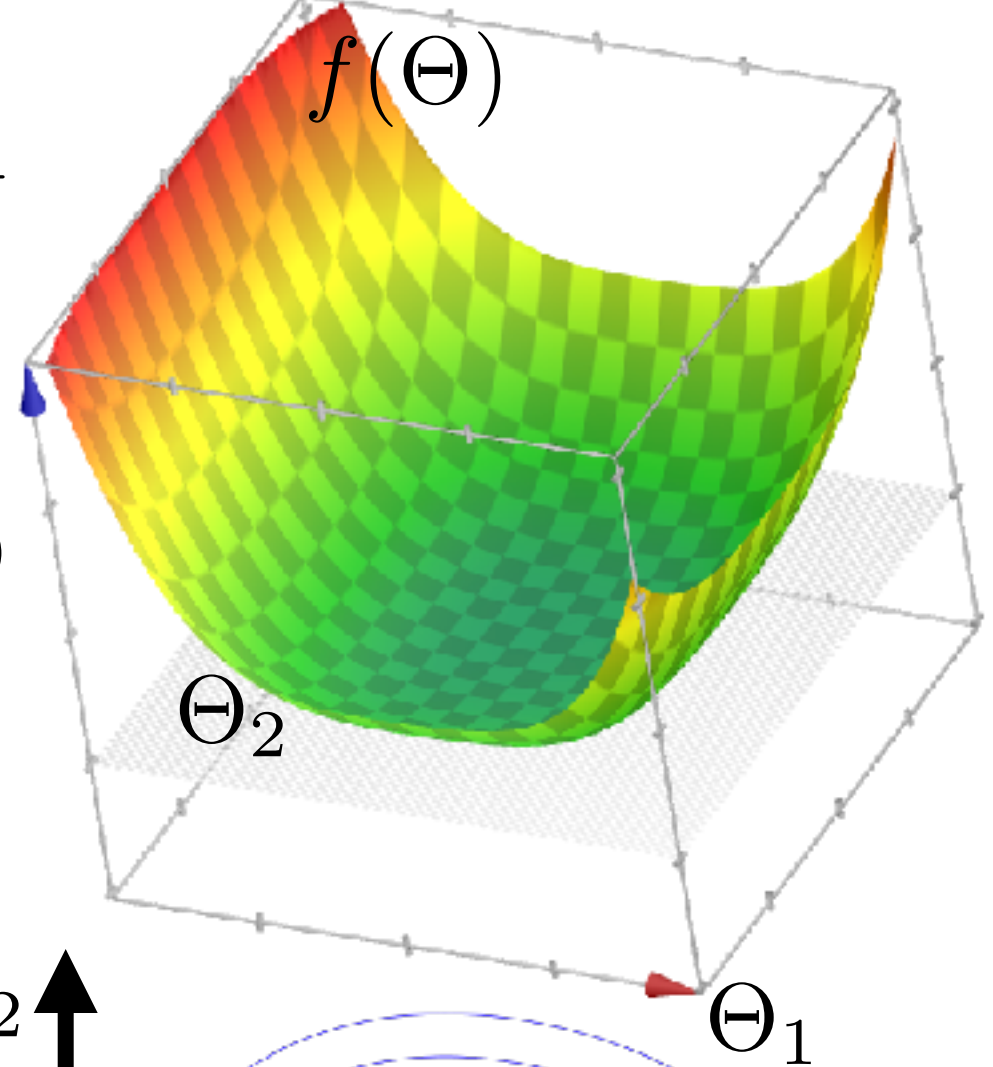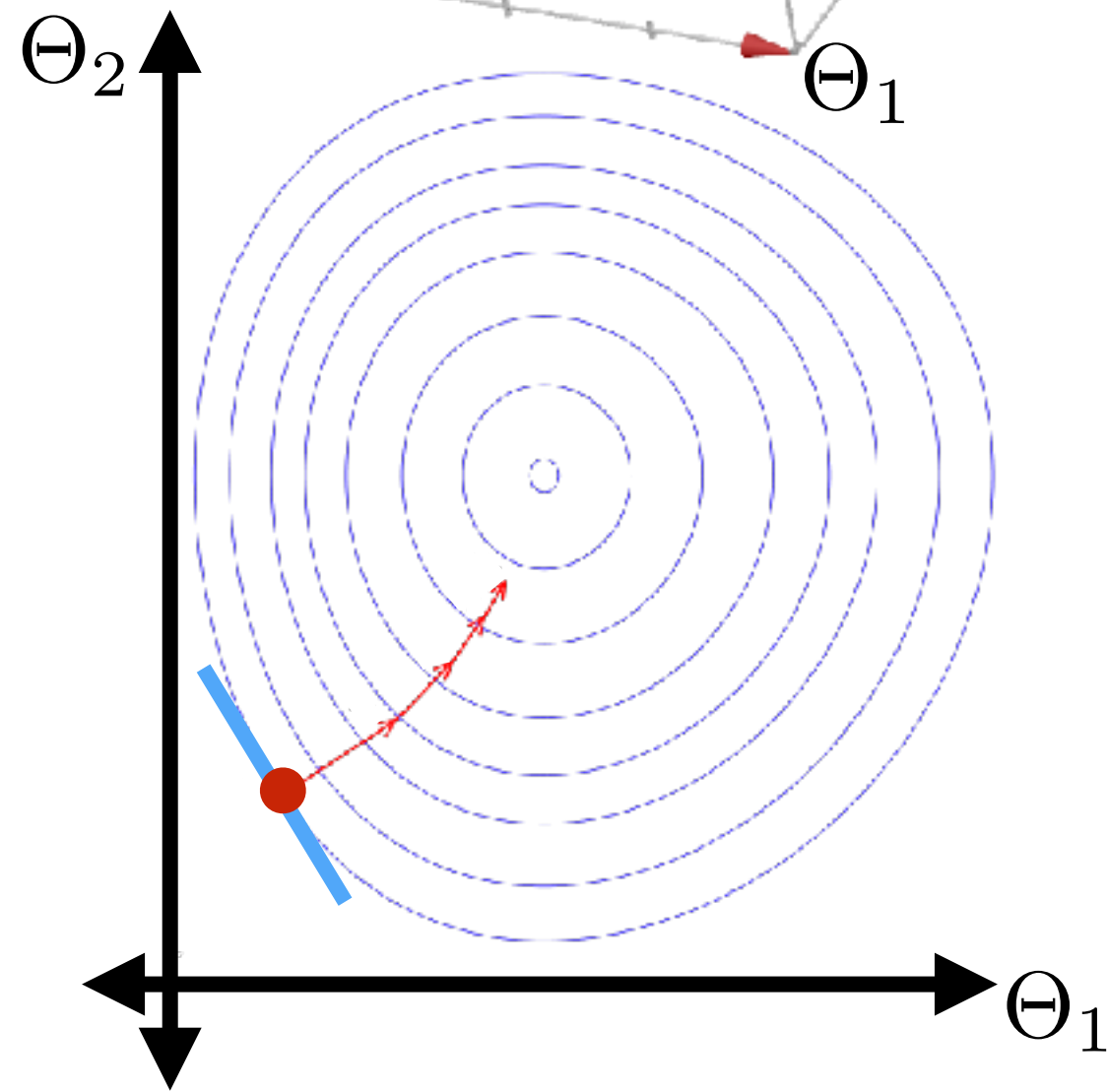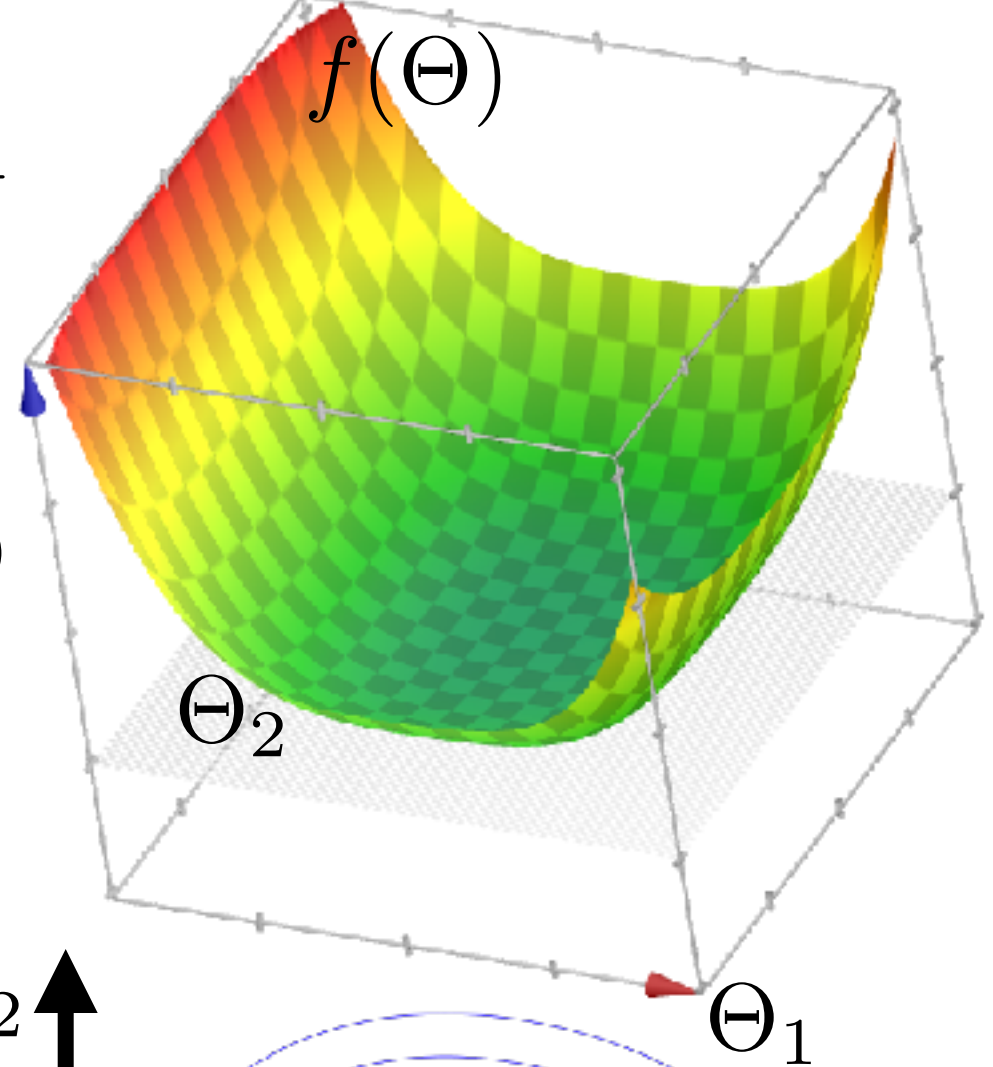
$\quad \text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\quad \text{Initialize t = 0}$

**repeat**

$\quad \text{t = t + 1}$

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$



3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\text{Gradient-Descent}(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

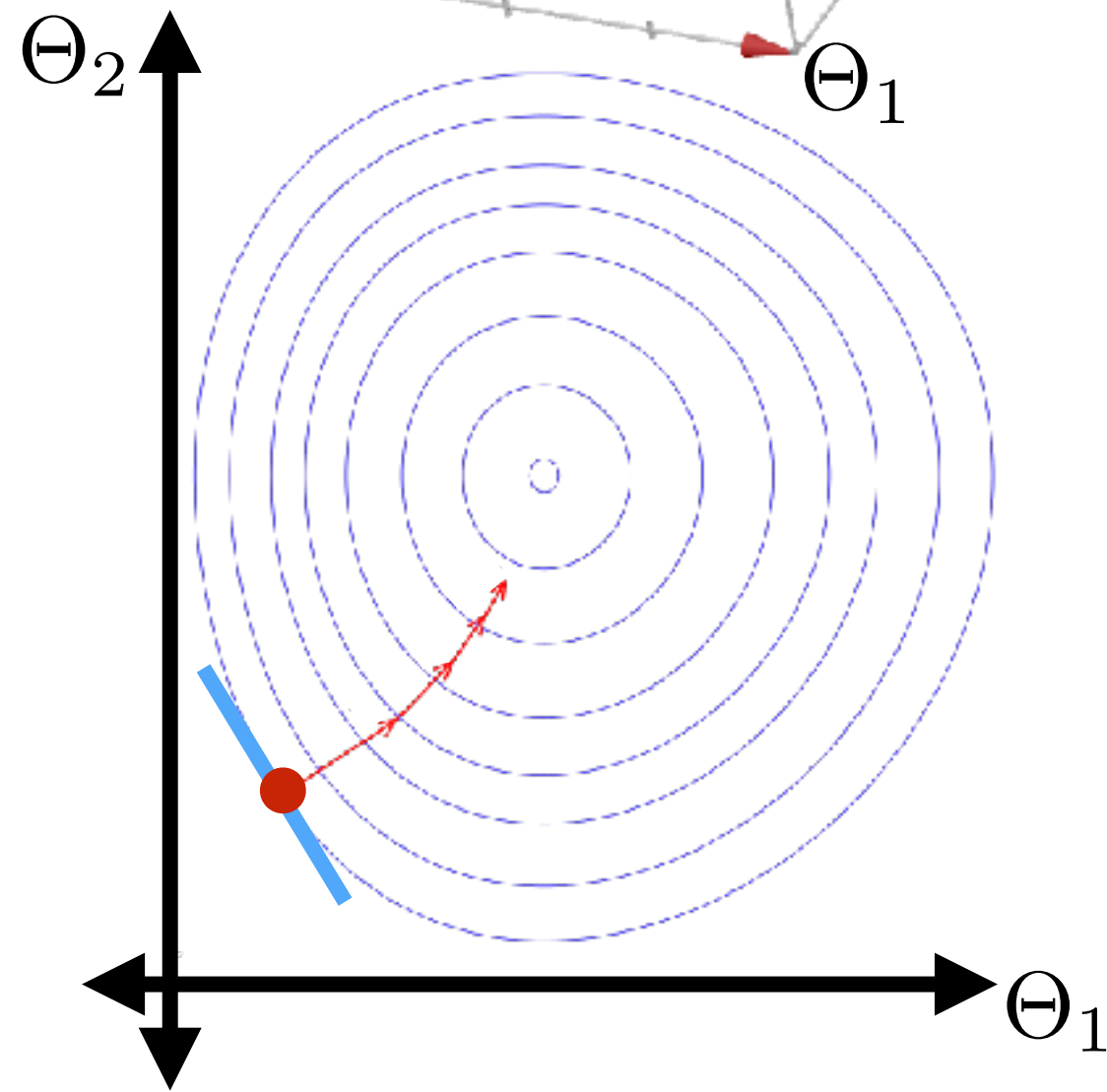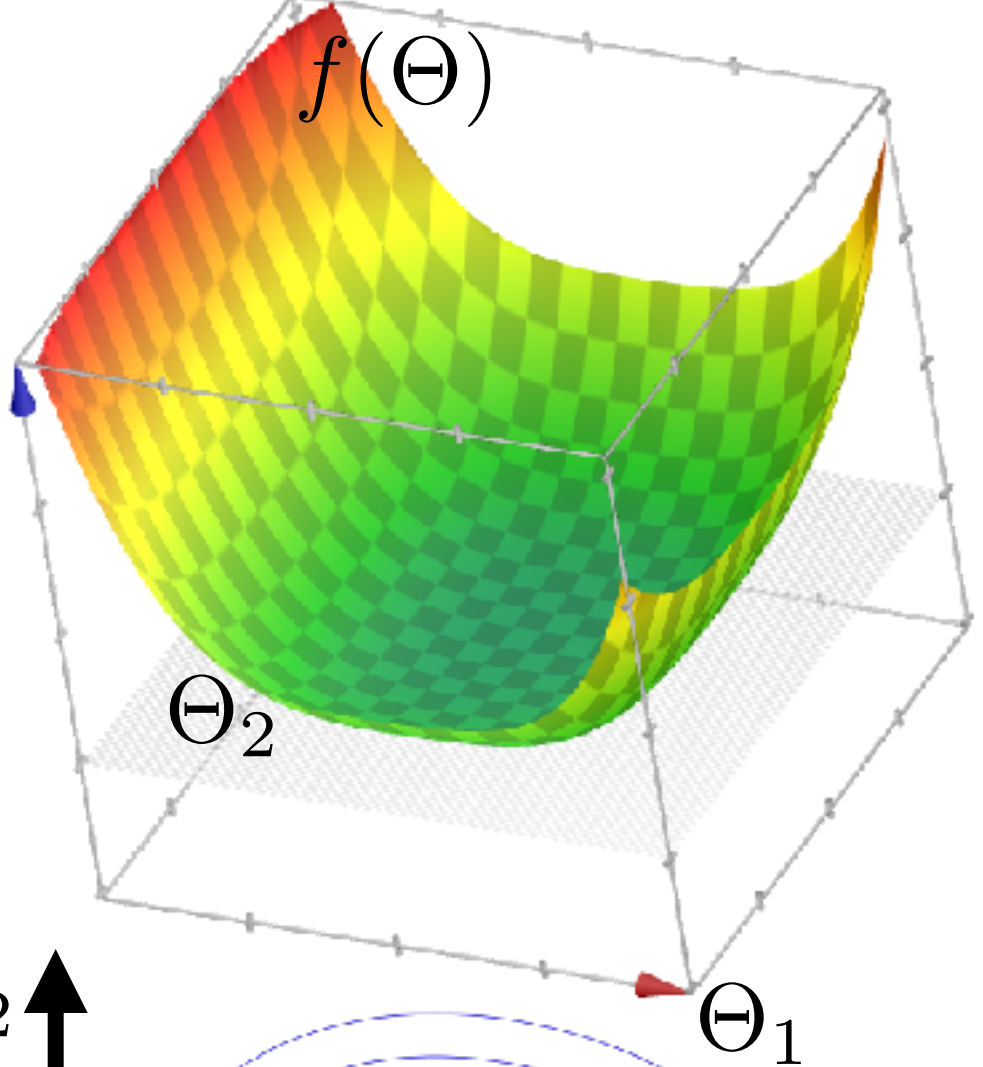$\text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\text{Initialize t = 0}$

**repeat**

$\quad \text{t = t + 1}$

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$



3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[\dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m}\right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\text{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$
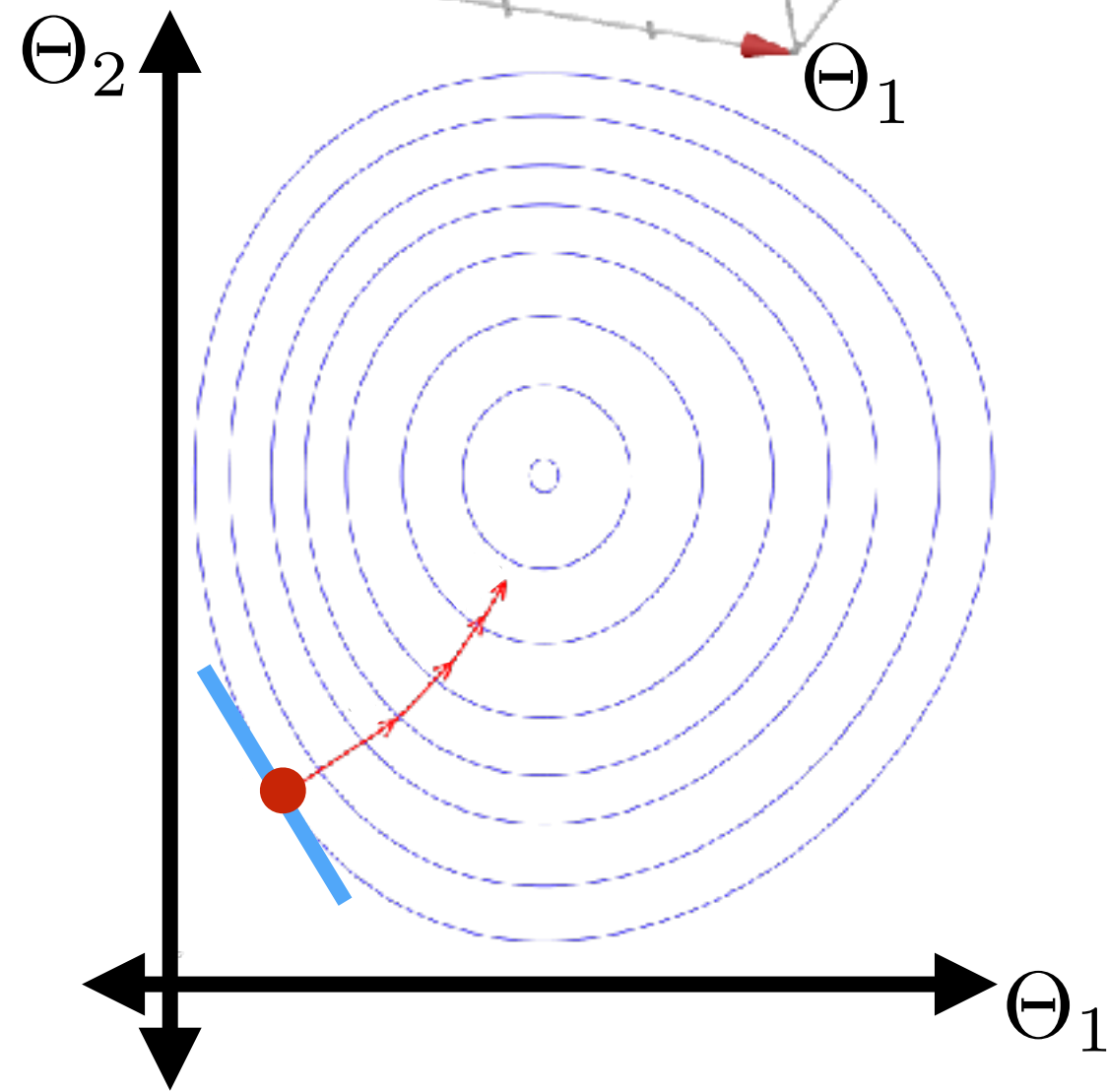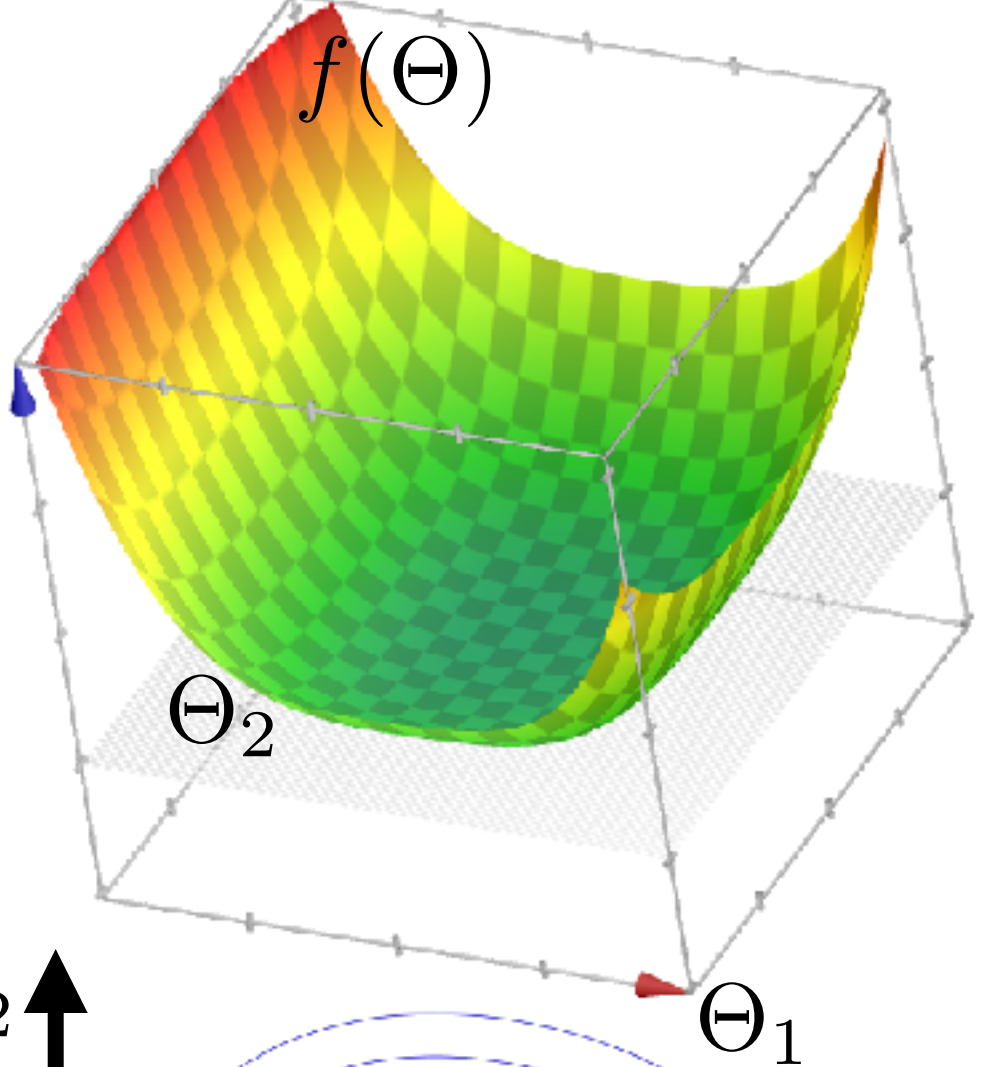
$\text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\text{Initialize } \texttt{t = 0}$

**repeat**

$\quad \texttt{t = t + 1}$

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$



3

# Gradient descent

- Gradient $\nabla_{\Theta} f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^{\top}$
  - with $\Theta \in \mathbb{R}^m$

$\texttt{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_{\Theta} f, \epsilon)$

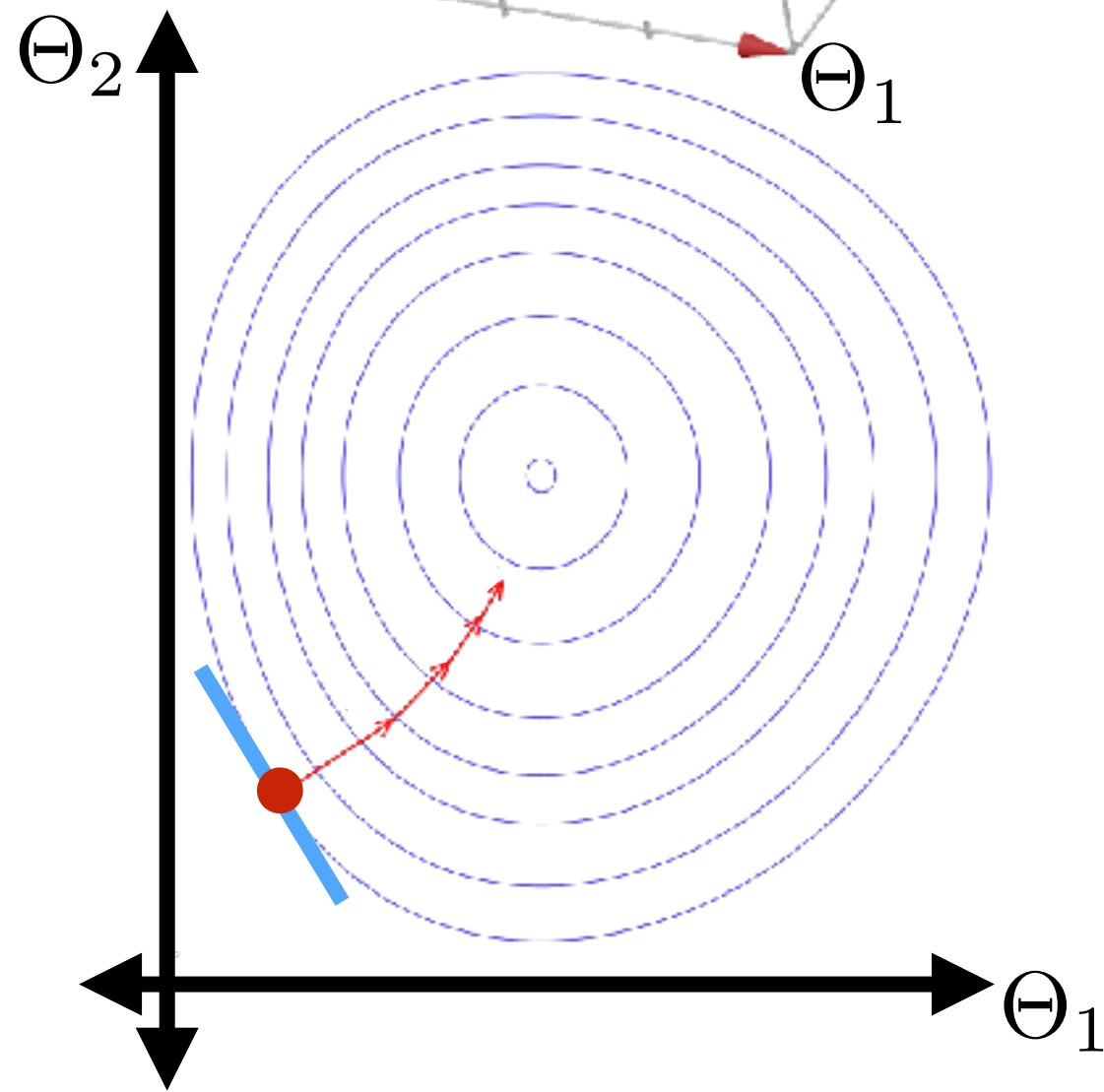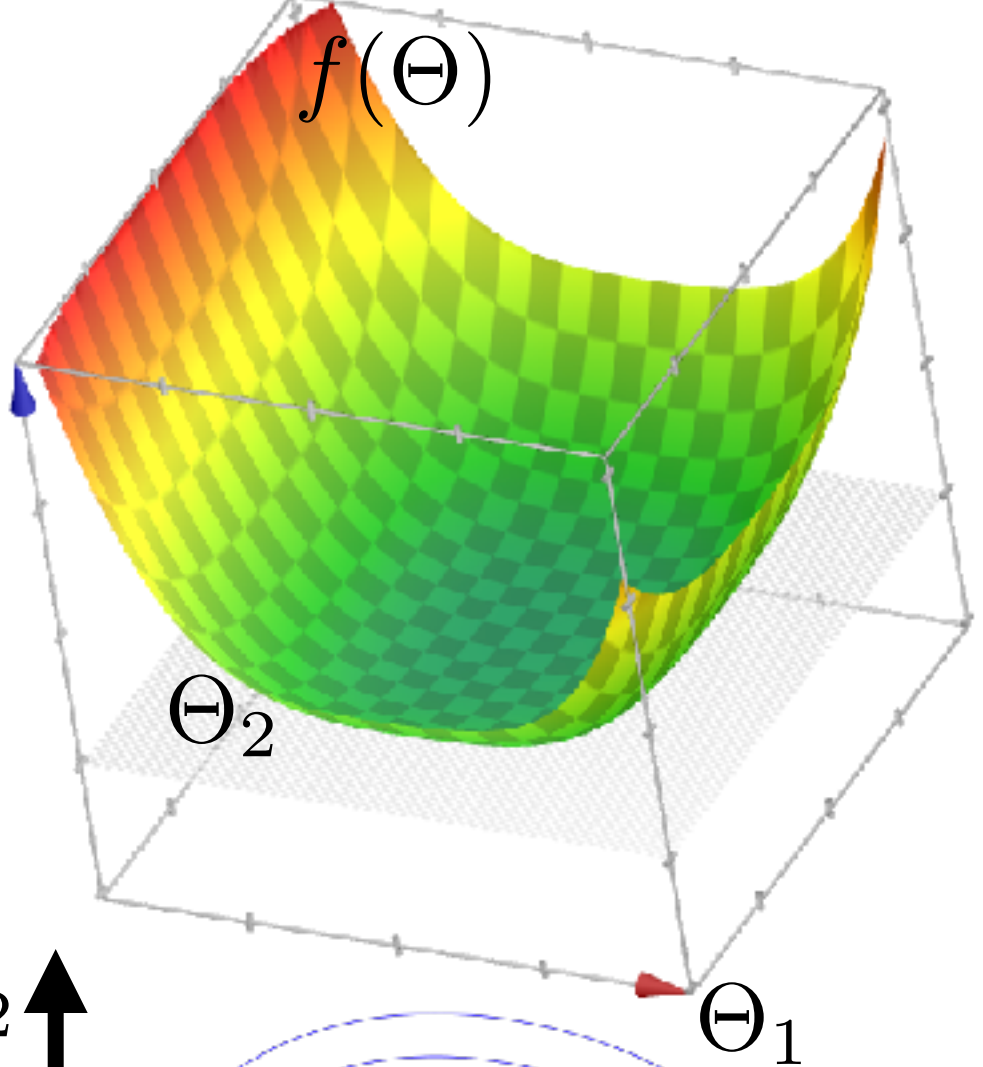$\texttt{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\texttt{Initialize t = 0}$

**repeat**

$\quad \texttt{t = t + 1}$

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_{\Theta} f(\Theta^{(t-1)})$

# Gradient descent

- Gradient $\nabla_\Theta f = \left[\dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m}\right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\text{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

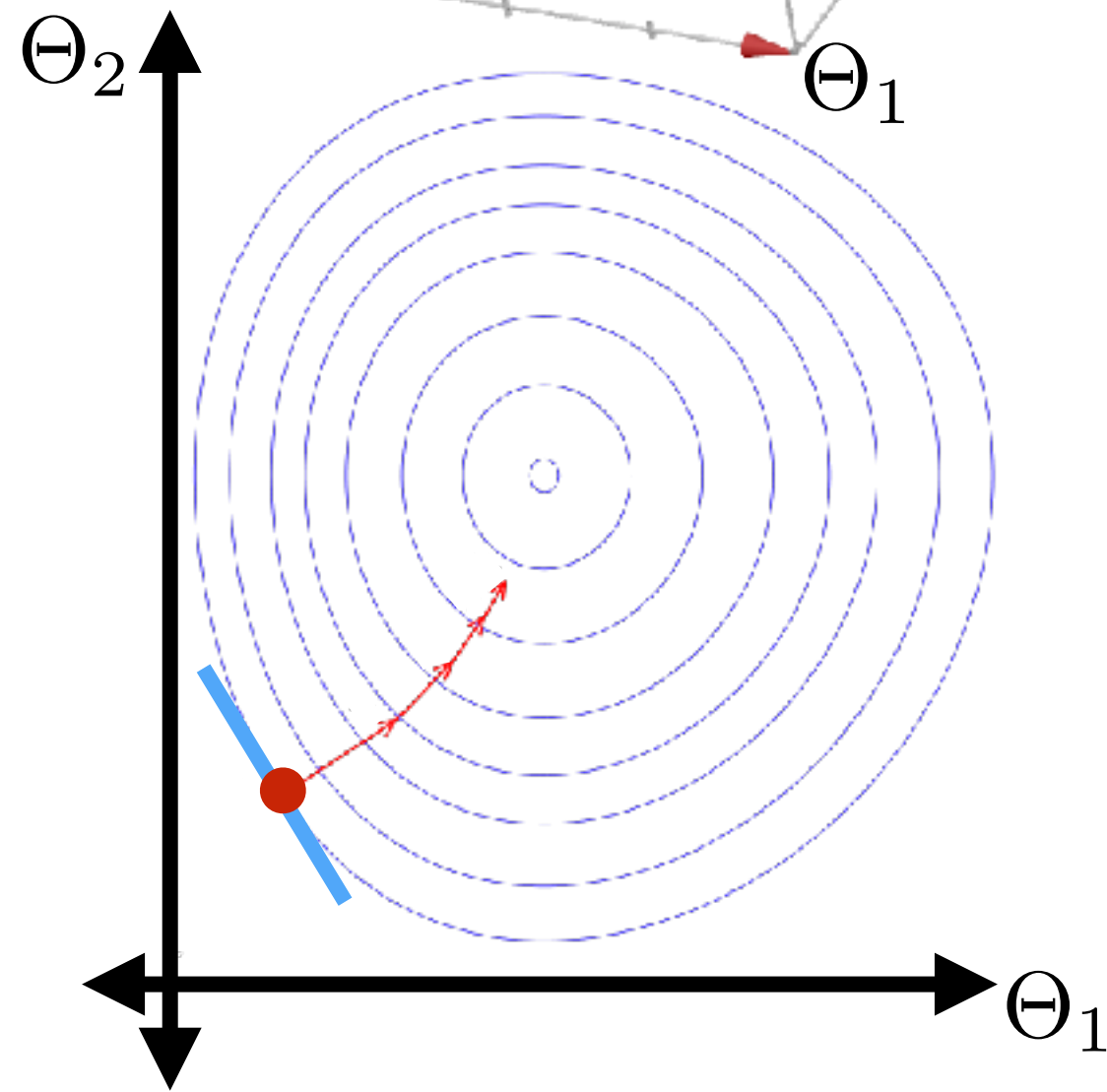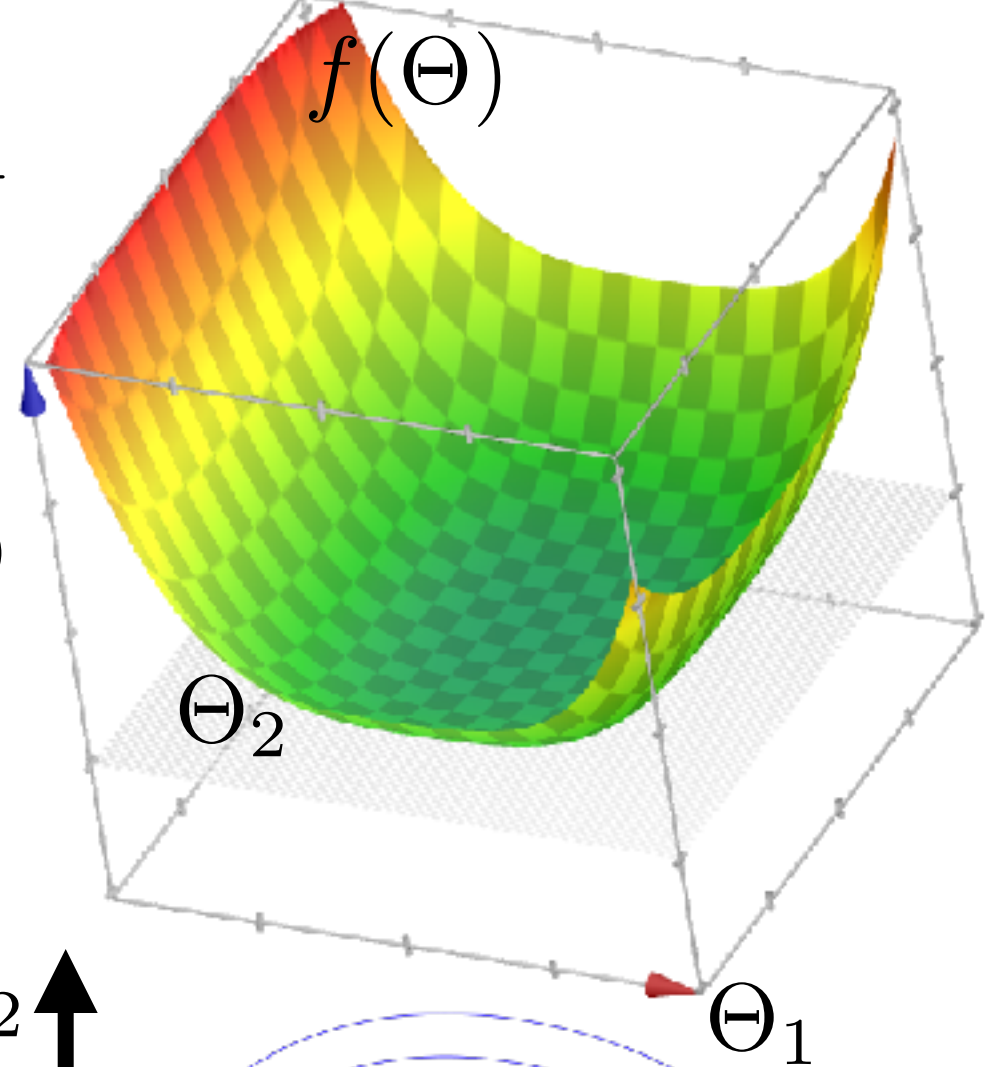$\text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\text{Initialize } \texttt{t = 0}$

**repeat**

$\quad \texttt{t = t + 1}$

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[\dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m}\right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\text{Gradient-Descent}(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

$\quad \text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

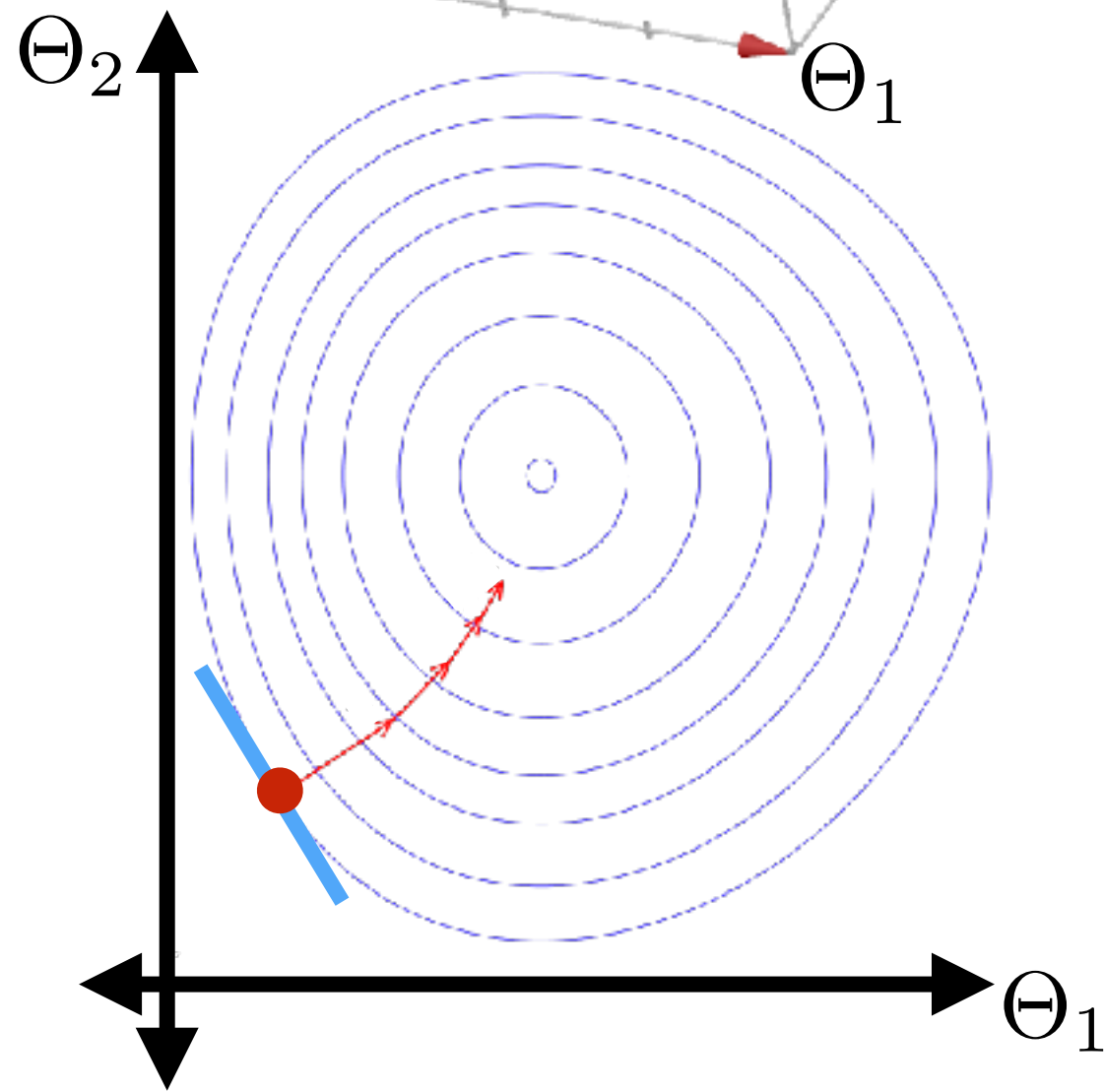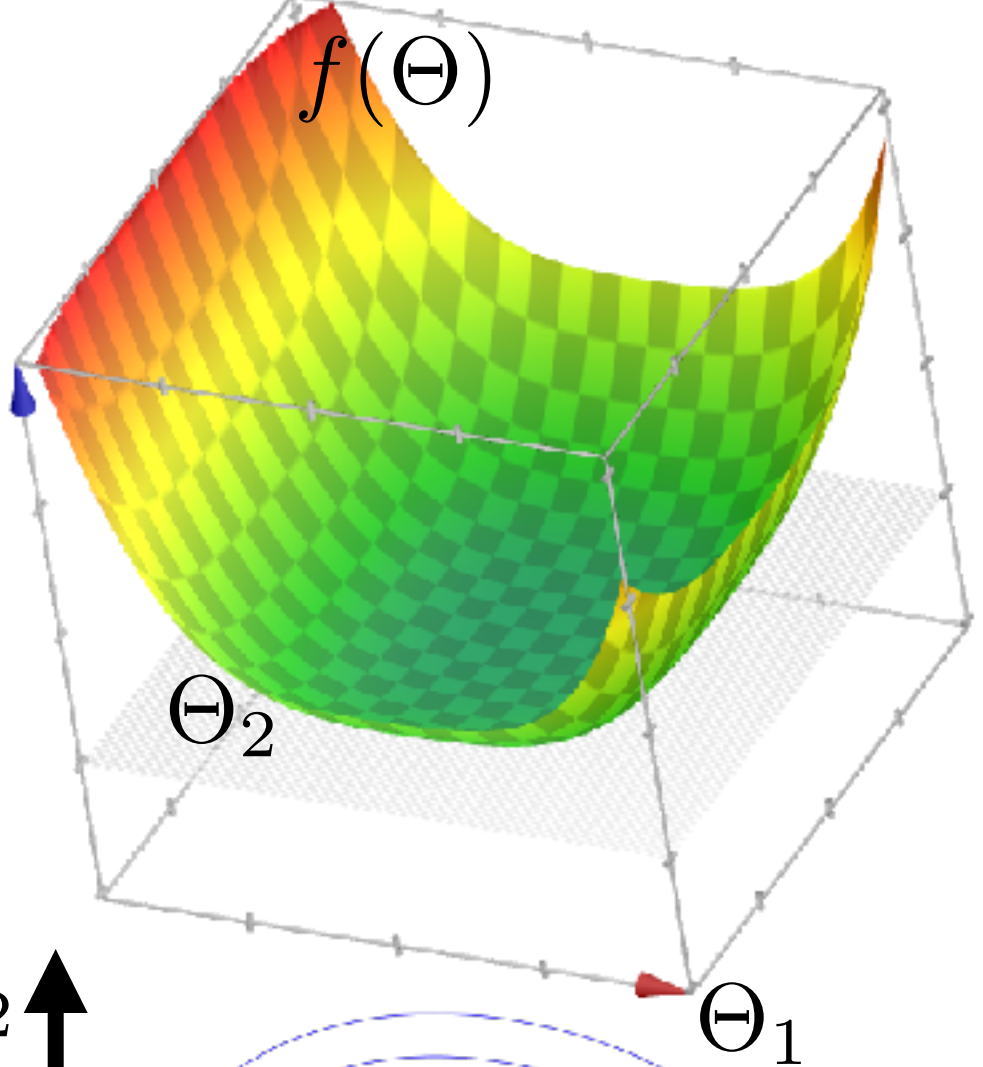$\quad \text{Initialize } \texttt{t = 0}$

**repeat**

$\quad \texttt{t = t + 1}$

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until**



3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$f(\Theta)$

$\Theta_2$

$\Theta_1$

```
Gradient-Descent (Θ_init, η, f, ∇_Θ f, ε)
  Initialize Θ^(0) = Θ_init
  Initialize t = 0
```
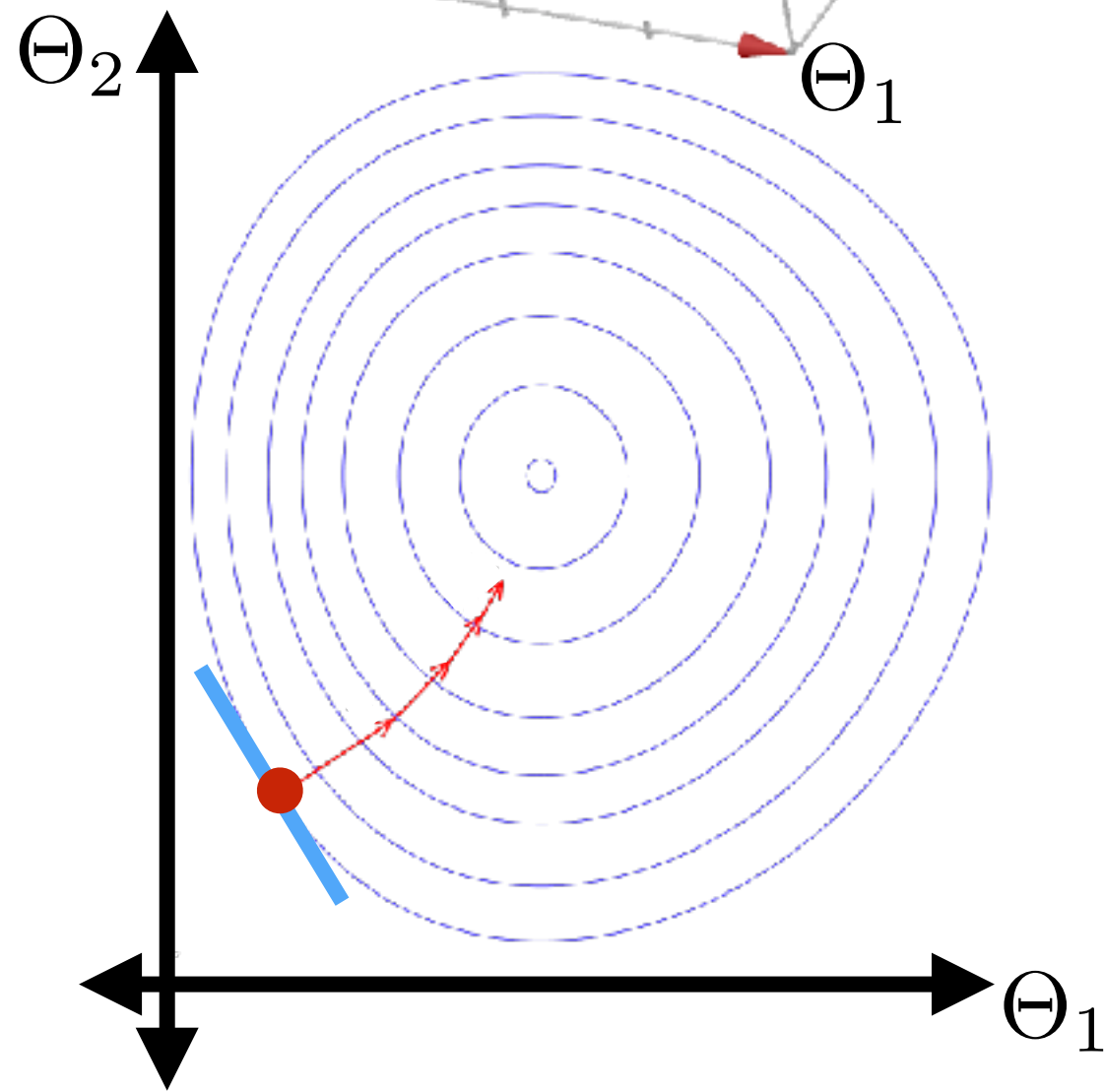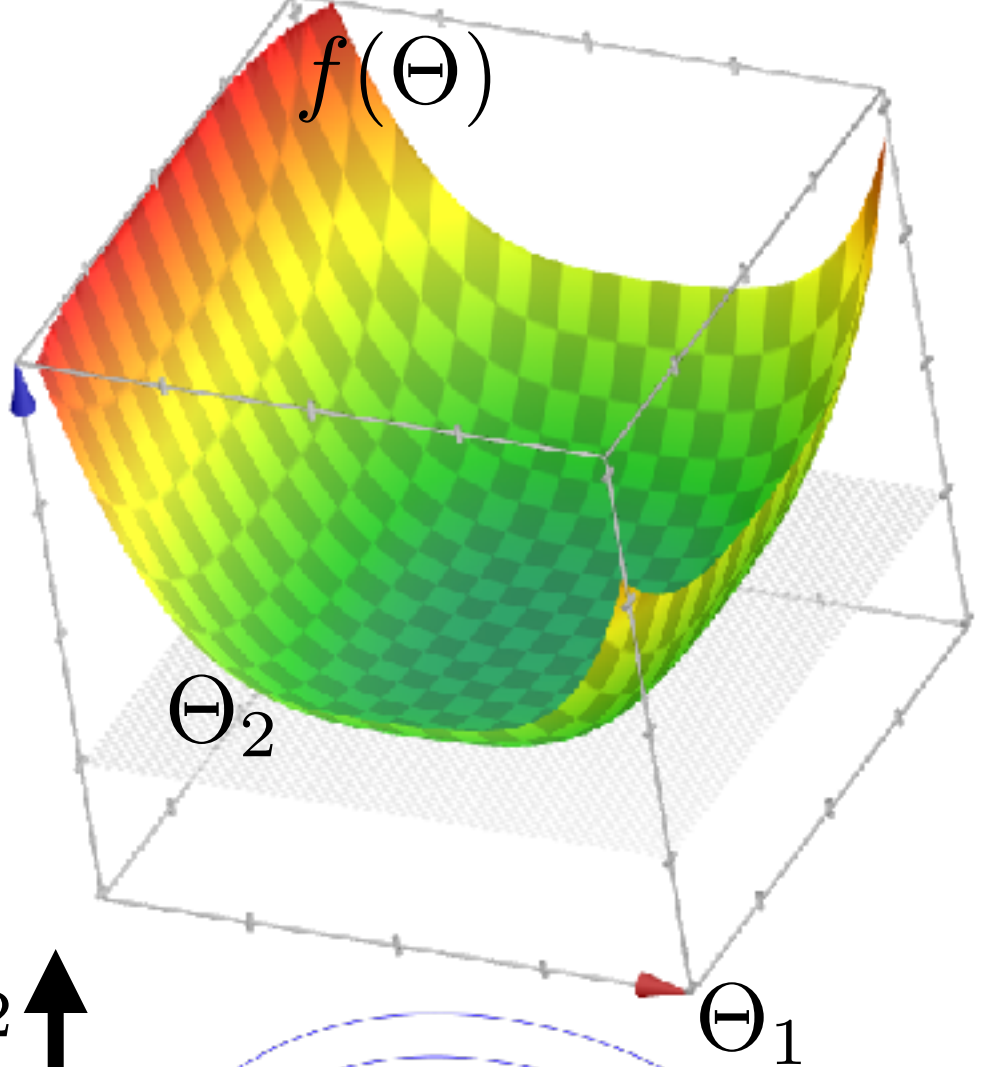
**repeat**

  $\quad$ `t = t + 1`

  $\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

$\Theta_2$

$\Theta_1$

# Gradient descent

- Gradient $\nabla_\Theta f = \left[\dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m}\right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\text{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

$\text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

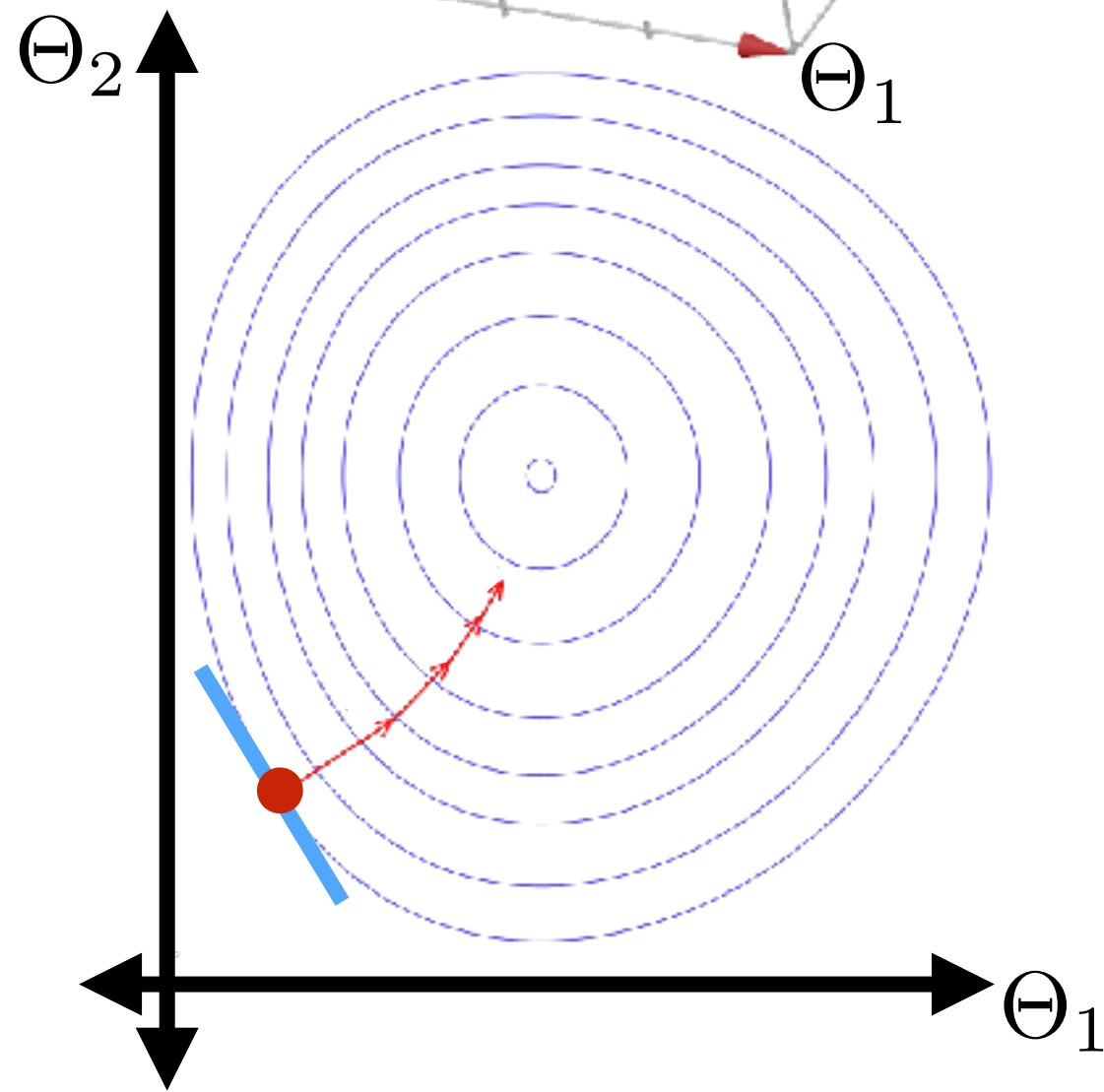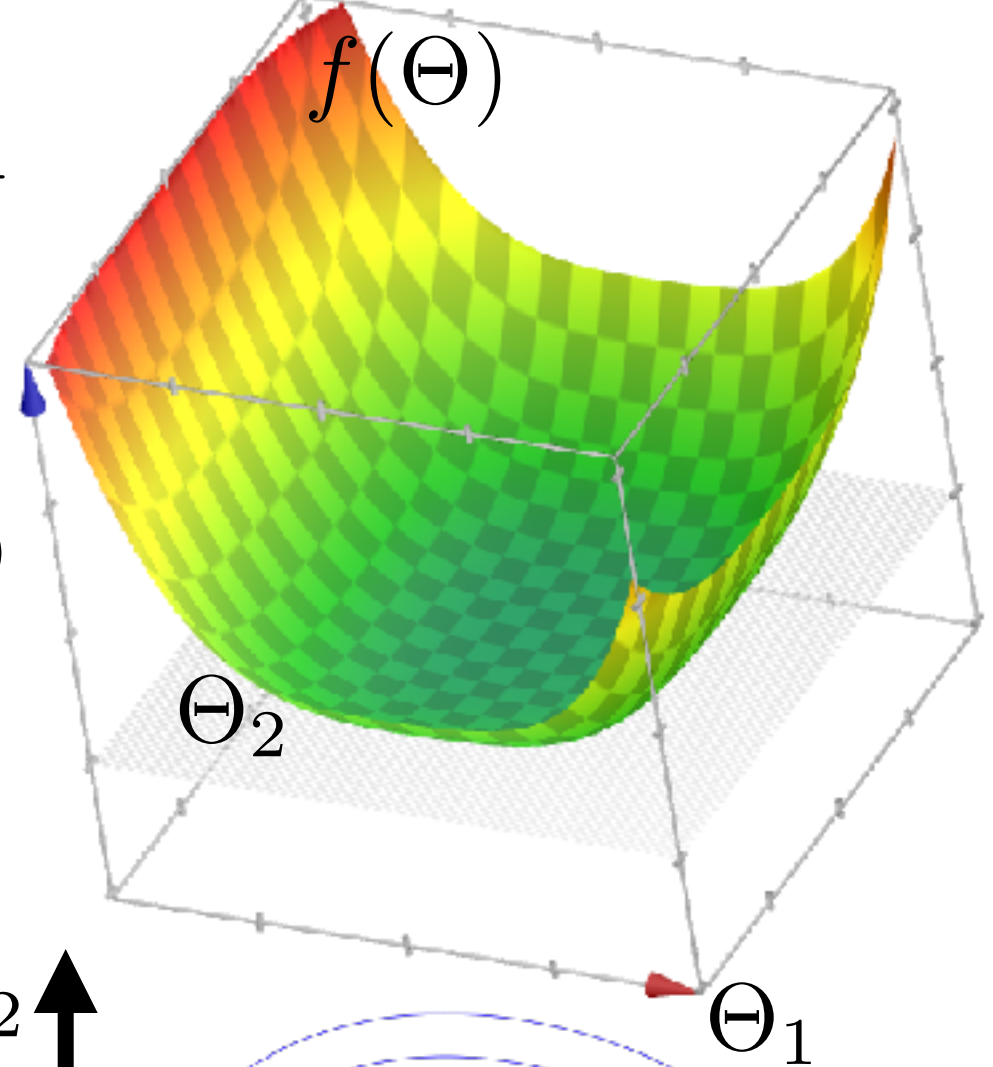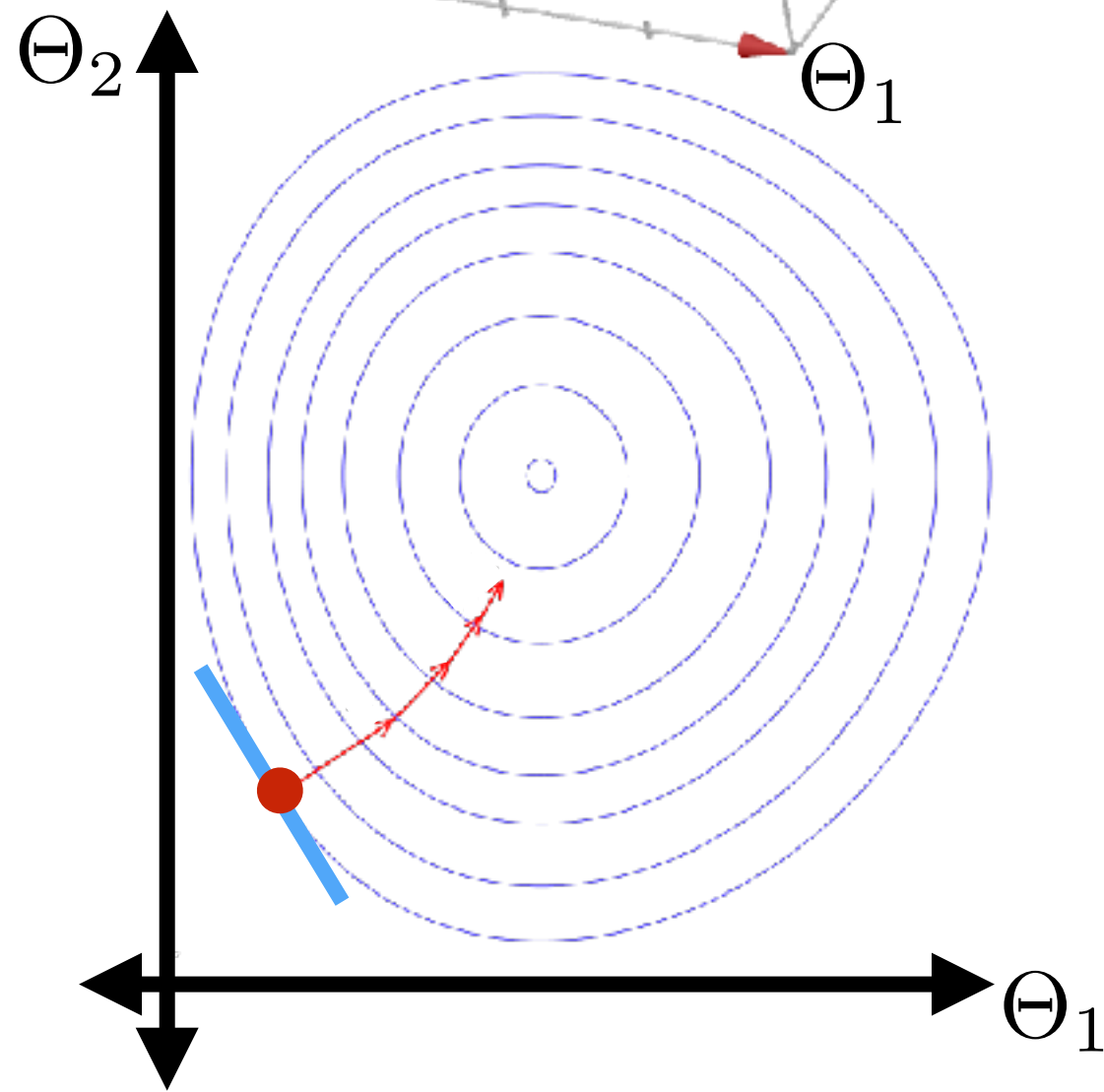$\text{Initialize } \texttt{t = 0}$

**repeat**

$\quad \texttt{t = t + 1}$

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[\dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m}\right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\texttt{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

$\texttt{Initialize}\ \Theta^{(0)} = \Theta_{\text{init}}$

$\texttt{Initialize t = 0}$

**repeat**

   $\texttt{t = t + 1}$

   $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \dots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

Gradient-Descent $(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

Initialize $\Theta^{(0)} = \Theta_{\text{init}}$

Initialize t = 0

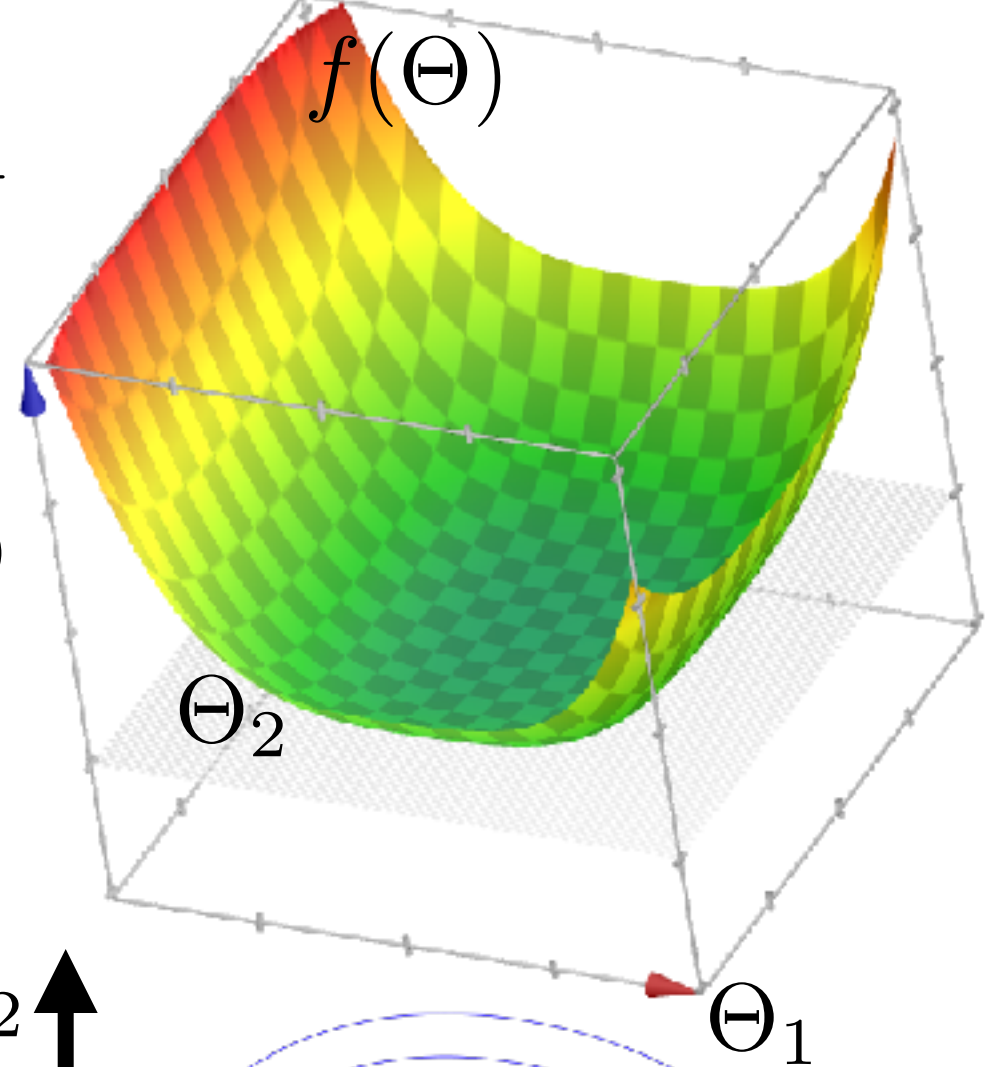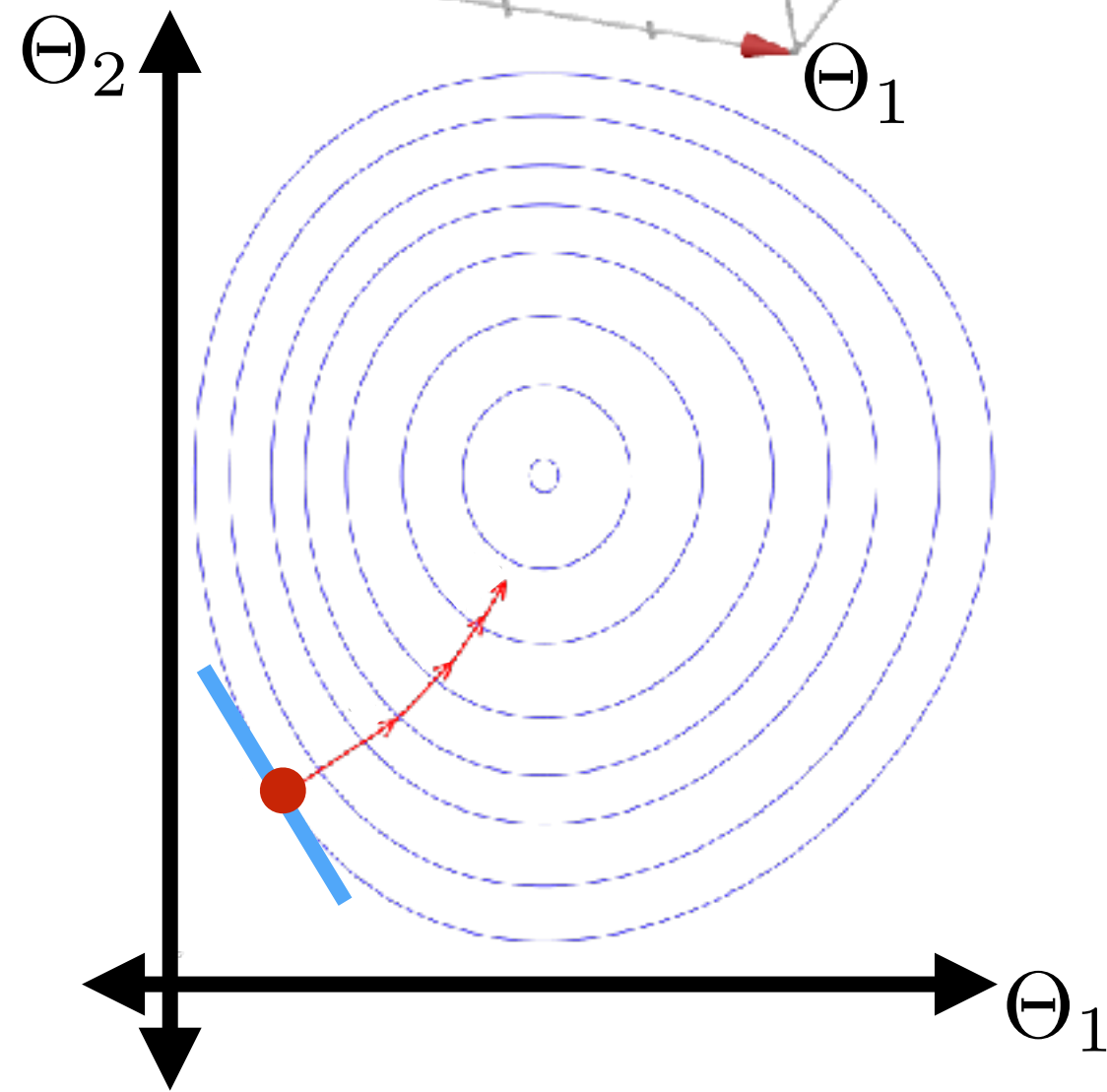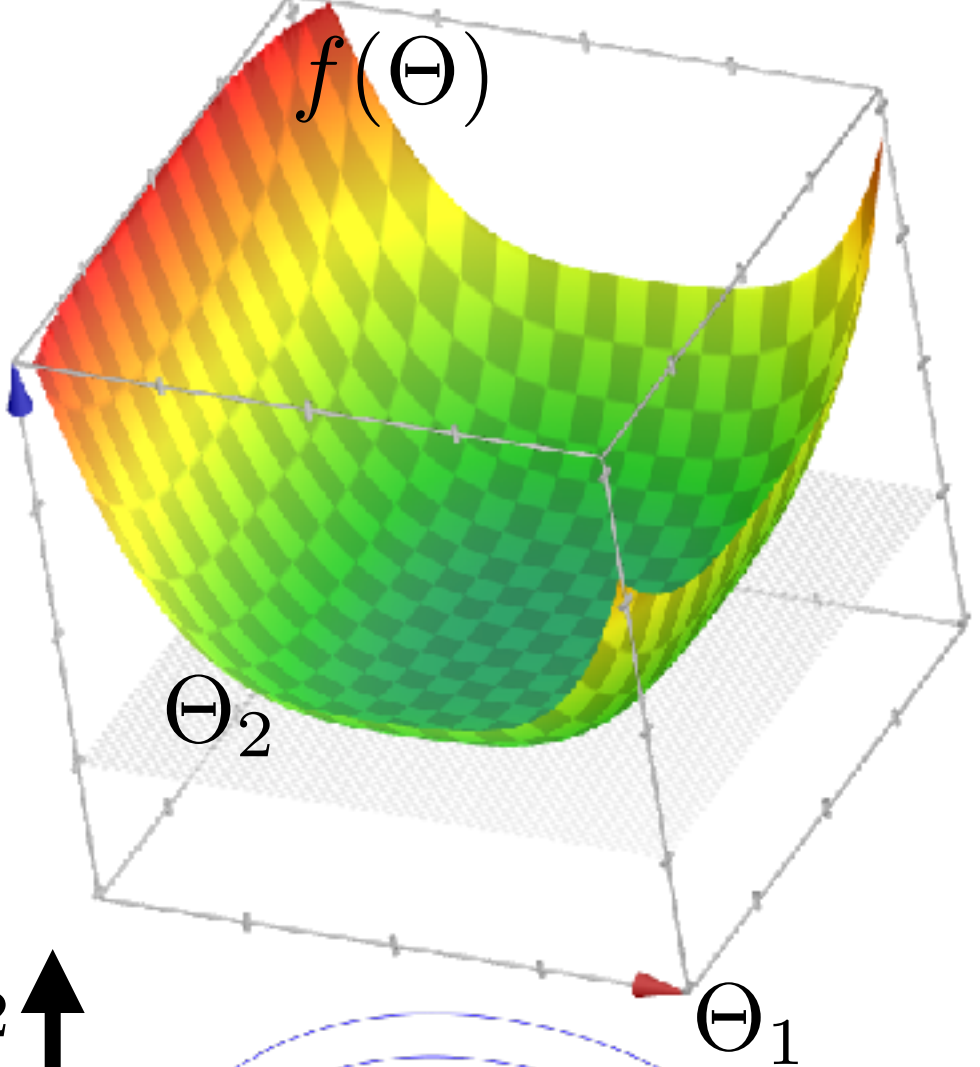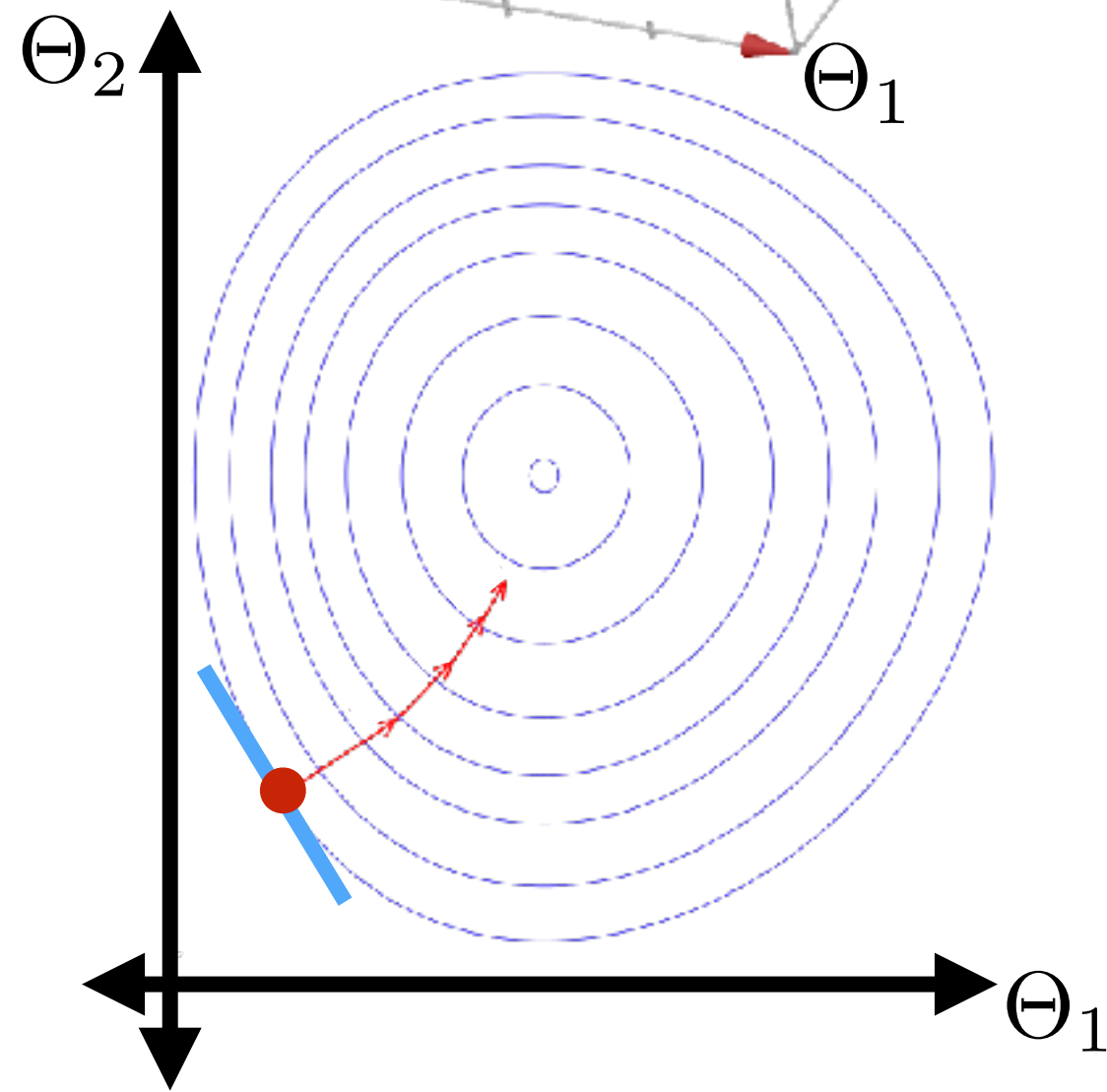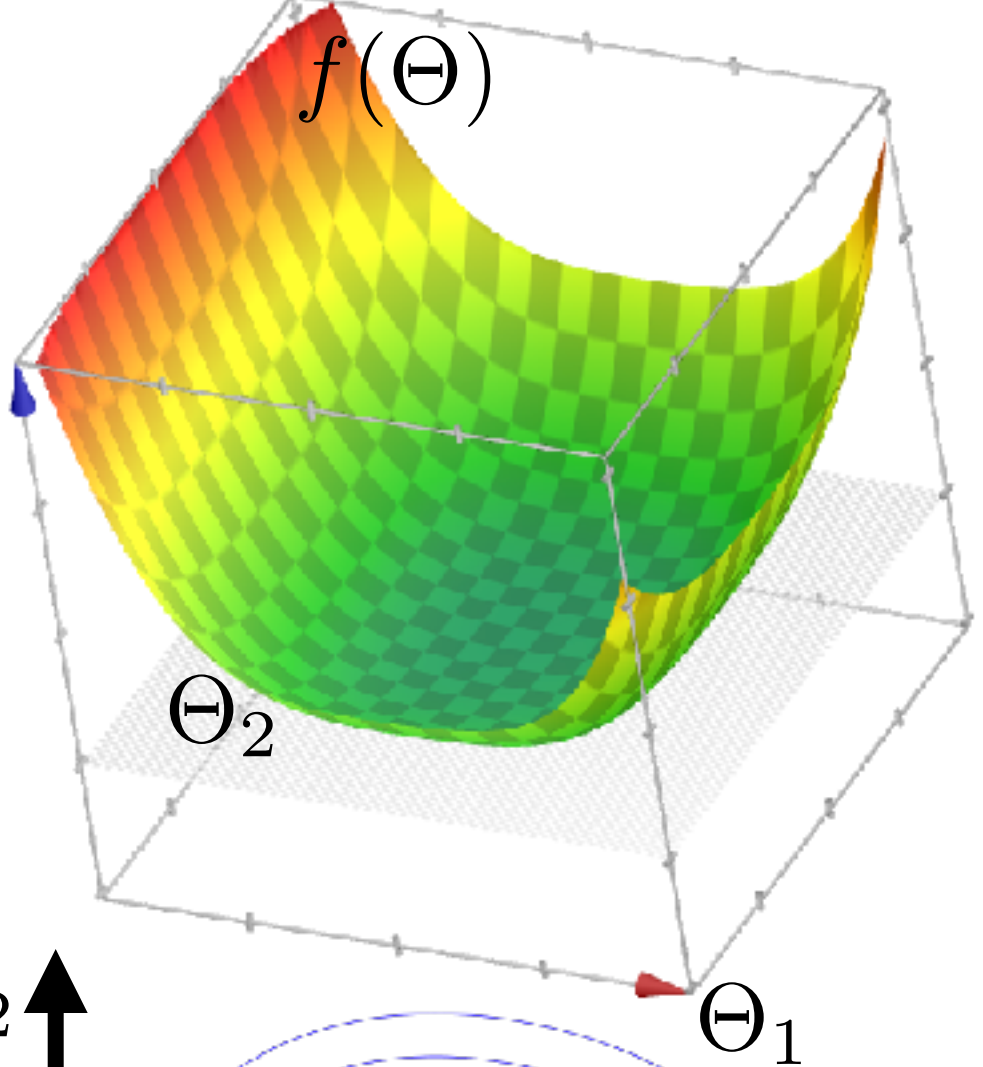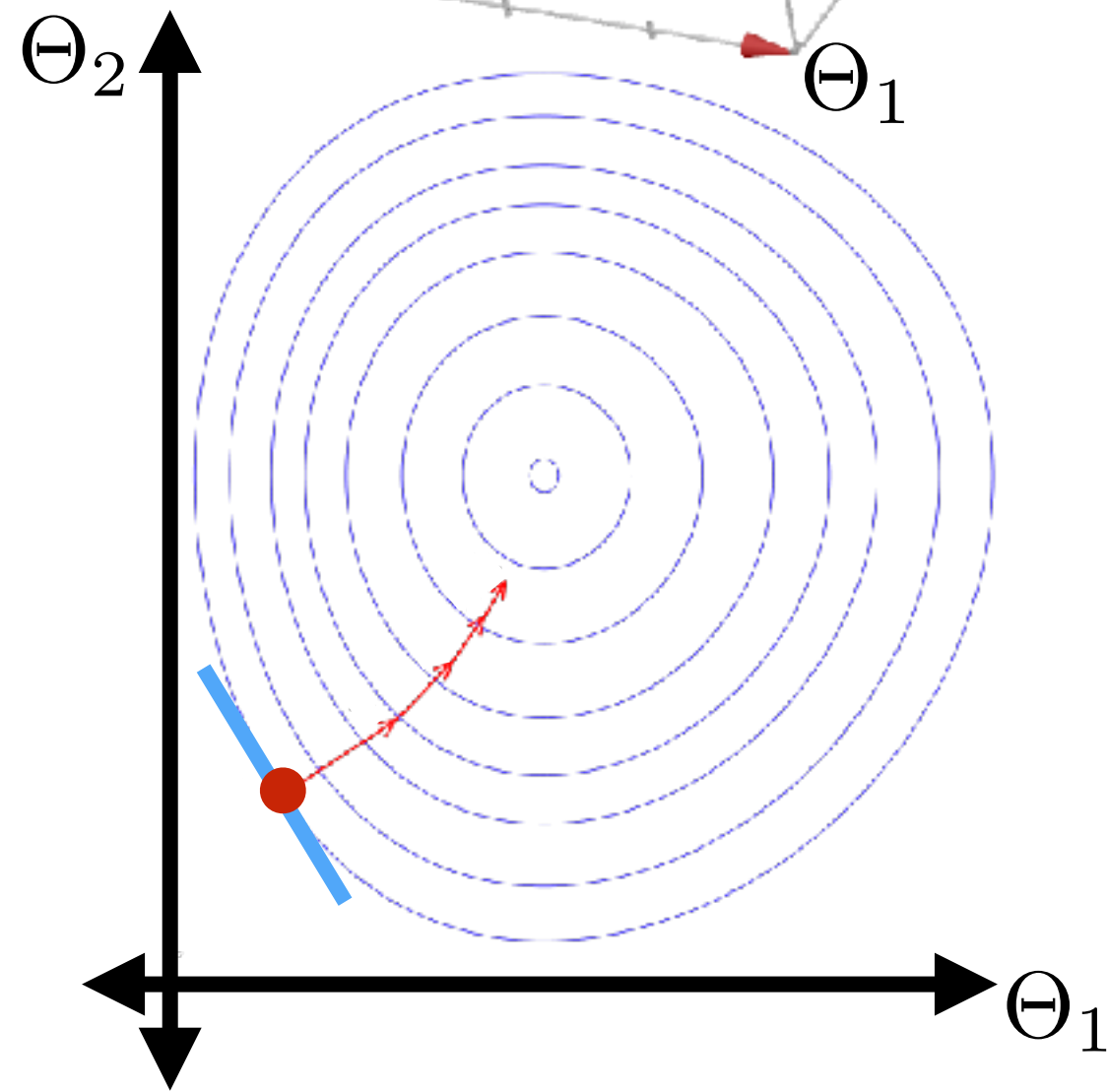**repeat**

  t = t + 1

  $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\text{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

$\text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\text{Initialize t = 0}$
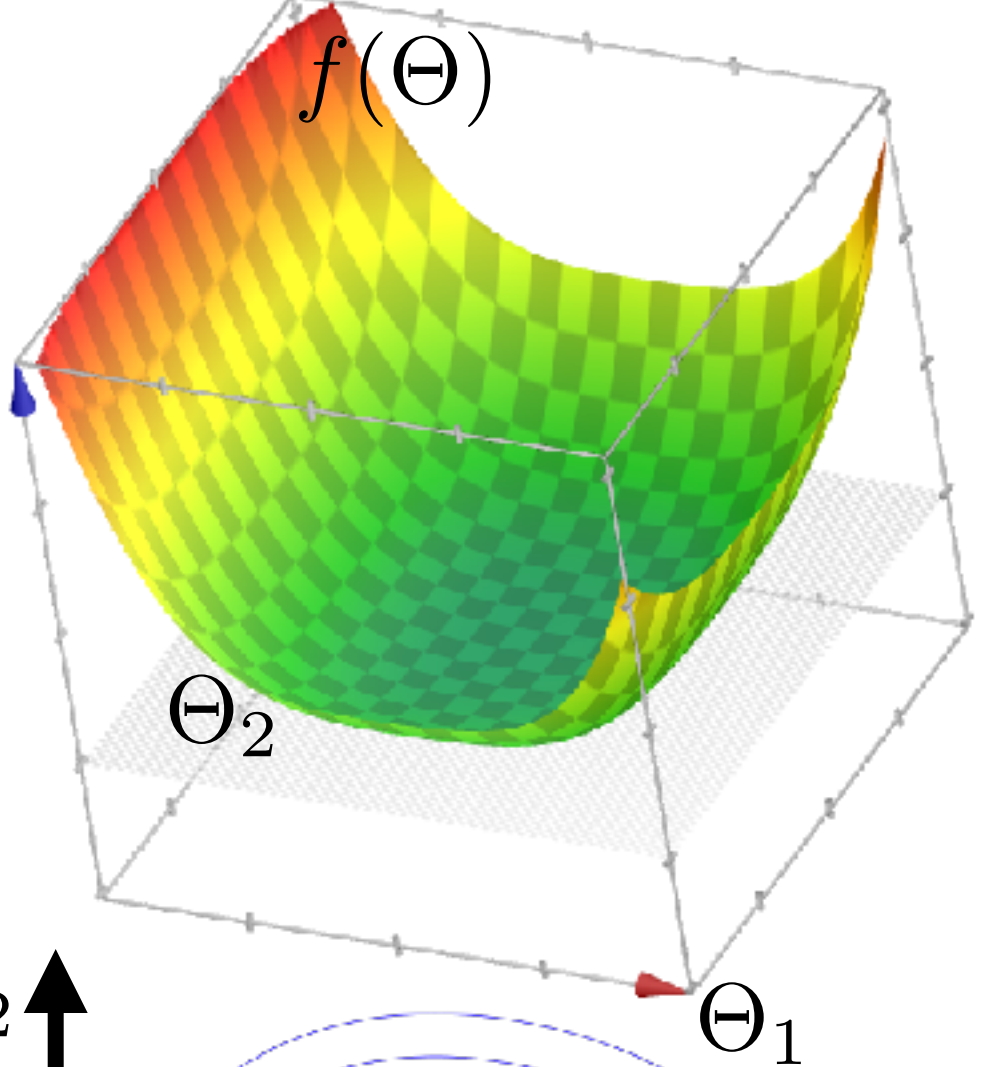
**repeat**

   $\text{t = t + 1}$

   $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[\dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m}\right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\text{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

$\text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\text{Initialize t = 0}$

**repeat**

$\quad \text{t = t + 1}$

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

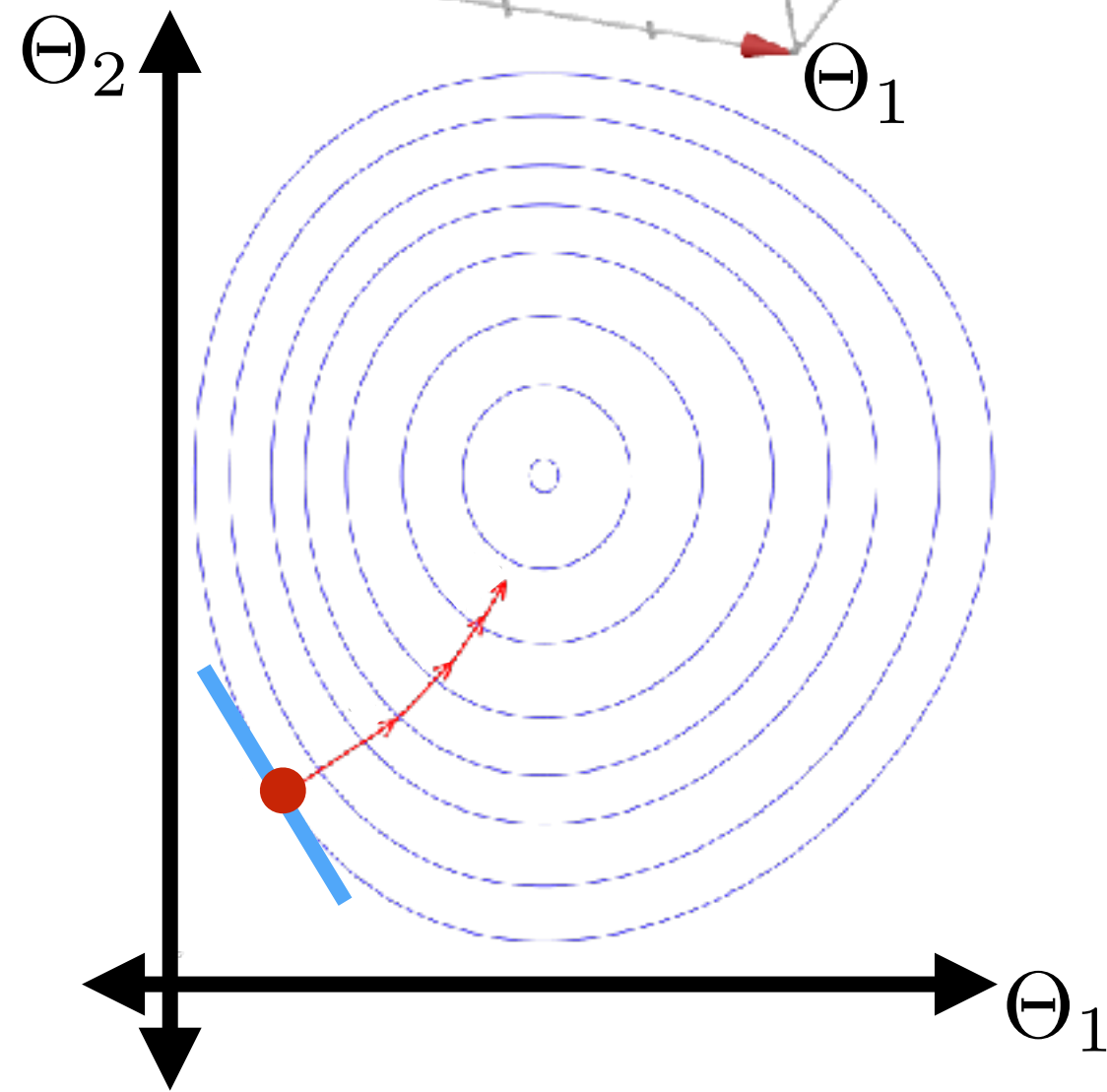**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\texttt{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

$\texttt{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\texttt{Initialize t = 0}$

**repeat**

    $\texttt{t = t + 1}$

    $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

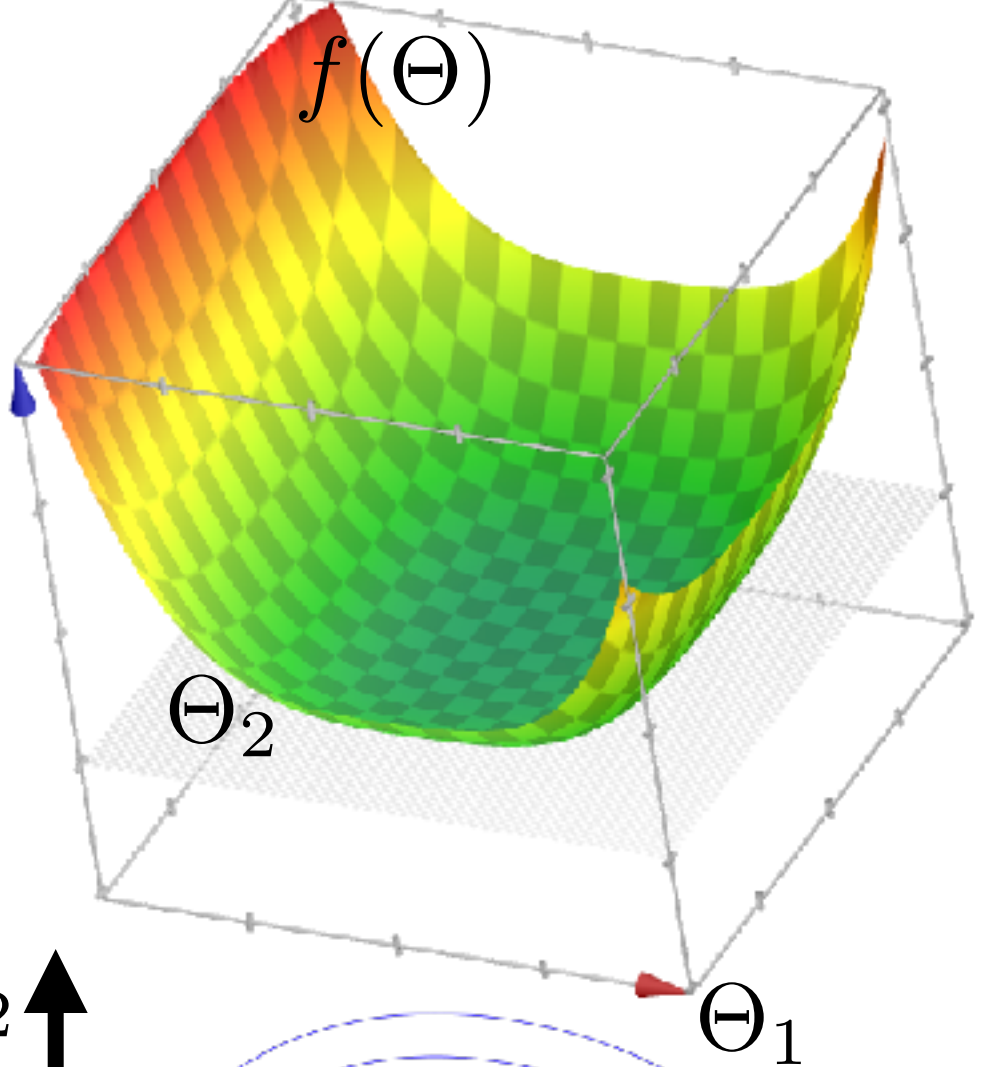**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$
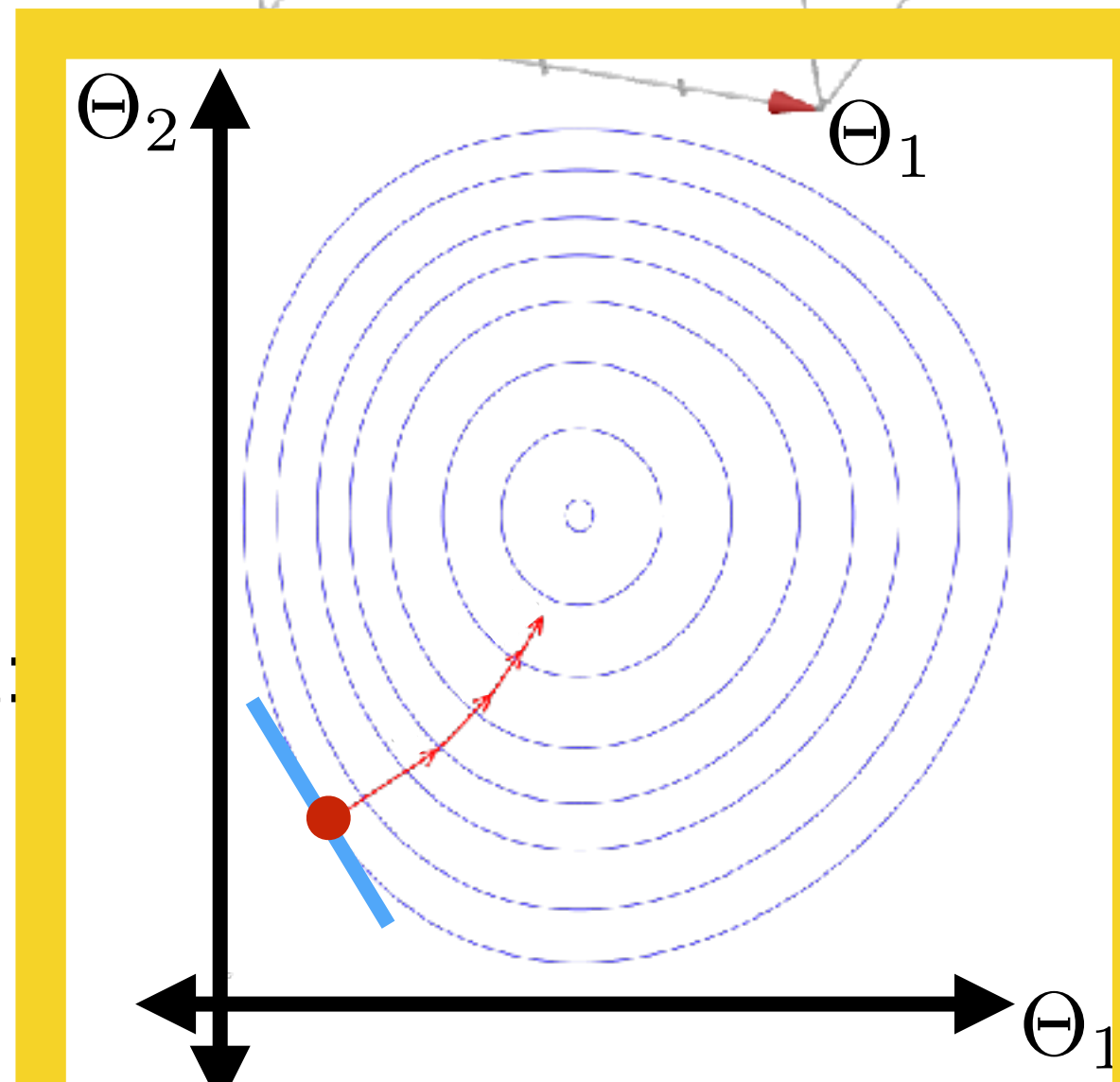
**Return** $\Theta^{(t)}$



3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\text{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

$\text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\text{Initialize } \texttt{t = 0}$

**repeat**

$\qquad \texttt{t = t + 1}$

$\qquad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

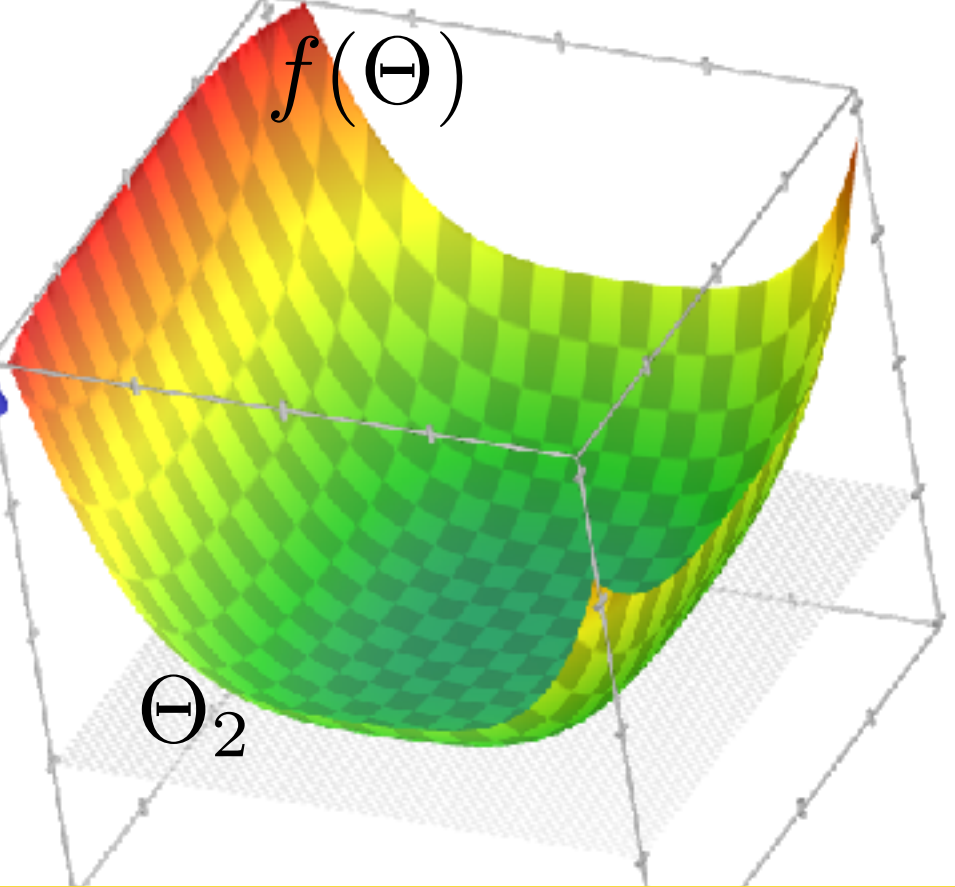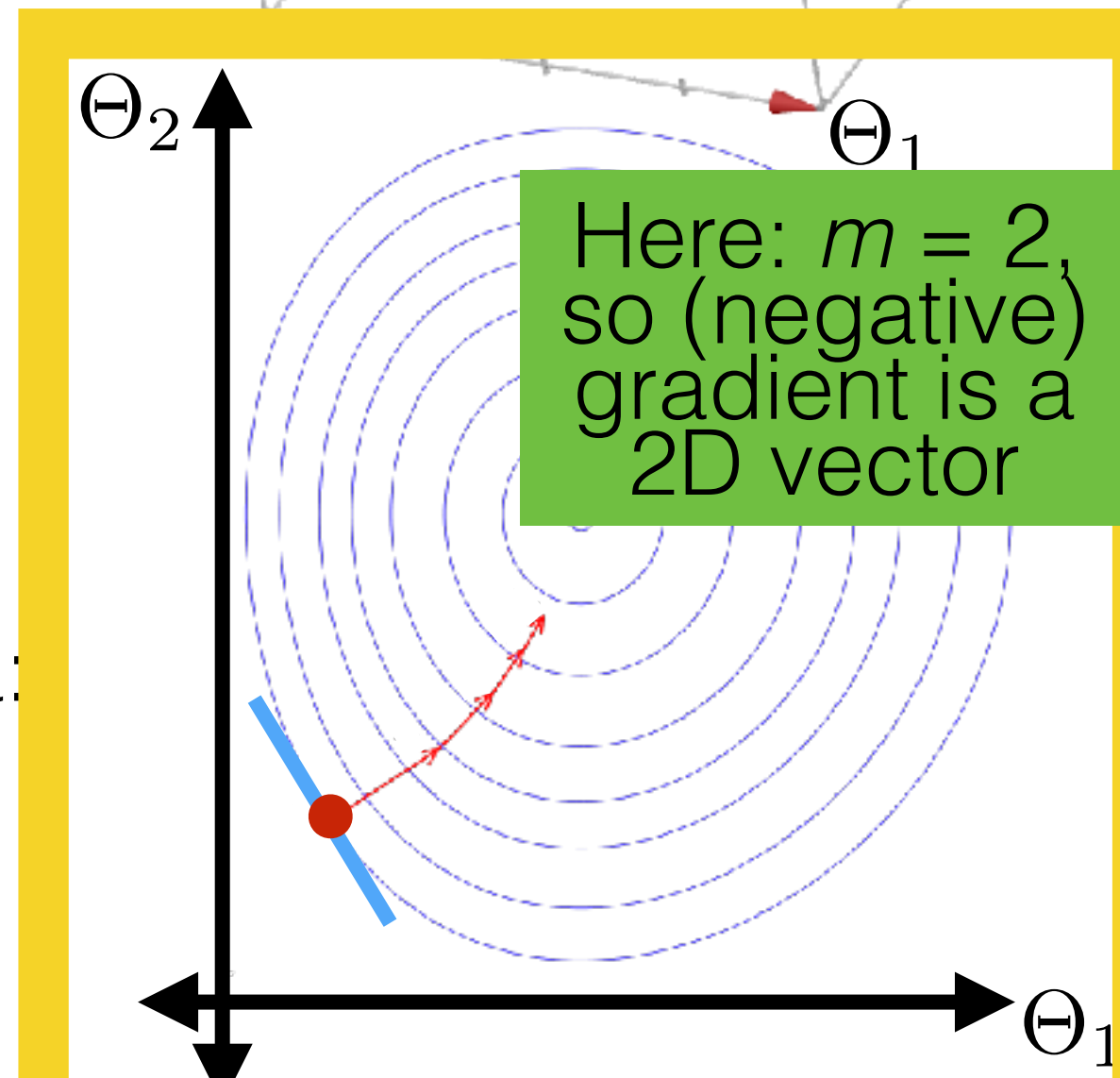**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

- Other possible stopping criteria:

# Gradient descent

- Gradient $\nabla_\Theta f = \left[\dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m}\right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\text{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

$\text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\text{Initialize t = 0}$

**repeat**

    $\text{t = t + 1}$

    $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

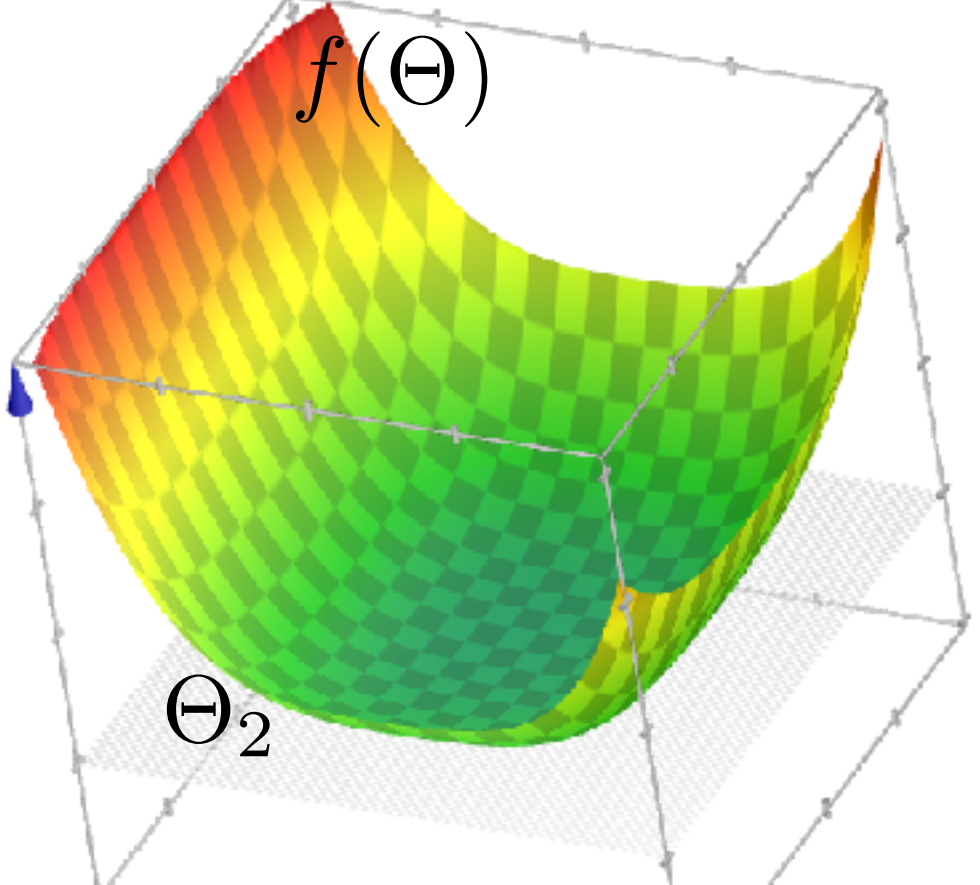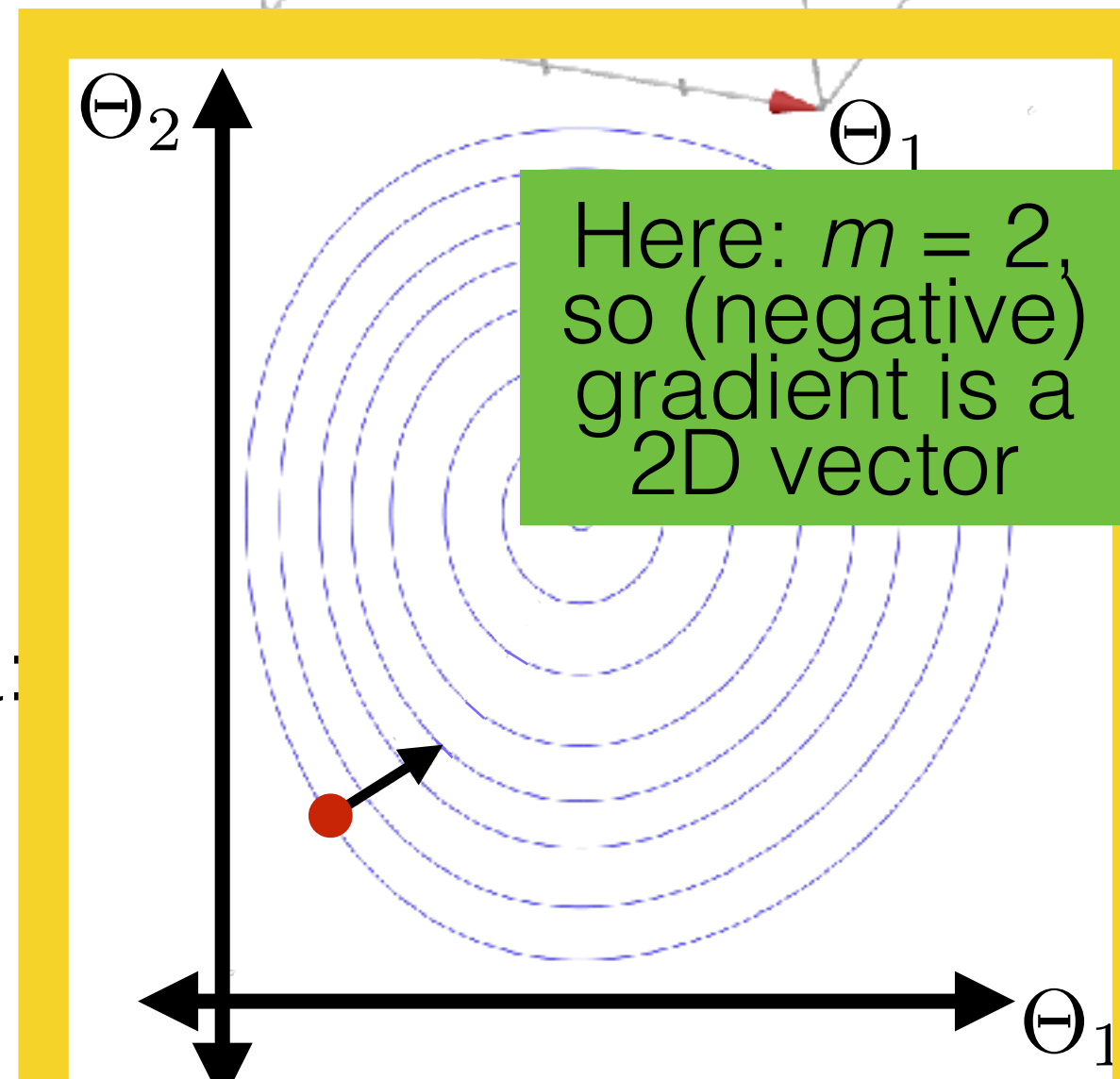**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

- Other possible stopping criteria:
  - Max number of iterations $T$

# Gradient descent



- Gradient $\nabla_\Theta f = \left[\dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m}\right]^\top$
  - with $\Theta \in \mathbb{R}^m$

```
Gradient-Descent (Θ_init, η, f, ∇_Θ f, ε)
  Initialize Θ^(0) = Θ_init
  Initialize t = 0
```

**repeat**

$\quad$ `t = t + 1`

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

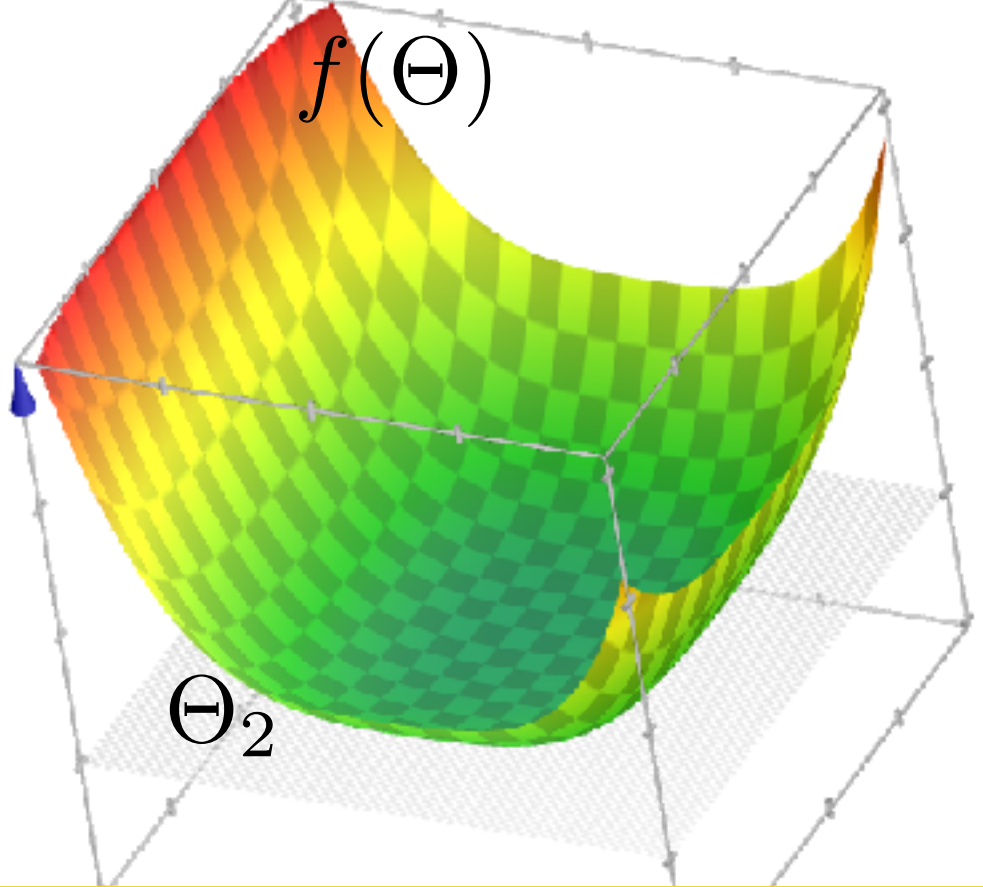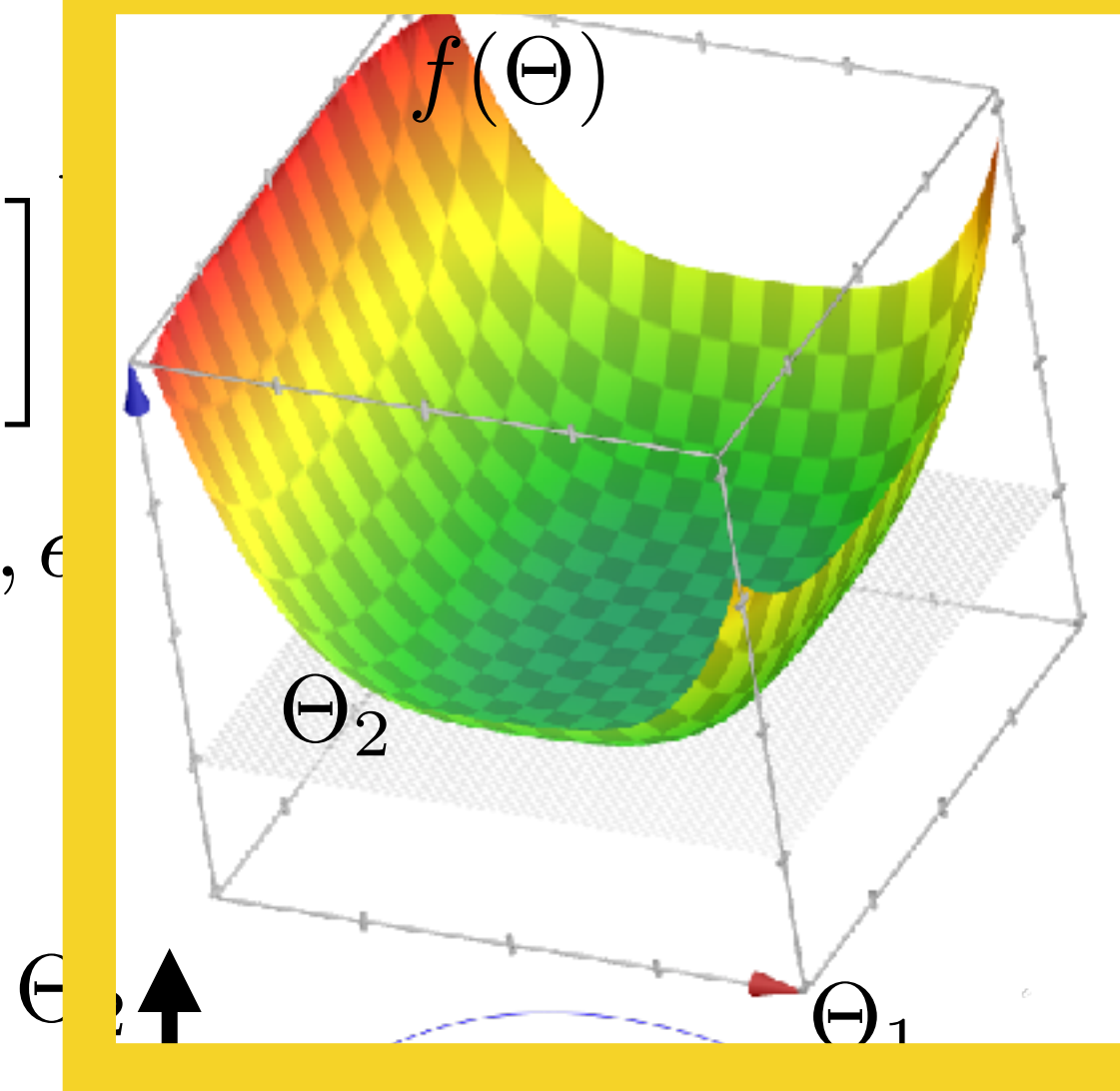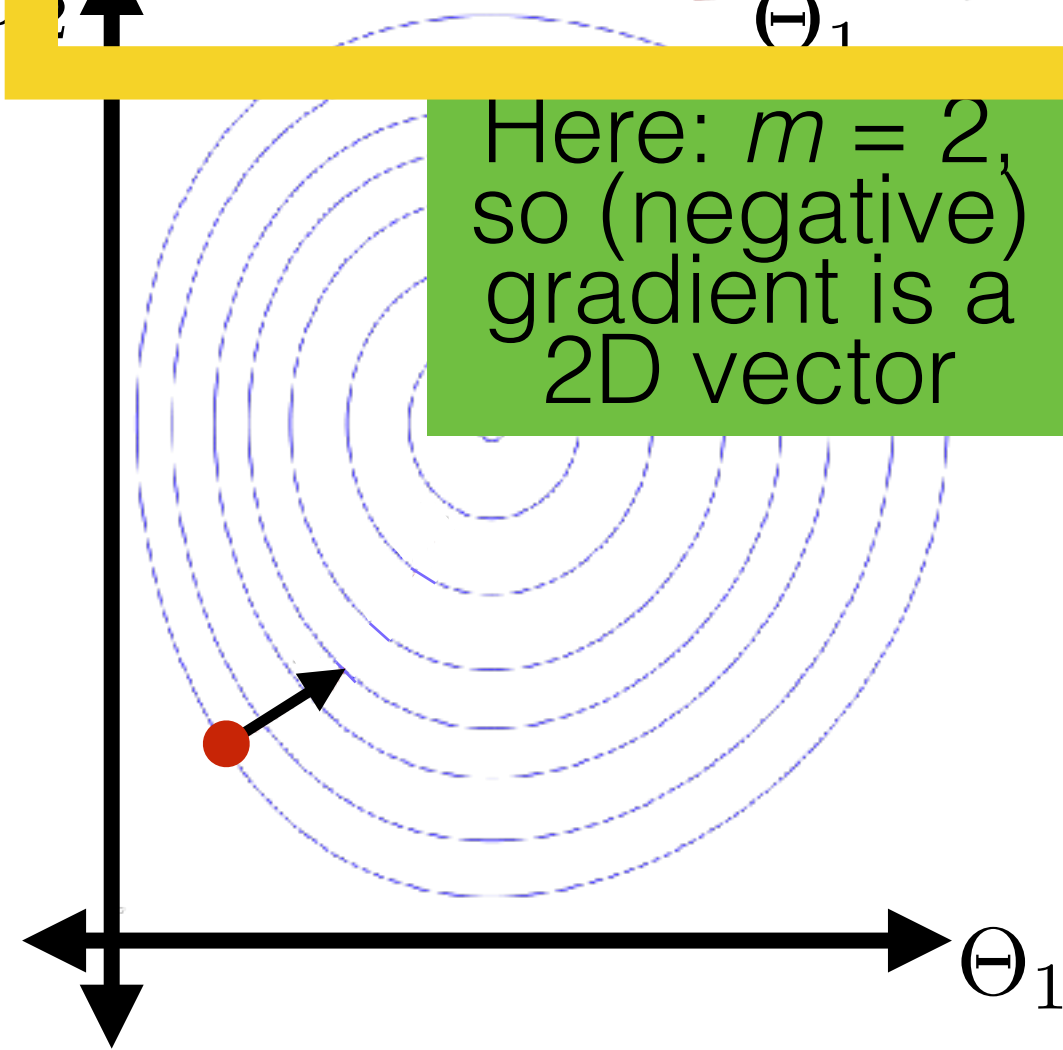**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

- Other possible stopping criteria:
  - Max number of iterations $T$
  - $\|\Theta^{(t)} - \Theta^{(t-1)}\| < \epsilon$



3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$\texttt{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon)$

$\texttt{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\texttt{Initialize t = 0}$

**repeat**

$\quad \texttt{t = t + 1}$

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$
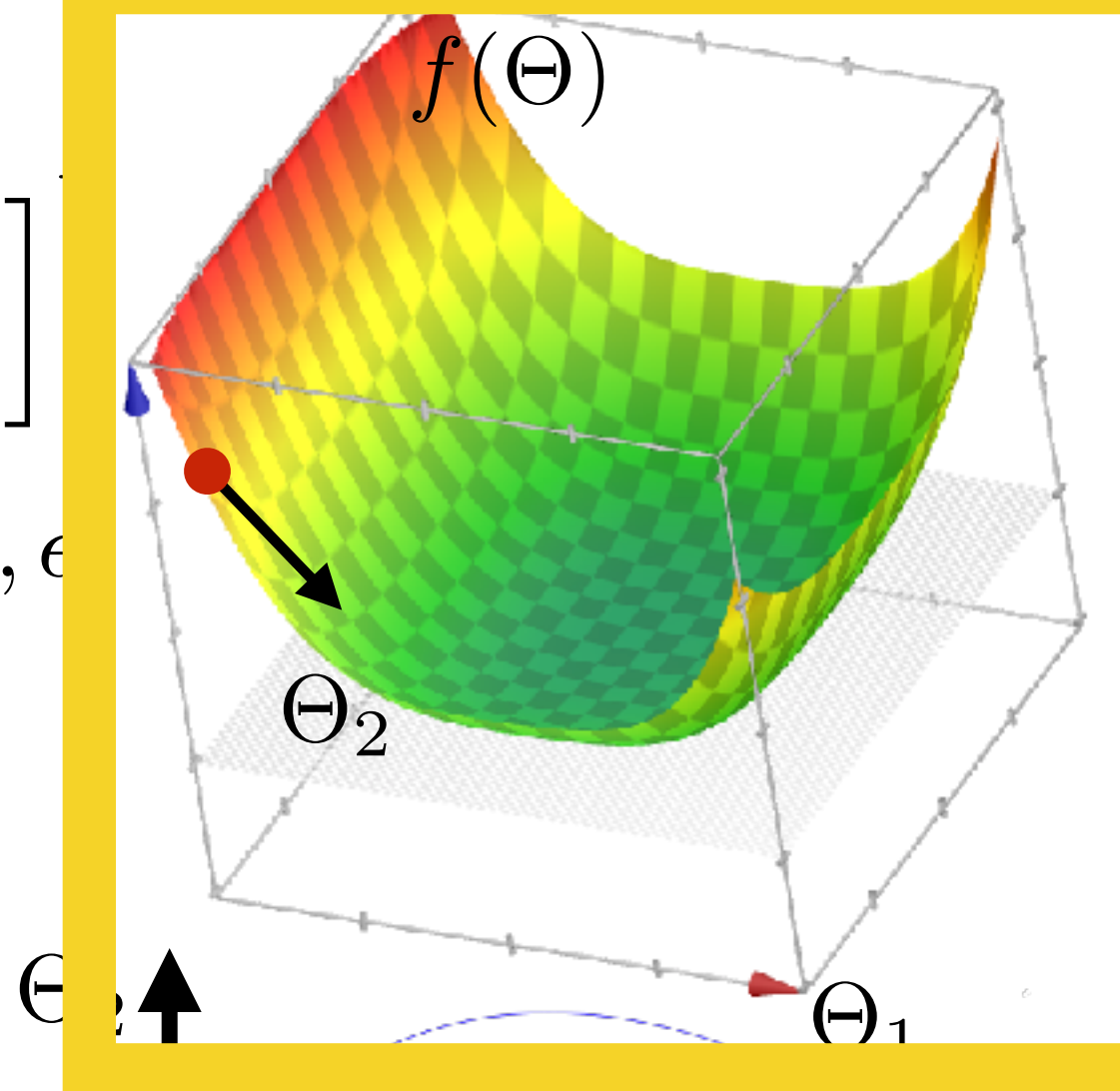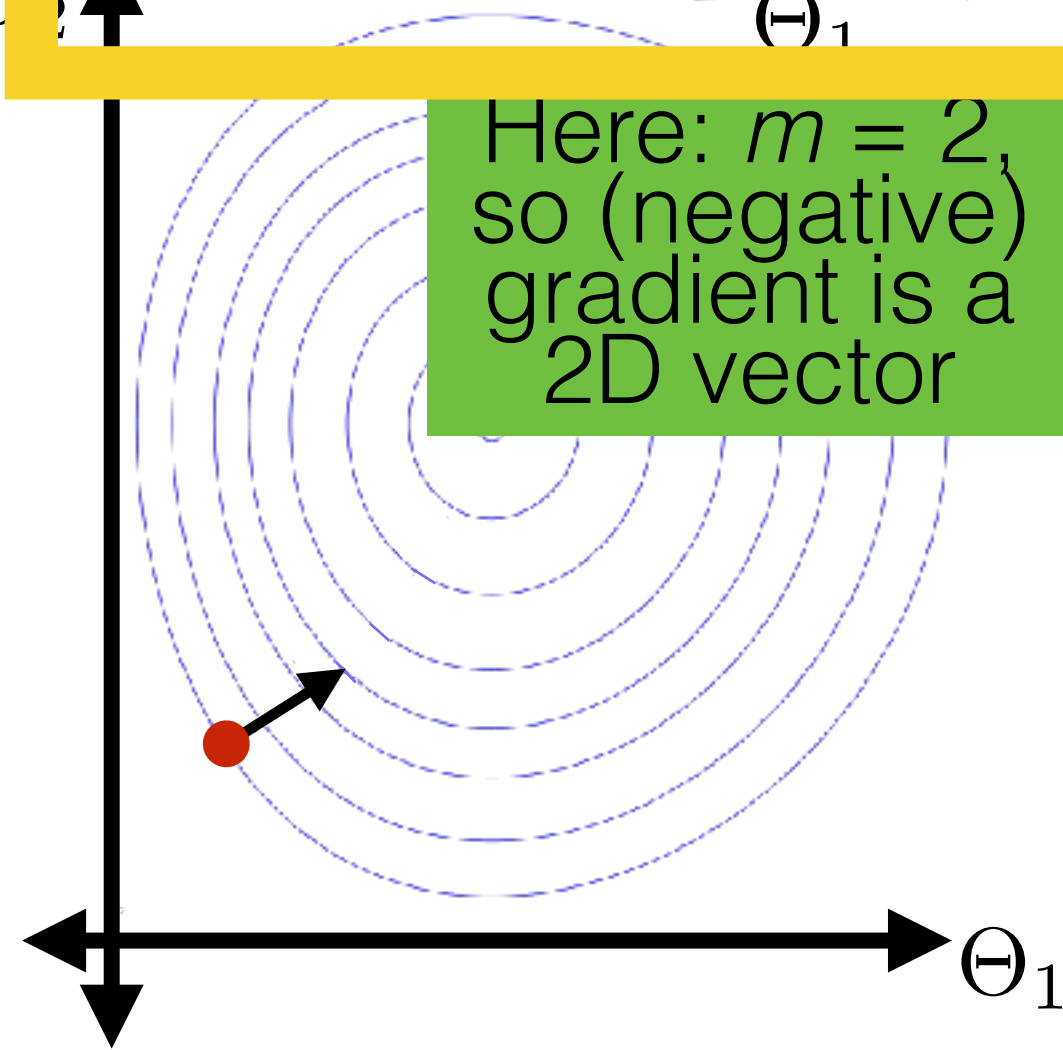
**Return** $\Theta^{(t)}$

- Other possible stopping criteria:
  - Max number of iterations $T$
  - $\|\Theta^{(t)} - \Theta^{(t-1)}\| < \epsilon$
  - $\|\nabla_\Theta f(\Theta^{(t)})\| < \epsilon$

3

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \dots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

```
Gradient-Descent (Θ_init, η, f, ∇_Θf, ε)
  Initialize Θ^(0) = Θ_init
  Initialize t = 0
```
**repeat**
    `t = t + 1`

$$\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$$

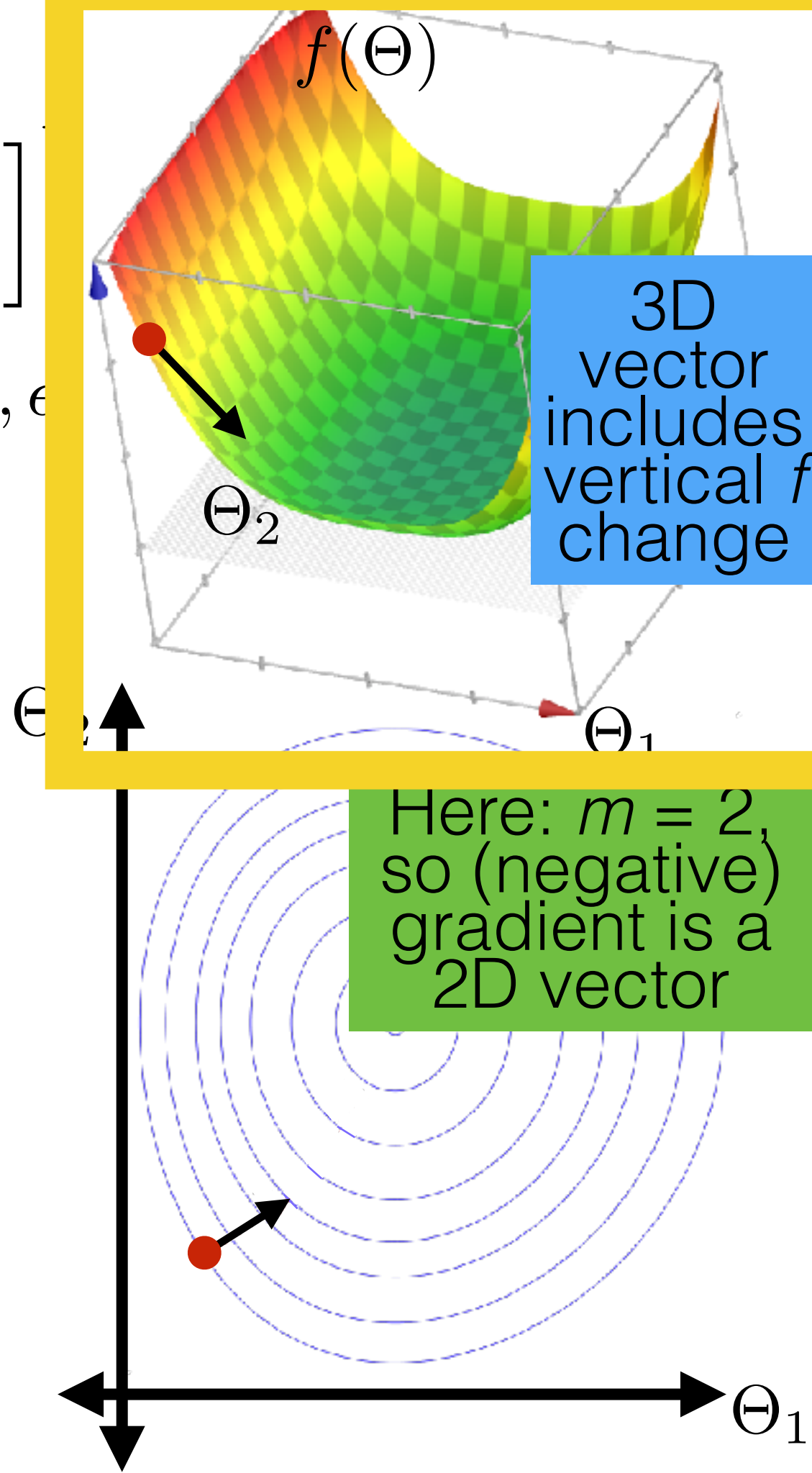**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

- Other possible stopping criteria:
  - Max number of iterations $T$
  - $\| \Theta^{(t)} - \Theta^{(t-1)} \| < \epsilon$
  - $\| \nabla_\Theta f(\Theta^{(t)}) \| < \epsilon$

4

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

```
Gradient-Descent (Θ_init, η, f, ∇_Θ f, ε)
  Initialize Θ^(0) = Θ_init
  Initialize t = 0
```
**repeat**
```
    t = t + 1
```
$$\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$$
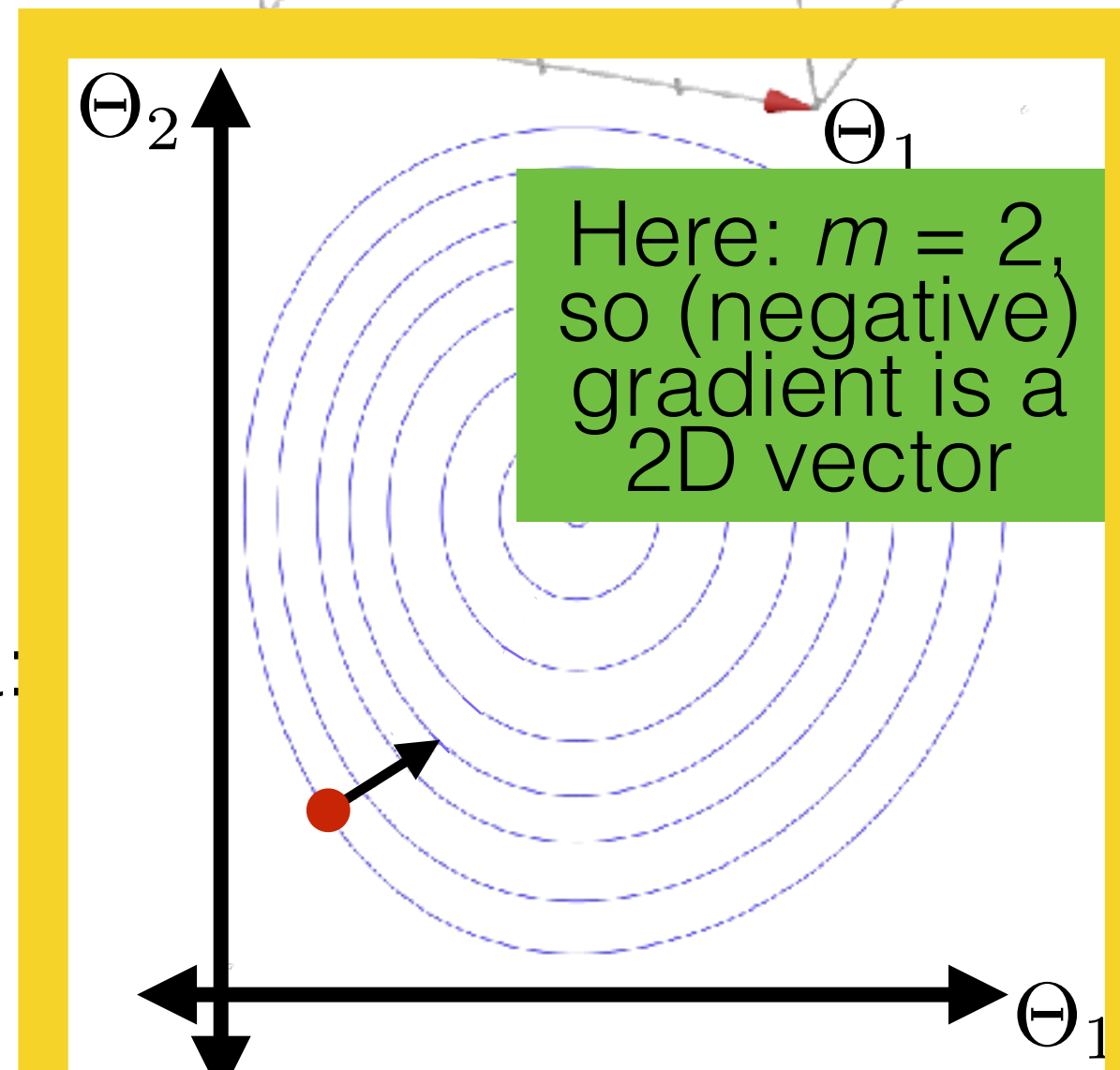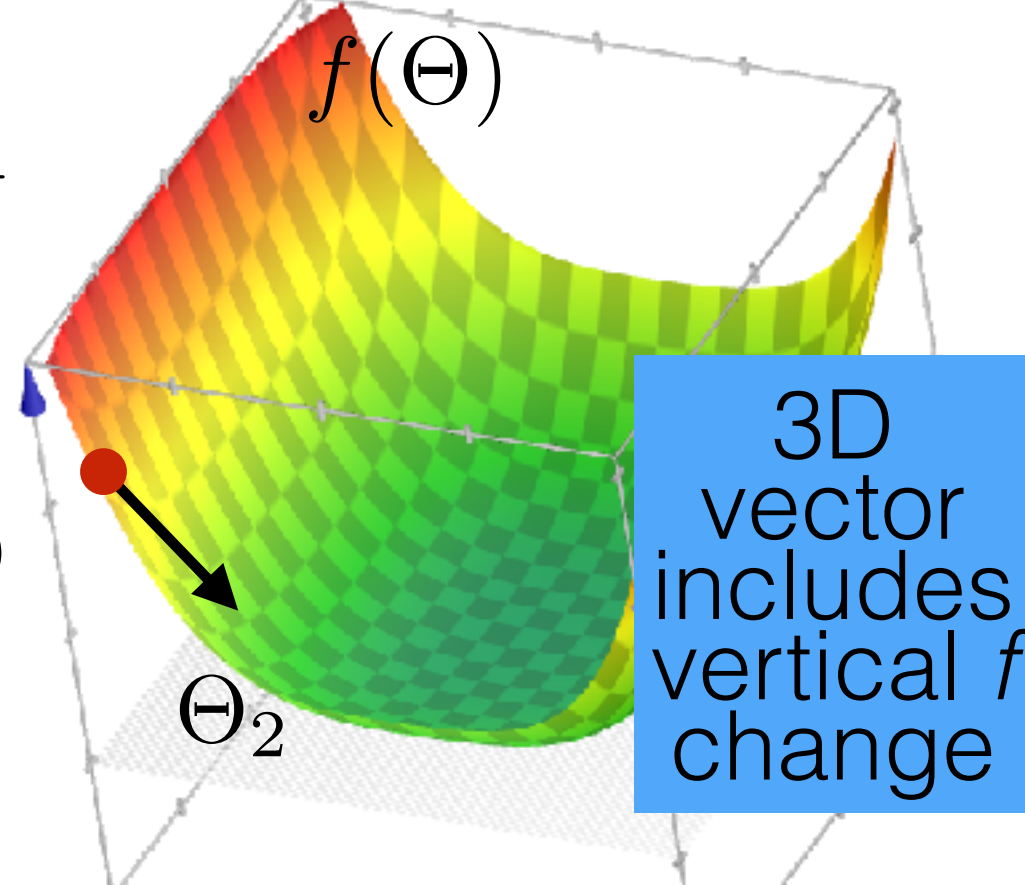**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$
**Return** $\Theta^{(t)}$

- Other possible stopping criteria:
  - Max number of iterations $T$
  - $\| \Theta^{(t)} - \Theta^{(t-1)} \| < \epsilon$
  - $\| \nabla_\Theta f(\Theta^{(t)}) \| < \epsilon$



$f(\Theta)$

$\Theta_2$

$\Theta_1$



$\Theta_2$

$\Theta_1$

Here: $m = 2$, so (negative) gradient is a 2D vector

4

# Gradient descent



- Gradient $\nabla_\Theta f = \left[\dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m}\right]^\top$
  - with $\Theta \in \mathbb{R}^m$

```
Gradient-Descent (Θ_init, η, f, ∇_Θ f, ε)
  Initialize Θ^(0) = Θ_init
  Initialize t = 0
```
**repeat**
```
    t = t + 1
```
$$\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$$
**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$
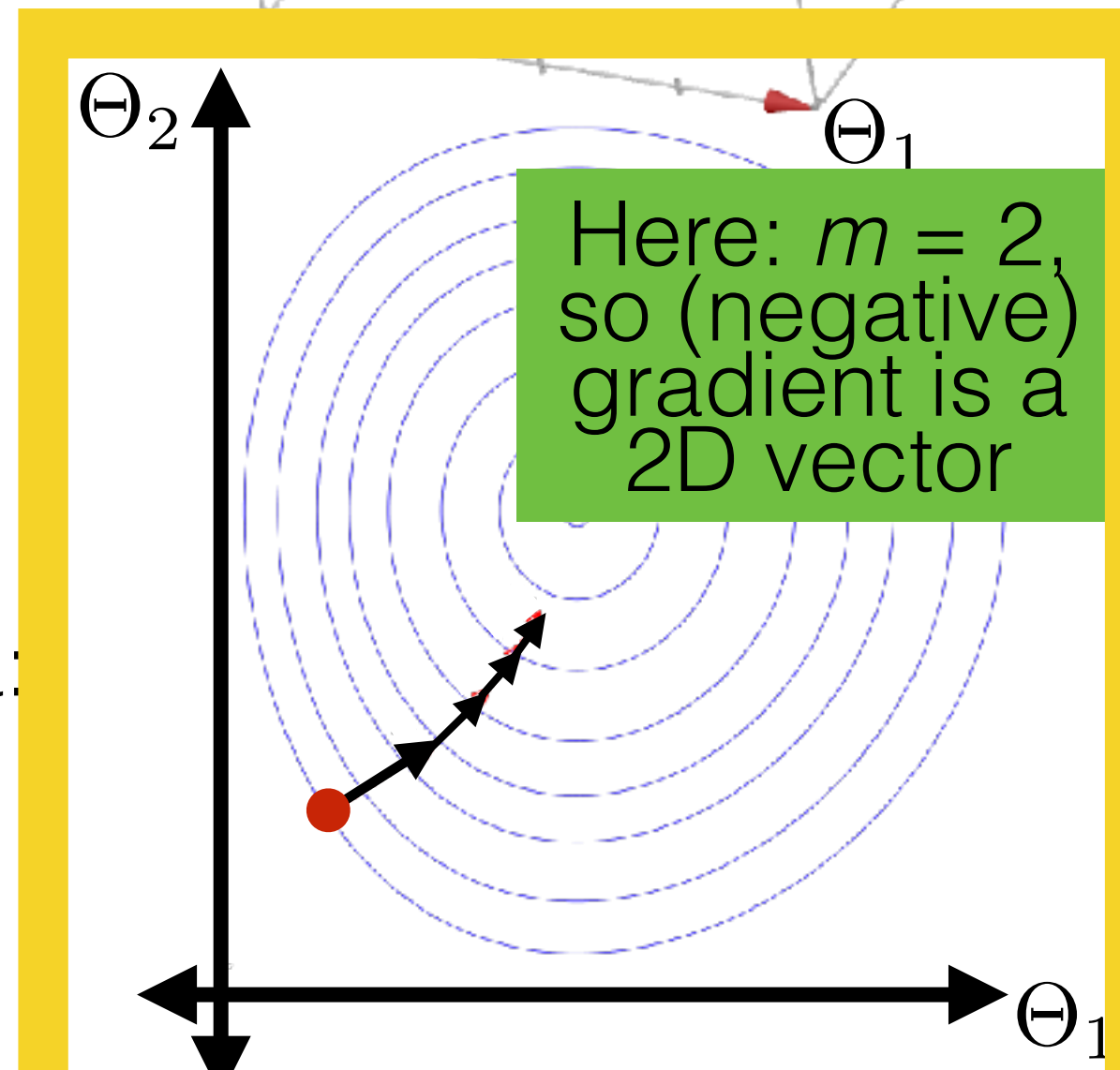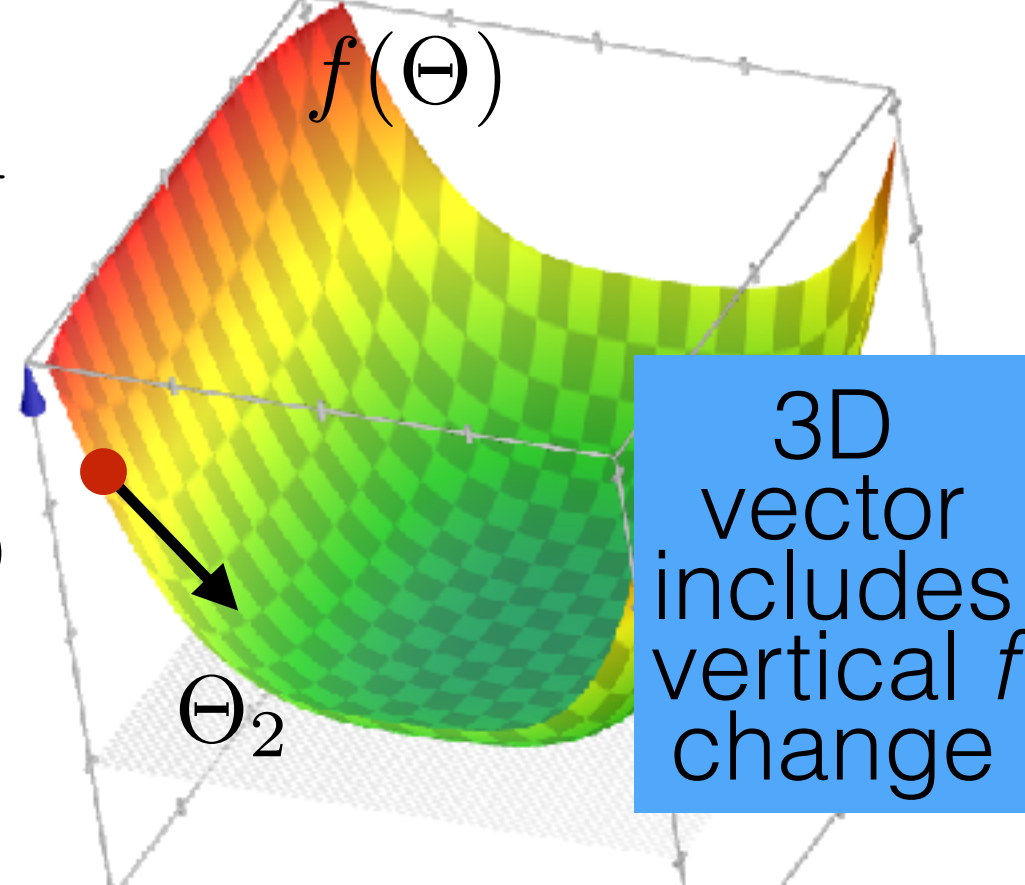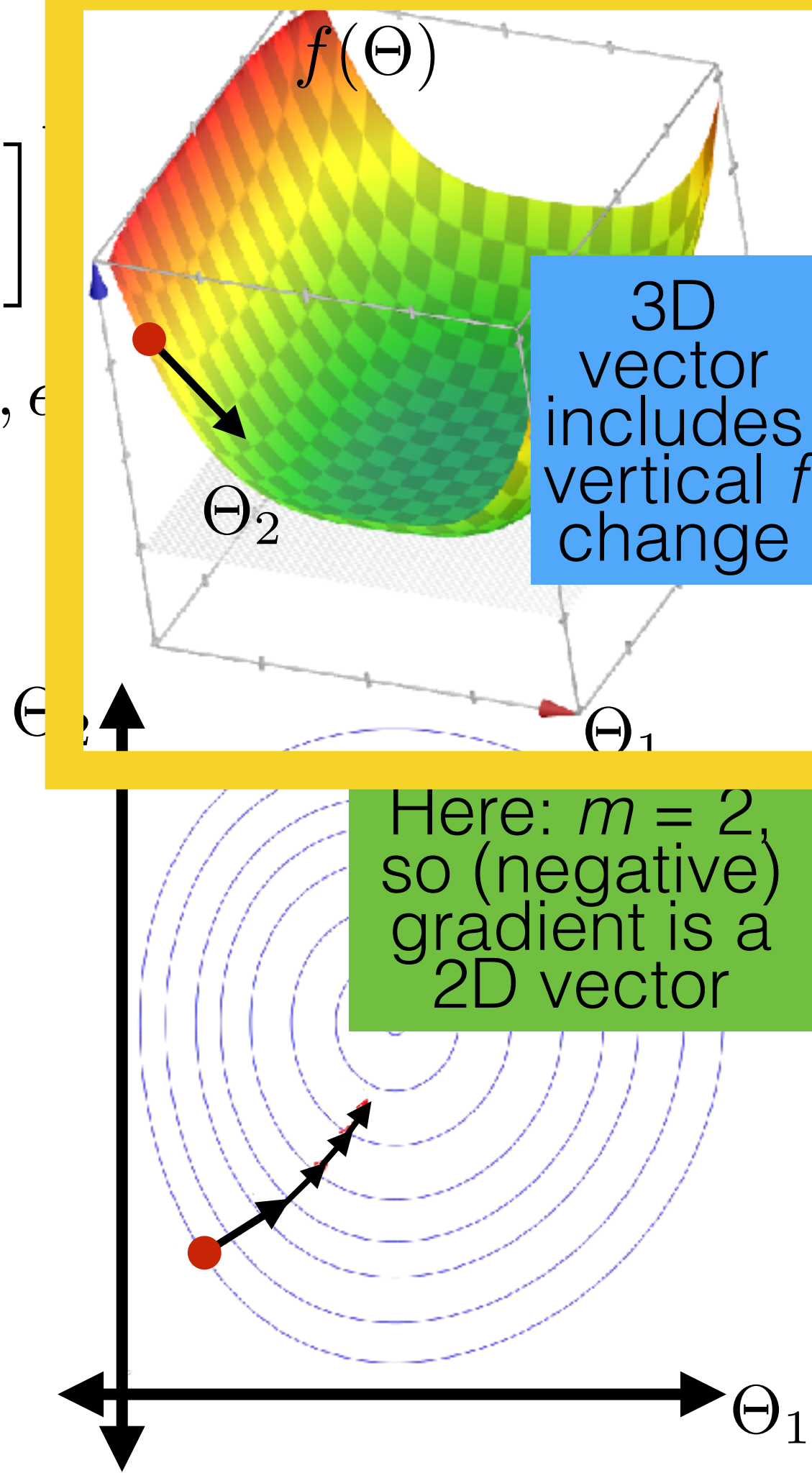**Return** $\Theta^{(t)}$

- Other possible stopping criteria:
  - Max number of iterations $T$
  - $\|\Theta^{(t)} - \Theta^{(t-1)}\| < \epsilon$
  - $\|\nabla_\Theta f(\Theta^{(t)})\| < \epsilon$



Here: $m = 2$, so (negative) gradient is a 2D vector

4

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]$
  - with $\Theta \in \mathbb{R}^m$

$\text{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon$

$\quad \text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\quad \text{Initialize } \texttt{t = 0}$
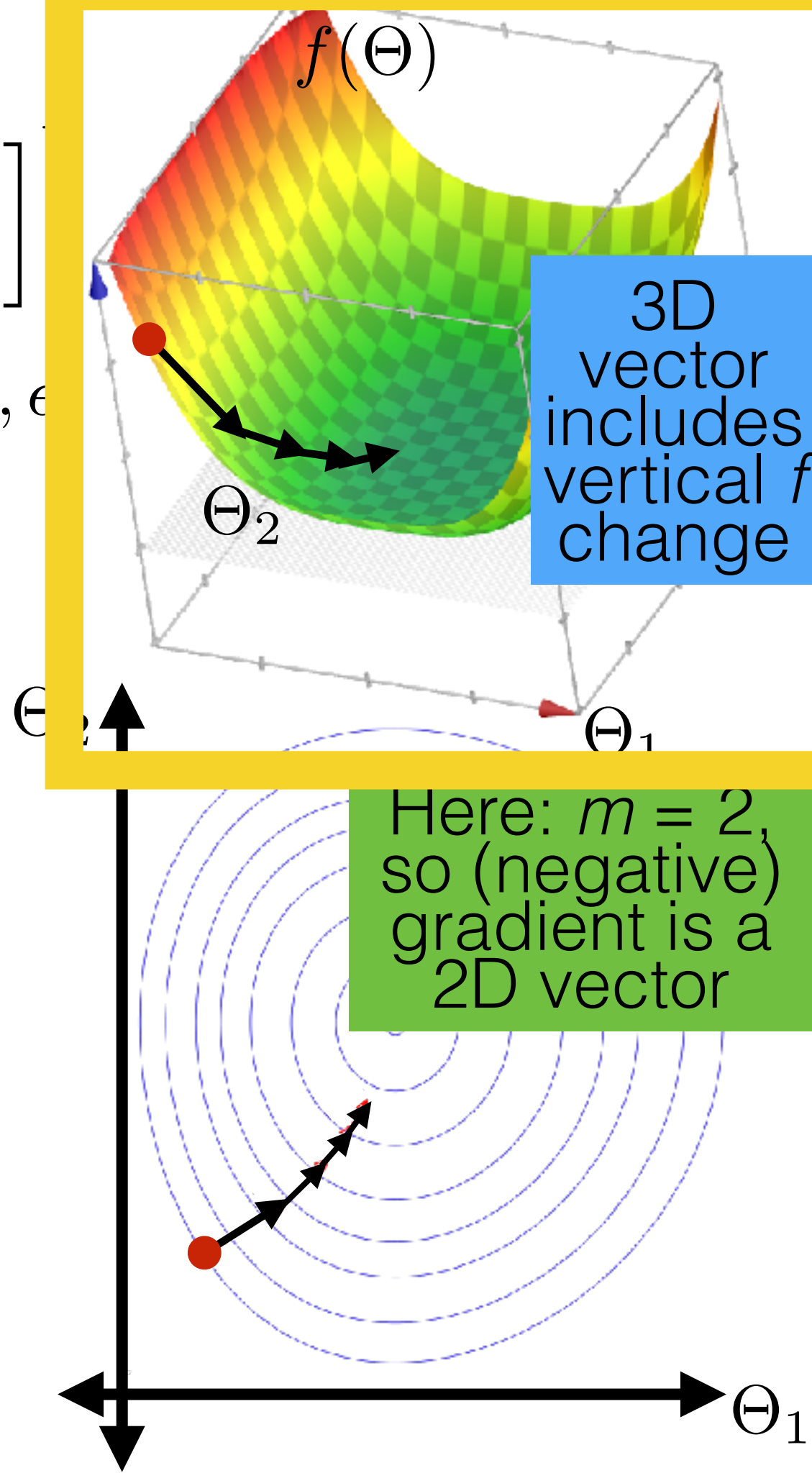
**repeat**

$\quad \texttt{t = t + 1}$

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

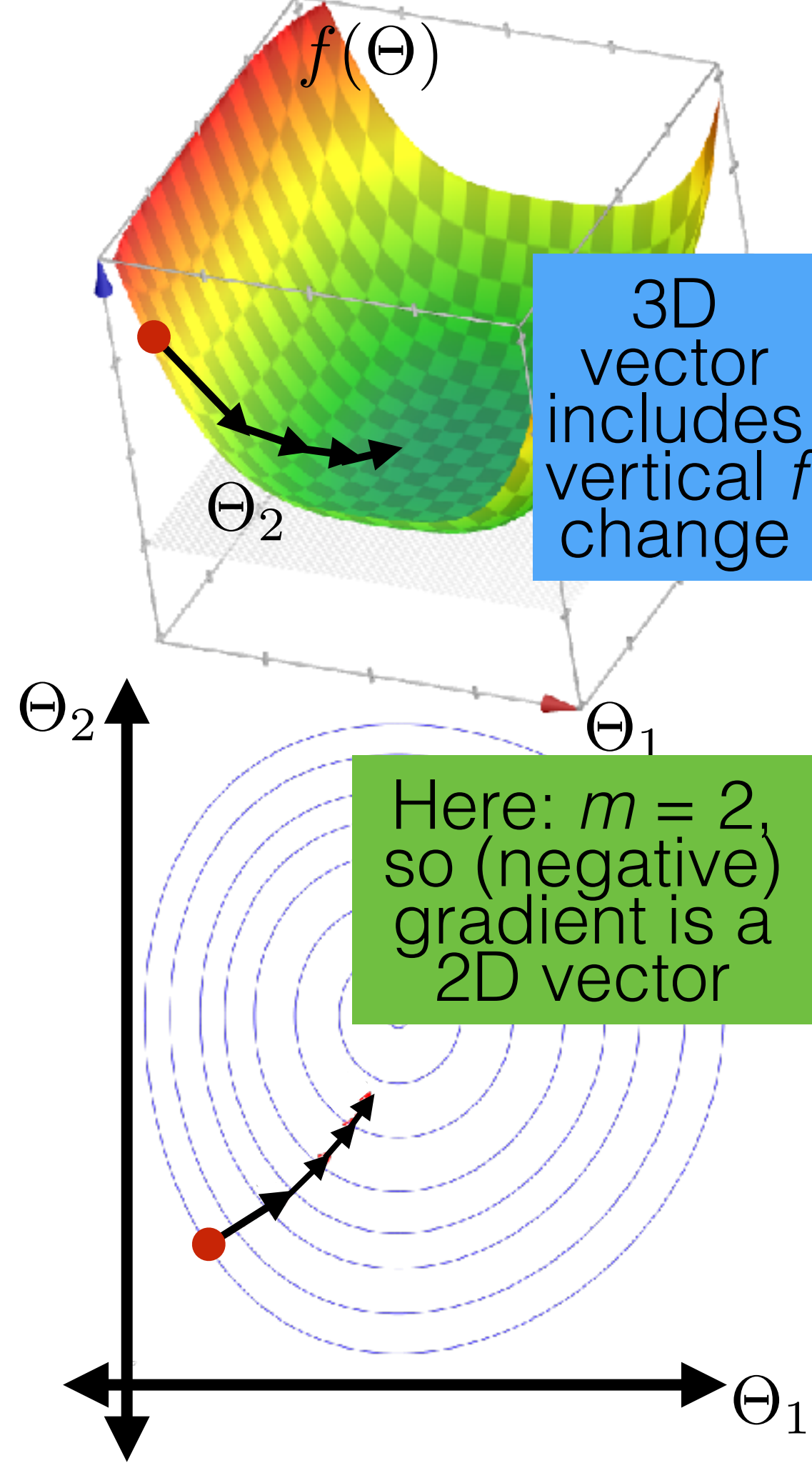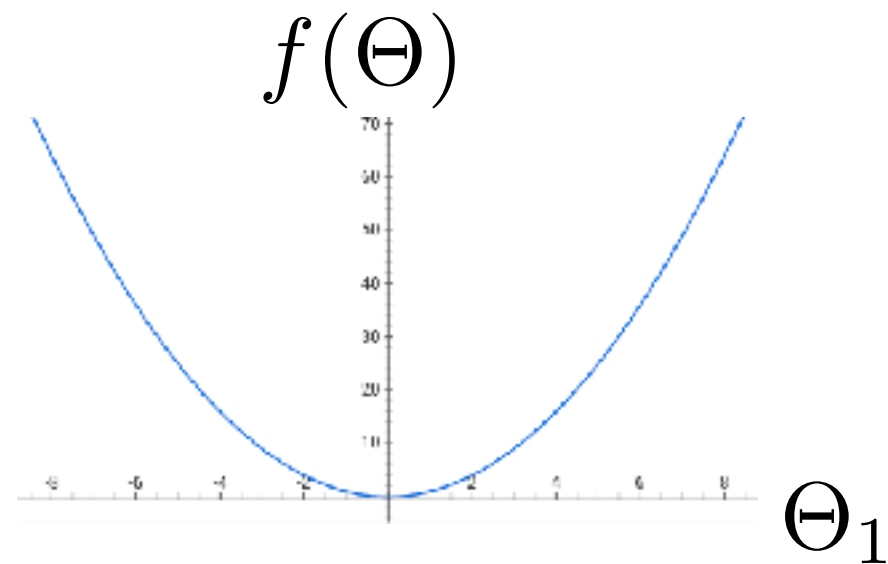**Return** $\Theta^{(t)}$

- Other possible stopping criteria:
  - Max number of iterations $T$
  - $\| \Theta^{(t)} - \Theta^{(t-1)} \| < \epsilon$
  - $\| \nabla_\Theta f(\Theta^{(t)}) \| < \epsilon$

4



$f(\Theta)$

$\Theta_2$

$\Theta_2$     $\Theta_1$

Here: $m = 2$, so (negative) gradient is a 2D vector

$\Theta_1$

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]$
  - with $\Theta \in \mathbb{R}^m$

```
Gradient-Descent (Θinit, η, f, ∇Θf, ε
  Initialize Θ⁽⁰⁾ = Θinit
  Initialize t = 0
```
**repeat**
  ```
  t = t + 1
  ```
  $$\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$$
**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$
**Return** $\Theta^{(t)}$

- Other possible stopping criteria:
  - Max number of iterations $T$
  - $\|\Theta^{(t)} - \Theta^{(t-1)}\| < \epsilon$
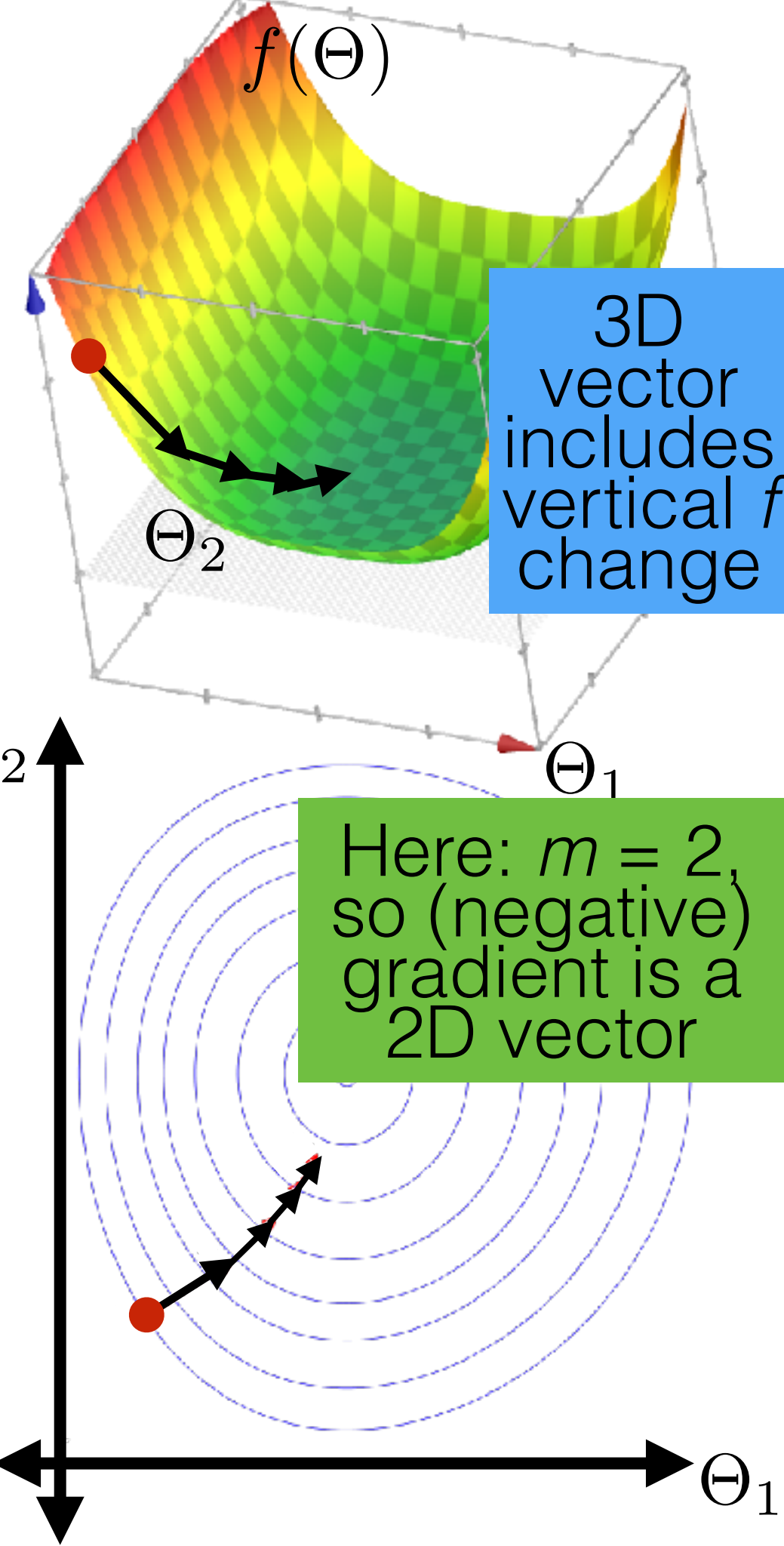  - $\|\nabla_\Theta f(\Theta^{(t)})\| < \epsilon$

4



$f(\Theta)$

$\Theta_2$

$\Theta_1$

Here: $m = 2$, so (negative) gradient is a 2D vector

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]$
  - with $\Theta \in \mathbb{R}^m$

$\texttt{Gradient-Descent}(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon$

$\quad \texttt{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\quad \texttt{Initialize t = 0}$

**repeat**

$\quad \texttt{t = t + 1}$

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

- Other possible stopping criteria:
  - Max number of iterations $T$
  - $\| \Theta^{(t)} - \Theta^{(t-1)} \| < \epsilon$
  - $\| \nabla_\Theta f(\Theta^{(t)}) \| < \epsilon$

4



$f(\Theta)$

$\Theta_2$

$\Theta_2$

$\Theta_1$

3D vector includes vertical $f$ change
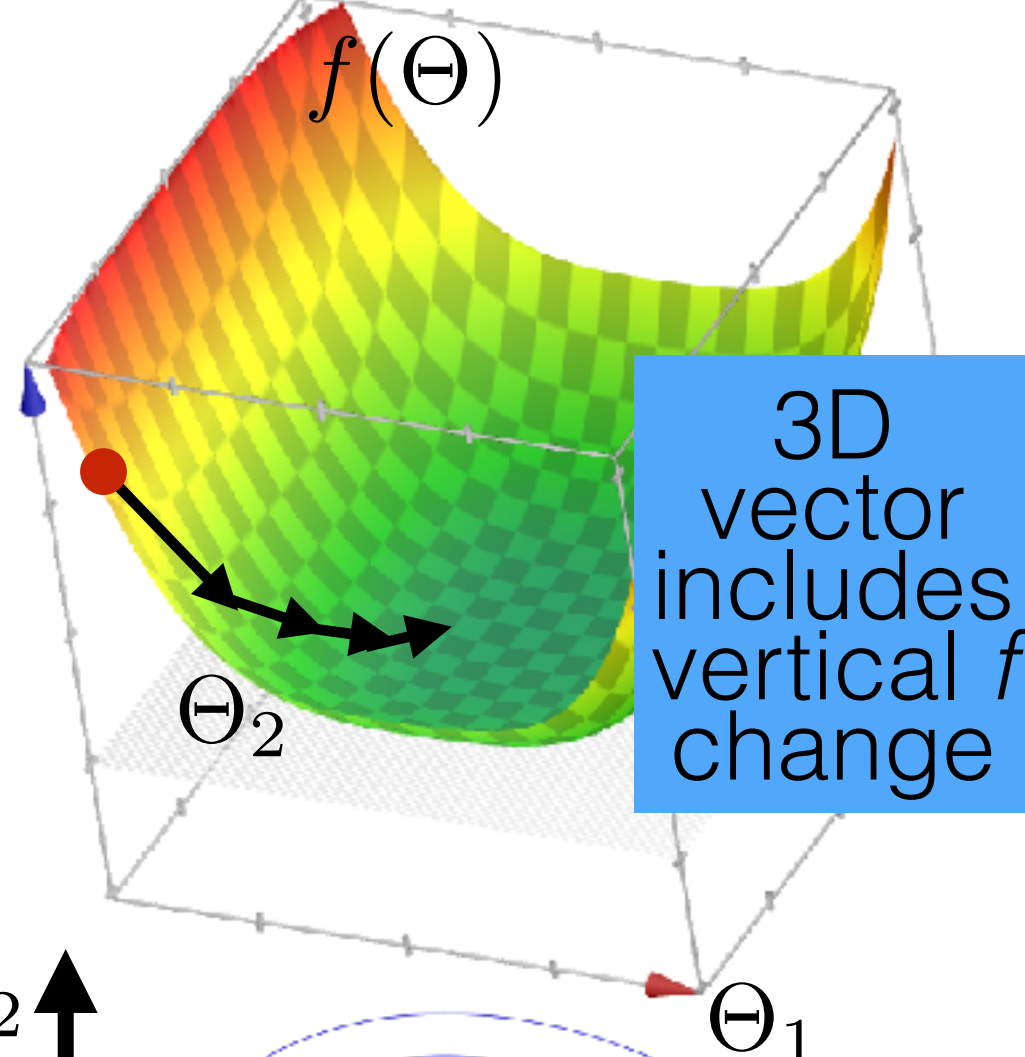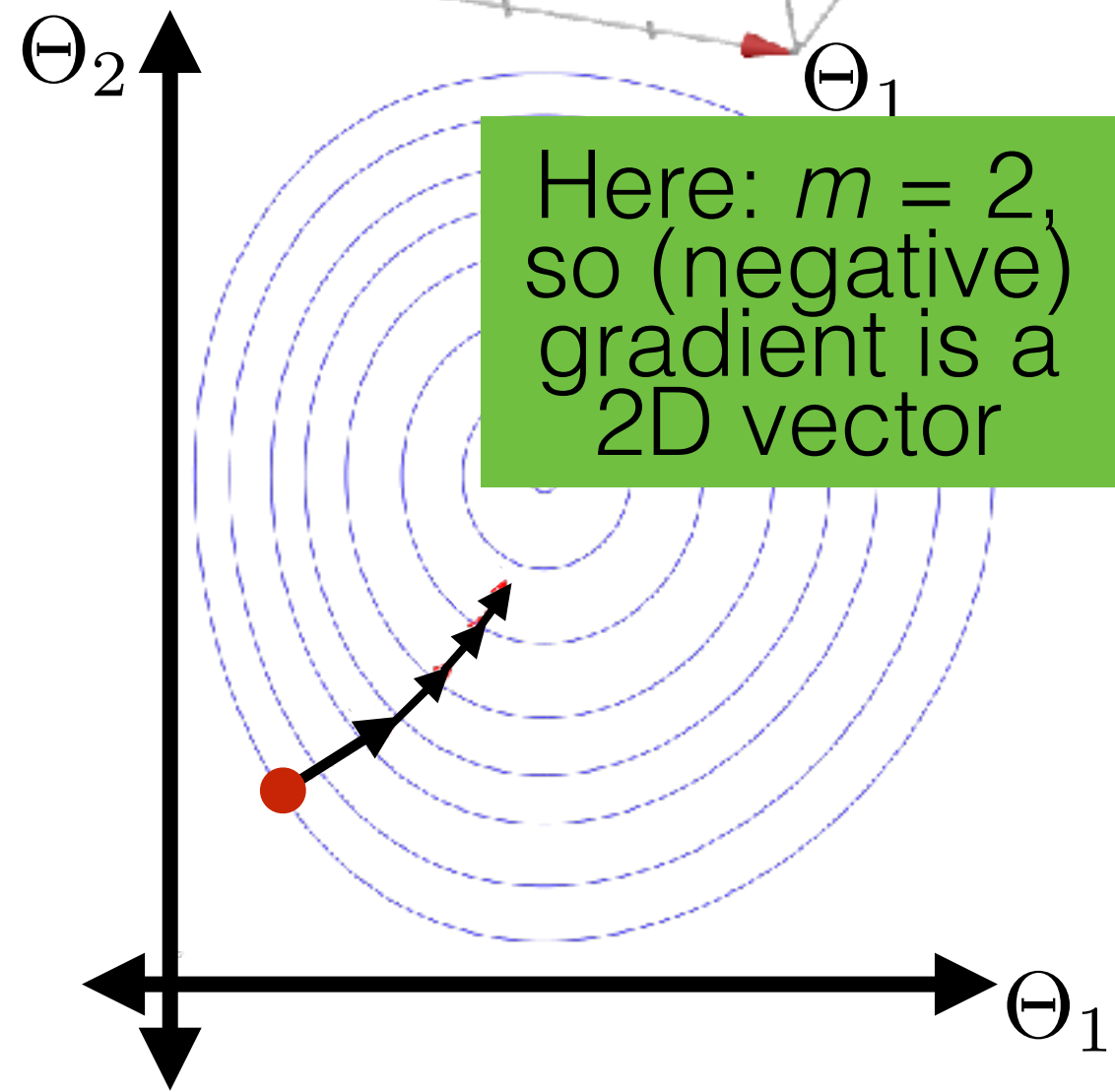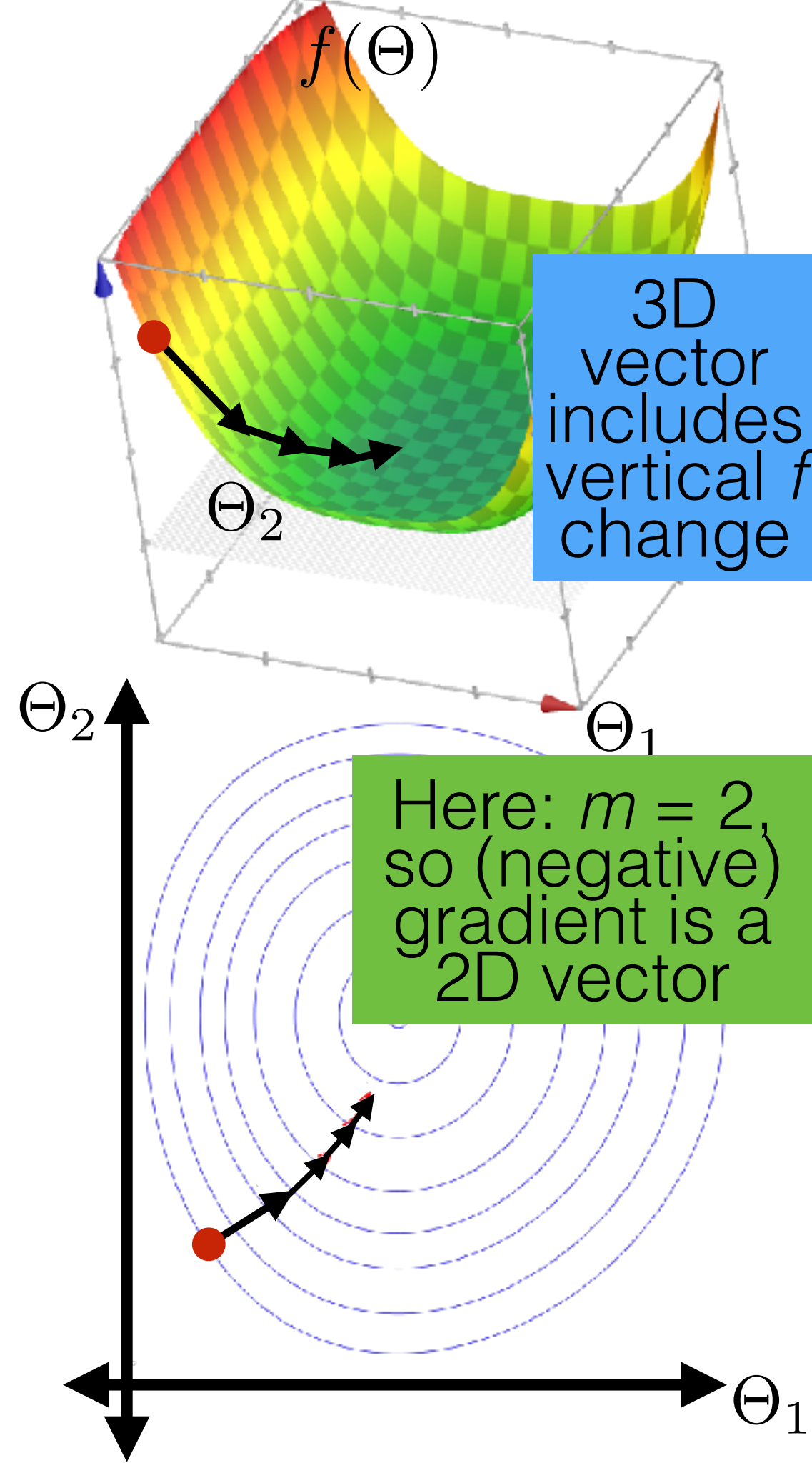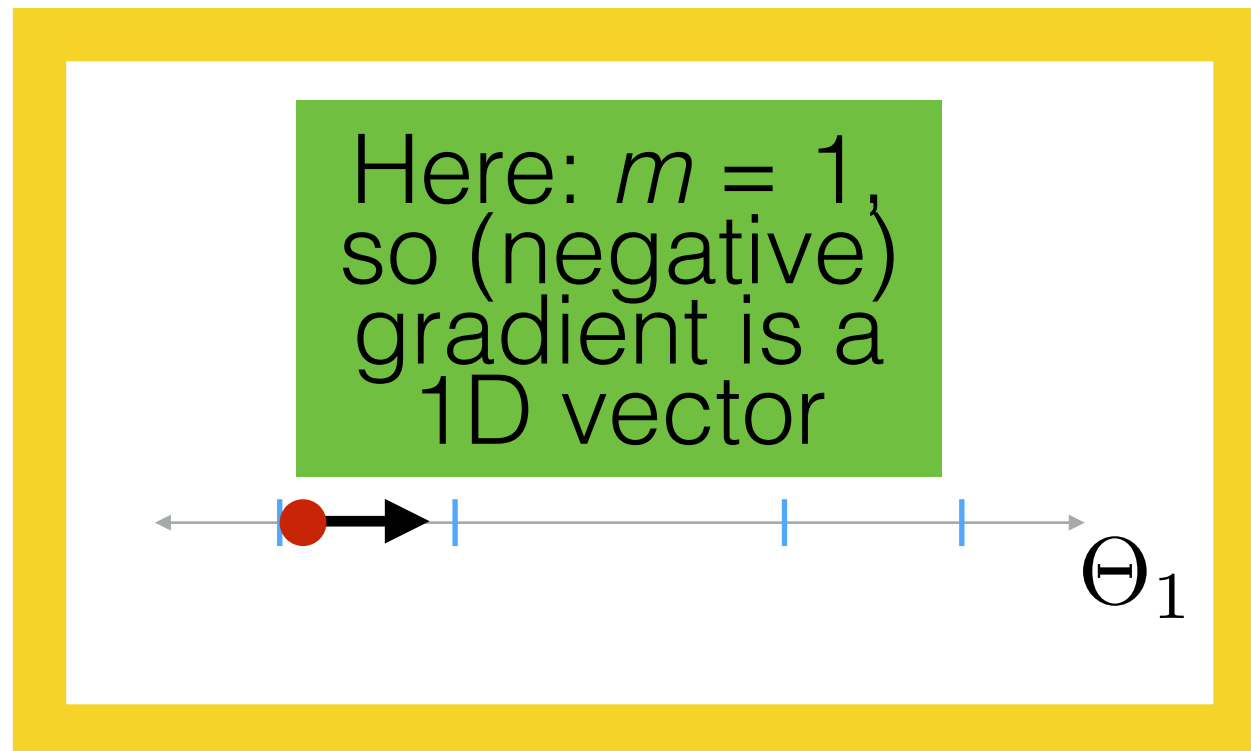
Here: $m = 2$, so (negative) gradient is a 2D vector
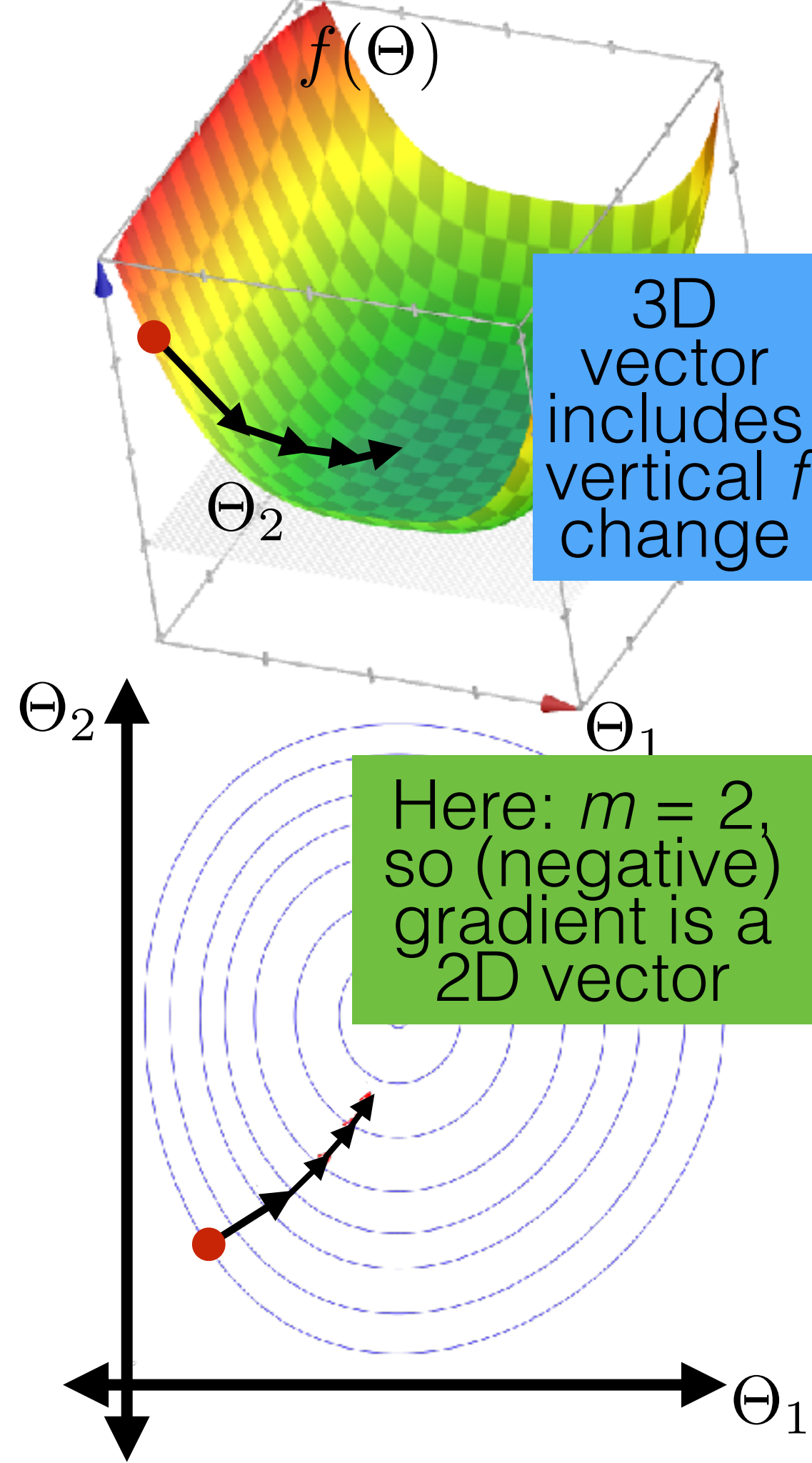
$\Theta_1$

# Gradient descent

- Gradient $\nabla_\Theta f = \left[\dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m}\right]^\top$
  - with $\Theta \in \mathbb{R}^m$

```
Gradient-Descent (Θ_init, η, f, ∇_Θ f, ε)
  Initialize Θ^(0) = Θ_init
  Initialize t = 0
```

**repeat**

$\qquad$ `t = t + 1`

$$\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$$

**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

- Other possible stopping criteria:
  - Max number of iterations $T$
  - $\| \Theta^{(t)} - \Theta^{(t-1)} \| < \epsilon$
  - $\| \nabla_\Theta f(\Theta^{(t)}) \| < \epsilon$

4

$f(\Theta)$

$\Theta_2$

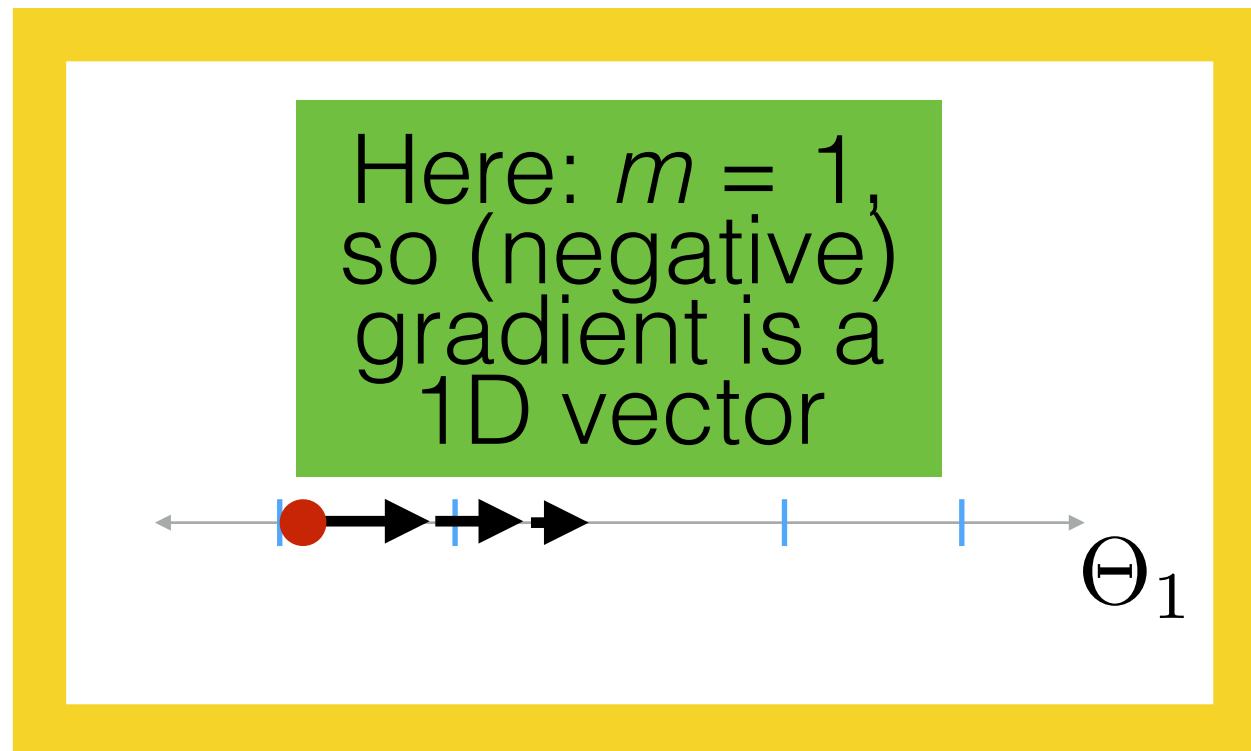3D vector includes vertical *f* change

$\Theta_2$

$\Theta_1$

Here: $m = 2$, so (negative) gradient is a 2D vector

$\Theta_1$

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]^\top$
  - with $\Theta \in \mathbb{R}^m$

$f(\Theta)$

3D vector includes vertical $f$ change

```
Gradient-Descent (Θ_init, η, f, ∇_Θ f, ε)
   Initialize Θ^(0) = Θ_init
   Initialize t = 0
```

**repeat**

```
   t = t + 1
```

$$\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$$

**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

$\Theta_2$

$\Theta_1$

Here: $m = 2$, so (negative) gradient is a 2D vector

- Other possible stopping criteria:
  - Max number of iterations $T$
  - $\| \Theta^{(t)} - \Theta^{(t-1)} \| < \epsilon$
  - $\| \nabla_\Theta f(\Theta^{(t)}) \| < \epsilon$

4

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \ldots, \dfrac{\partial f}{\partial \Theta_m} \right]$
  - with $\Theta \in \mathbb{R}^m$

```
Gradient-Descent (Θ_init, η, f, ∇_Θf, ε)
  Initialize Θ^(0) = Θ_init
  Initialize t = 0
```
**repeat**
  t = t + 1

$$\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$$

**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

- Other possible stopping criteria:
  - Max number of iterations $T$
  - $\| \Theta^{(t)} - \Theta^{(t-1)} \| < \epsilon$
  - $\| \nabla_\Theta f(\Theta^{(t)}) \| < \epsilon$

4



$f(\Theta)$

$\Theta_2$

$\Theta_2$   $\Theta_1$

3D vector includes vertical $f$ change

Here: $m = 2$, so (negative) gradient is a 2D vector

$\Theta_1$

# Gradient descent

- Gradient $\nabla_\Theta f = \left[ \dfrac{\partial f}{\partial \Theta_1}, \dots, \dfrac{\partial f}{\partial \Theta_m} \right]$
  - with $\Theta \in \mathbb{R}^m$

$\texttt{Gradient-Descent}\,(\Theta_{\text{init}}, \eta, f, \nabla_\Theta f, \epsilon$

$\texttt{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\texttt{Initialize t = 0}$

**repeat**

$\quad \texttt{t = t + 1}$

$\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** $\left| f(\Theta^{(t)}) - f(\Theta^{(t-1)}) \right| < \epsilon$

**Return** $\Theta^{(t)}$

- Other possible stopping criteria:
  - Max number of iterations $T$
  - $\| \Theta^{(t)} - \Theta^{(t-1)} \| < \epsilon$
  - $\| \nabla_\Theta f(\Theta^{(t)}) \| < \epsilon$

4



$f(\Theta)$

$\Theta_2$

3D vector includes vertical $f$ change

$\Theta_1$

Here: $m = 2$, so (negative) gradient is a 2D vector

$\Theta_1$

# Gradient descent



$f(\Theta)$

$\Theta_2$

3D vector includes vertical $f$ change

$\Theta_2$

$\Theta_1$

Here: $m = 2$, so (negative) gradient is a 2D vector

$\Theta_1$

4

# Gradient descent

$f(\Theta)$



$\Theta_1$

$f(\Theta)$

$\Theta_2$

3D vector includes vertical $f$ change

$\Theta_2$

$\Theta_1$

Here: $m = 2$, so (negative) gradient is a 2D vector

4

# Gradient descent

$f(\Theta)$



$\Theta_1$

$f(\Theta)$

3D vector includes vertical $f$ change

$\Theta_2$

$\Theta_1$

$\Theta_2$

Here: $m = 2$, so (negative) gradient is a 2D vector

$\Theta_1$

$\Theta_1$

4

# Gradient descent

$f(\Theta)$

$\Theta_1$



$f(\Theta)$

3D vector includes vertical $f$ change

$\Theta_2$

$\Theta_1$

$\Theta_2$

Here: $m = 2$, so (negative) gradient is a 2D vector

$\Theta_1$

$\Theta_1$

4

# Gradient descent

$f(\Theta)$

$\Theta_1$

$f(\Theta)$

$\Theta_2$

3D vector includes vertical $f$ change

$\Theta_2$

$\Theta_1$

Here: $m$ = 1, so (negative) gradient is a 1D vector

$\Theta_1$

Here: $m$ = 2, so (negative) gradient is a 2D vector

$\Theta_1$

4

# Gradient descent

$f(\Theta)$

$\Theta_1$

$f(\Theta)$

$f(\Theta)$

3D vector includes vertical *f* change

$\Theta_2$

$\Theta_1$

$\Theta_2$

Here: $m = 2$, so (negative) gradient is a 2D vector

Here: $m = 1$, so (negative) gradient is a 1D vector

$\Theta_1$

$\Theta_1$

4

# Gradient descent

$f(\Theta)$

$\Theta_1$

$f(\Theta)$

$f(\Theta)$

$\Theta_2$

3D vector includes vertical *f* change

$\Theta_1$

$\Theta_2$

Here: *m* = 2, so (negative) gradient is a 2D vector

Here: *m* = 1, so (negative) gradient is a 1D vector

$\Theta_1$

$\Theta_1$

4

# Gradient descent

$f(\Theta)$

$\Theta_1$

$f(\Theta)$

3D vector includes vertical $f$ change

$\Theta_2$

$\Theta_1$

Here: $m = 1$, so (negative) gradient is a 1D vector

$\Theta_1$

$\Theta_2$

Here: $m = 2$, so (negative) gradient is a 2D vector

$\Theta_1$

4

# Gradient descent

$f(\Theta)$

$\Theta_1$

$f(\Theta)$

3D vector includes vertical $f$ change

$\Theta_2$

$\Theta_1$

$\Theta_2$

Here: $m = 2$, so (negative) gradient is a 2D vector

Here: $m = 1$, so (negative) gradient is a 1D vector

$\Theta_1$

$\Theta_1$

4

# Gradient descent

$f(\Theta)$

$\Theta_1$

2D vector includes vertical $f$ change

$f(\Theta)$

$\Theta_2$

3D vector includes vertical $f$ change

$\Theta_1$

$\Theta_2$

Here: $m = 1$, so (negative) gradient is a 1D vector

$\Theta_1$

Here: $m = 2$, so (negative) gradient is a 2D vector

$\Theta_1$

4

# Gradient descent

$f(\Theta)$

$\Theta_1$

2D vector includes vertical $f$ change

$f(\Theta)$

3D vector includes vertical $f$ change

$\Theta_2$

$\Theta_1$

$\Theta_2$

Here: $m = 2$, so (negative) gradient is a 2D vector

Here: $m = 1$, so (negative) gradient is a 1D vector

$\Theta_1$

# Gradient descent

$f(\Theta)$

2D vector includes vertical $f$ change

$\Theta_1$

$f(\Theta)$

3D vector includes vertical $f$ change

$\Theta_2$

$\Theta_1$

$\Theta_2$

Here: $m = 1$, so (negative) gradient is a 1D vector

$\Theta_1$

Here: $m = 2$, so (negative) gradient is a 2D vector

$\Theta_1$

4

# Gradient descent

$f(\Theta)$



2D vector includes vertical $f$ change

$f(\Theta)$

3D vector includes vertical $f$ change

$\Theta_2$

$\Theta_1$

$\Theta_2$

$\Theta_1$

Here: $m = 1$, so (negative) gradient is a 1D vector

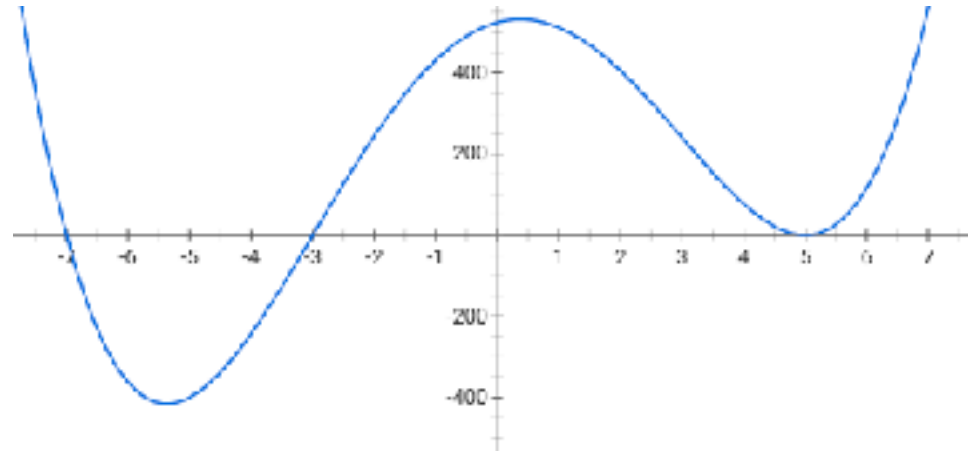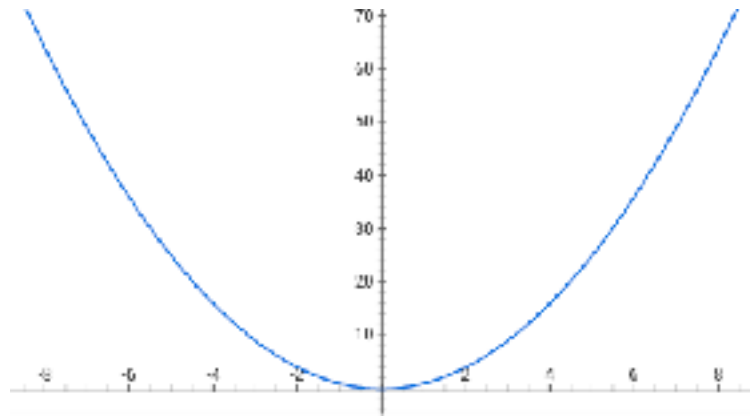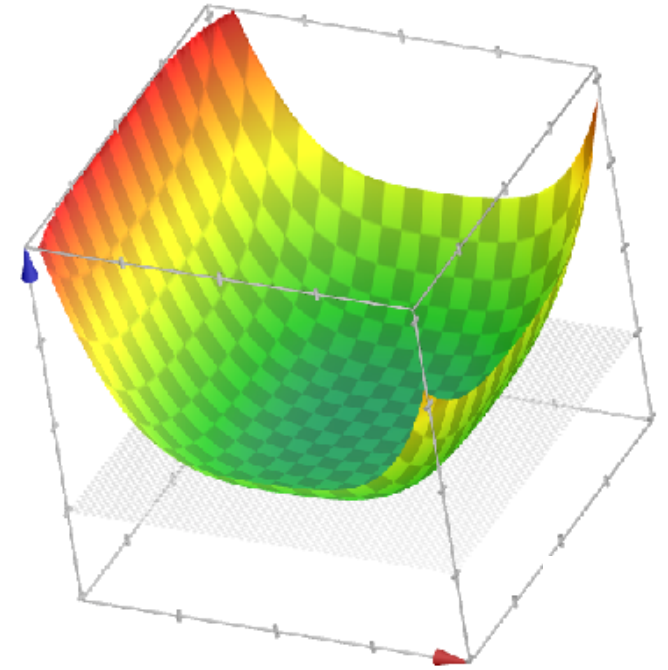Here: $m = 2$, so (negative) gradient is a 2D vector

$\Theta_1$

4

# Gradient descent properties

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
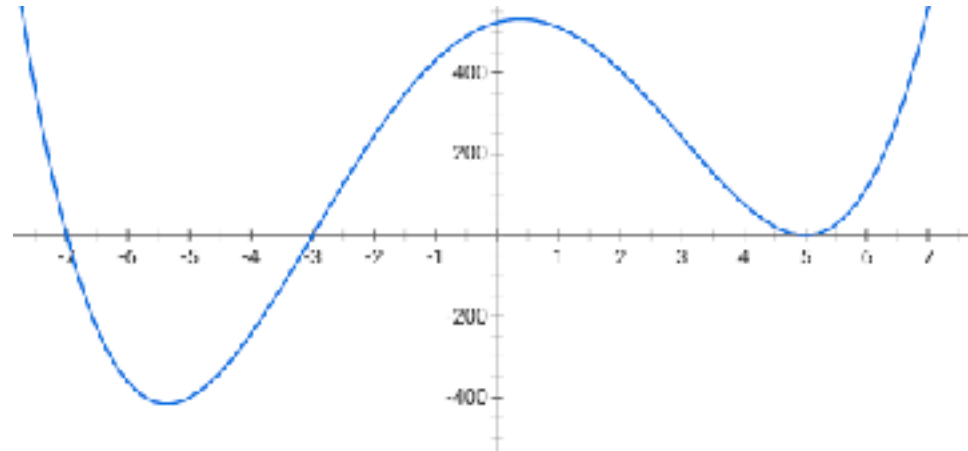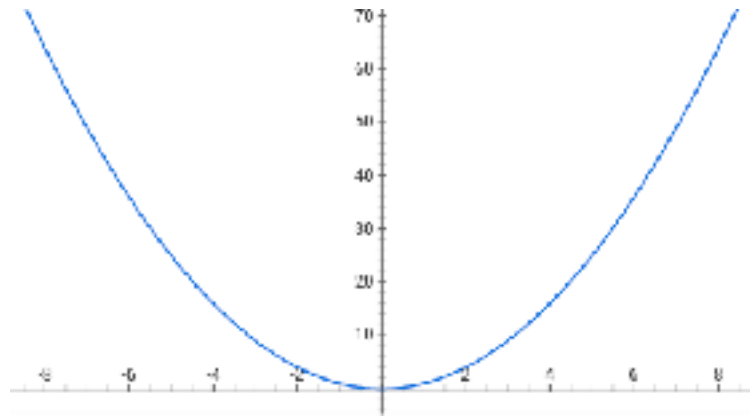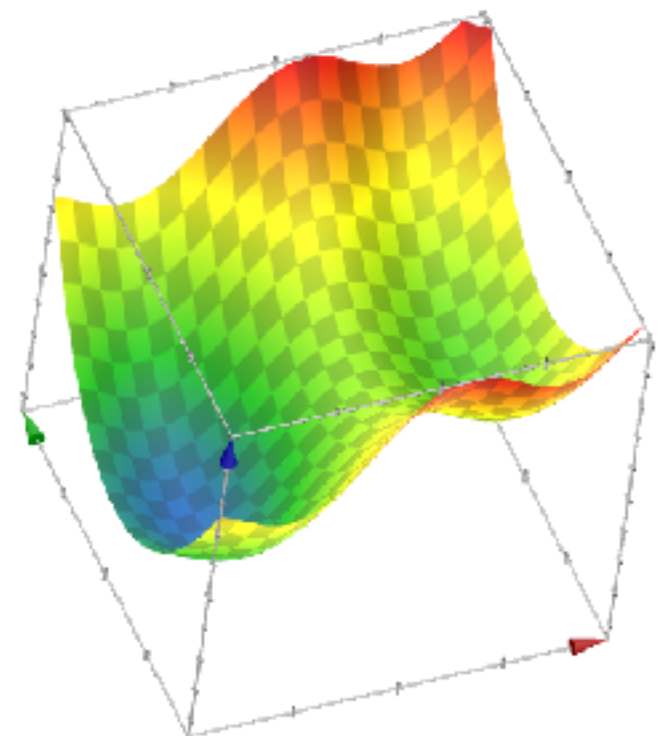
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
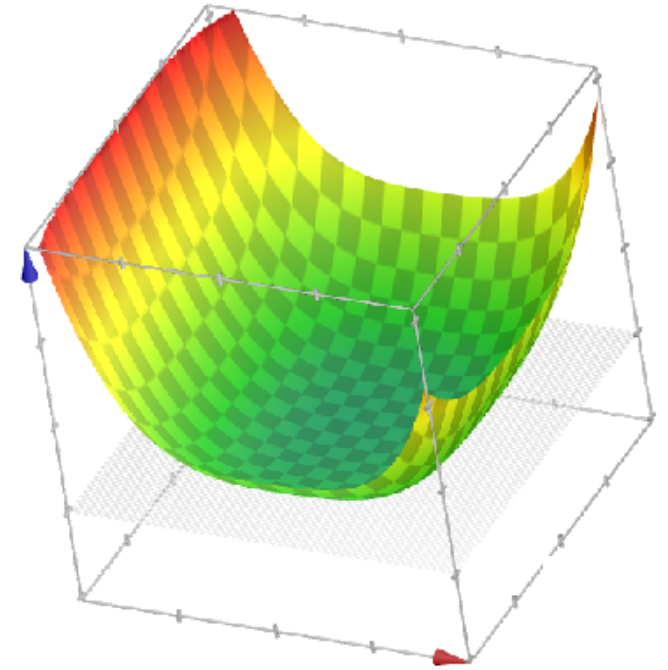
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
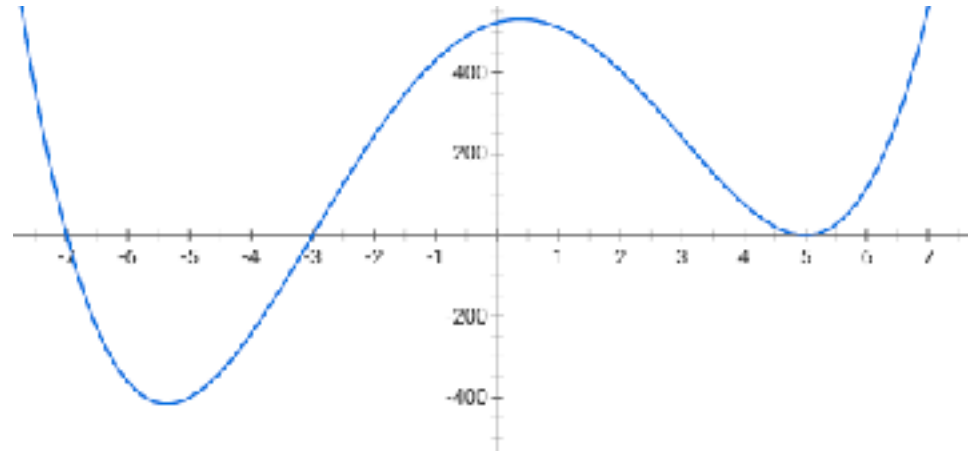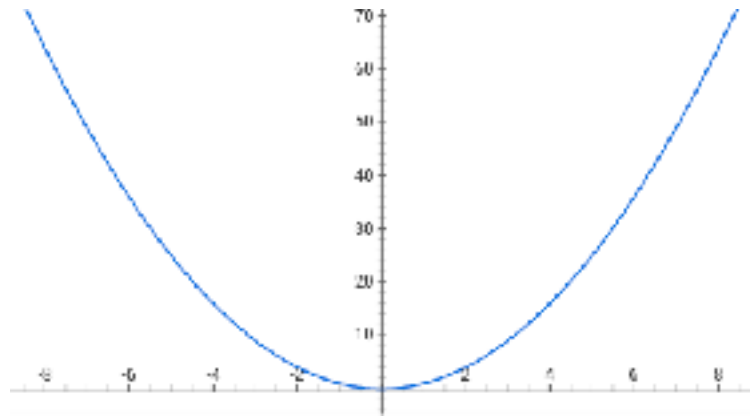
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
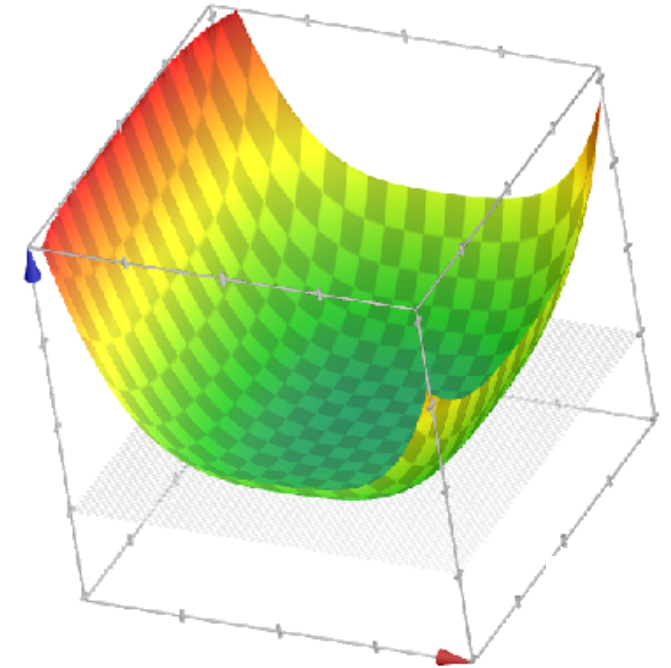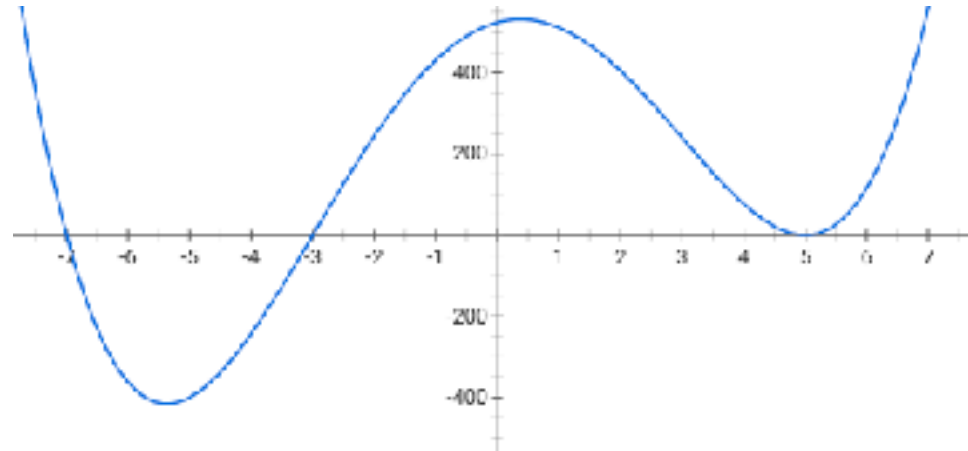
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
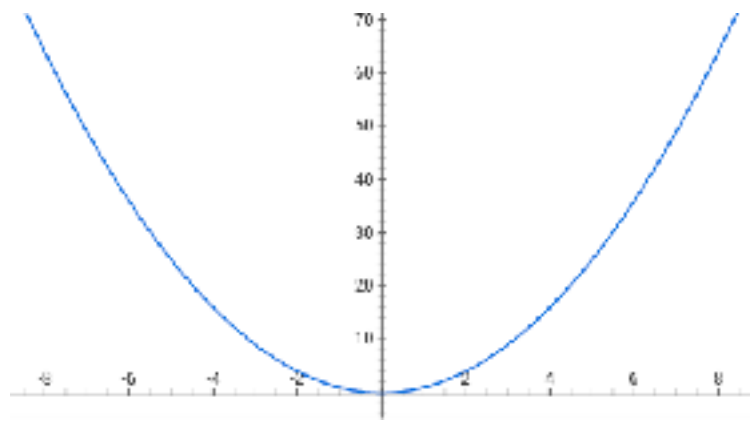
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
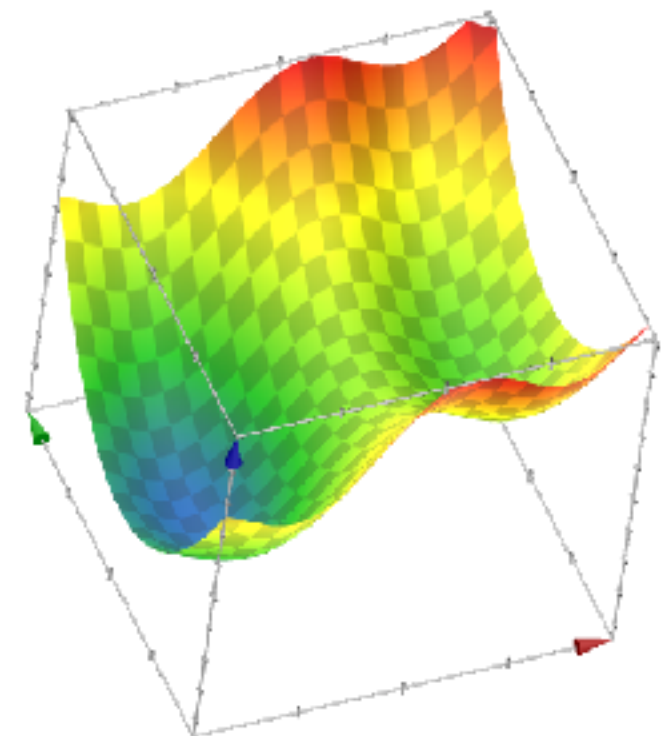
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
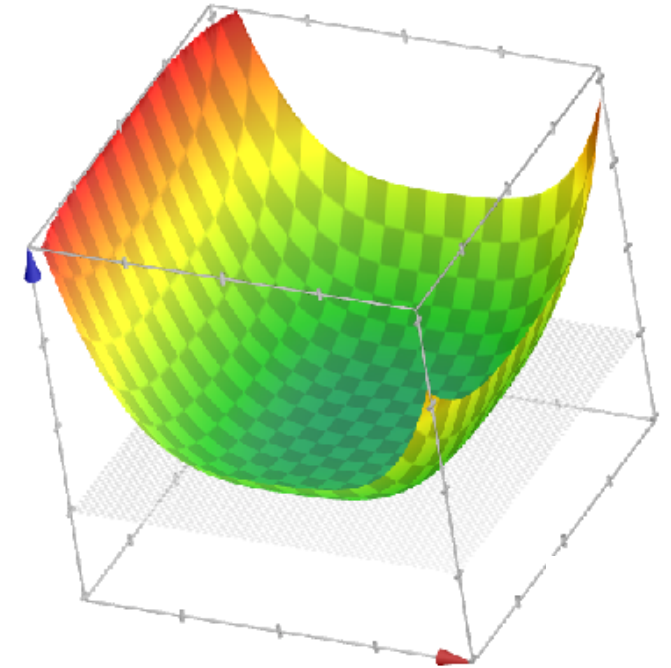
# Gradient descent properties

- A function *f* on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of *f* lies above or on the graph
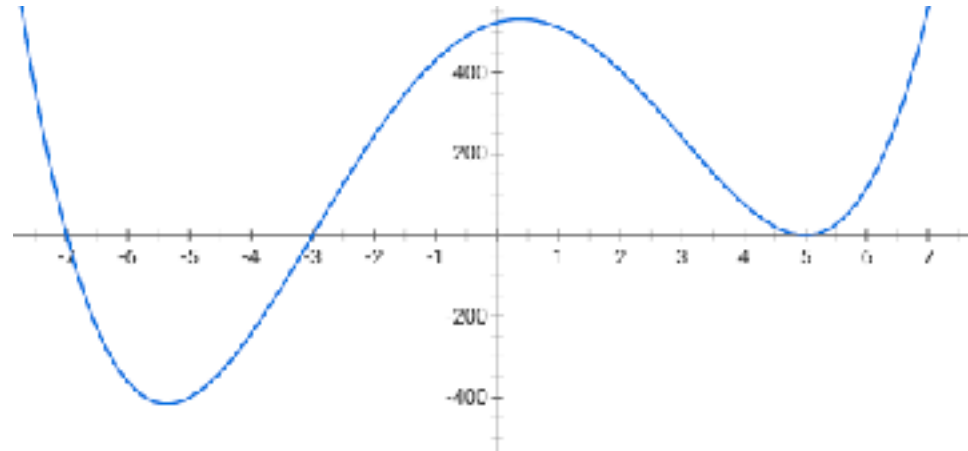
# Gradient descent properties

- A function _f_ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of _f_ lies above or on the graph
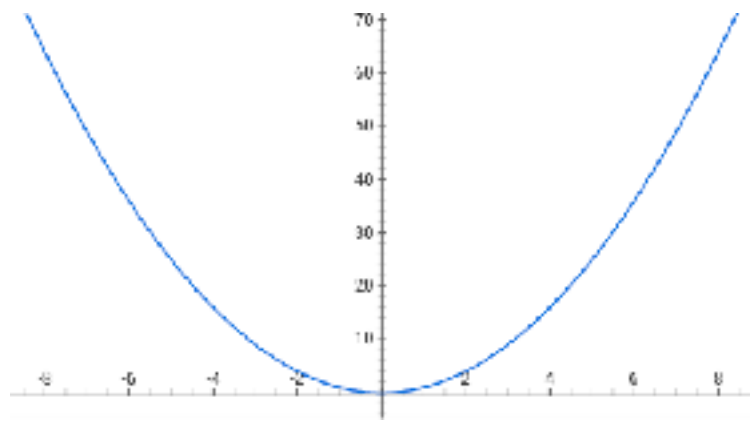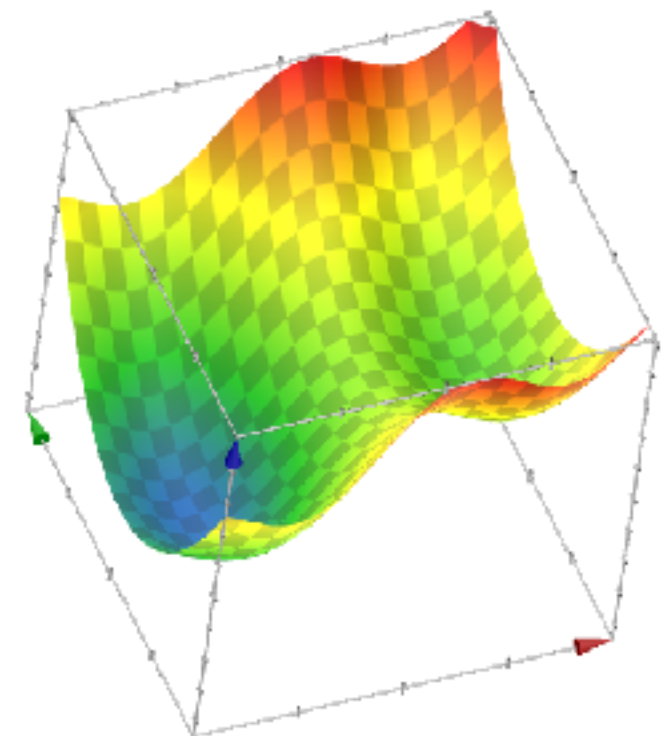
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
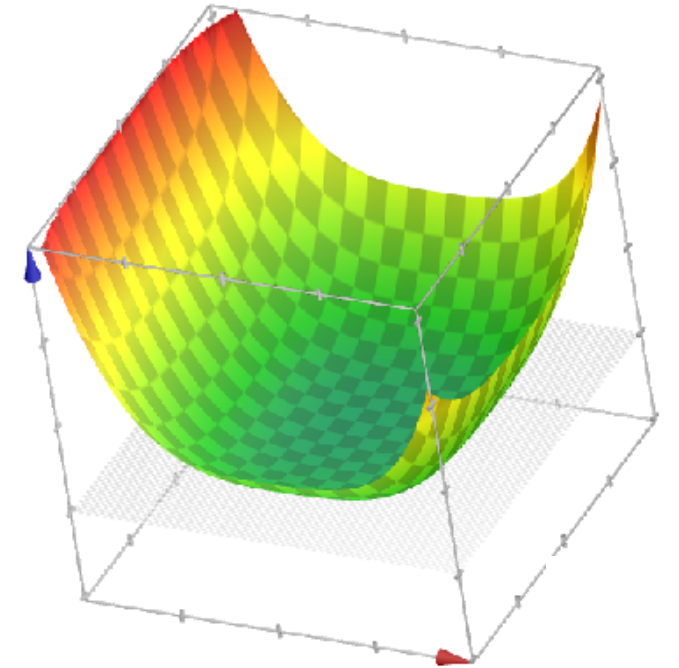
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
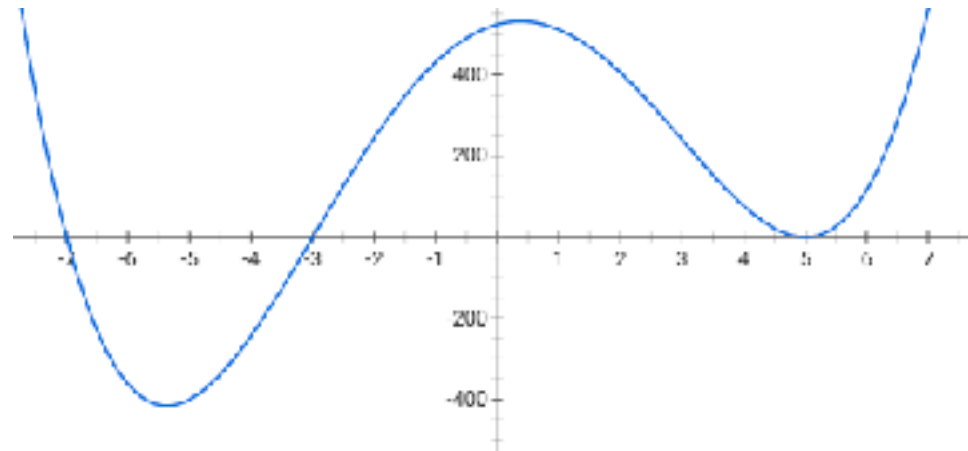
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
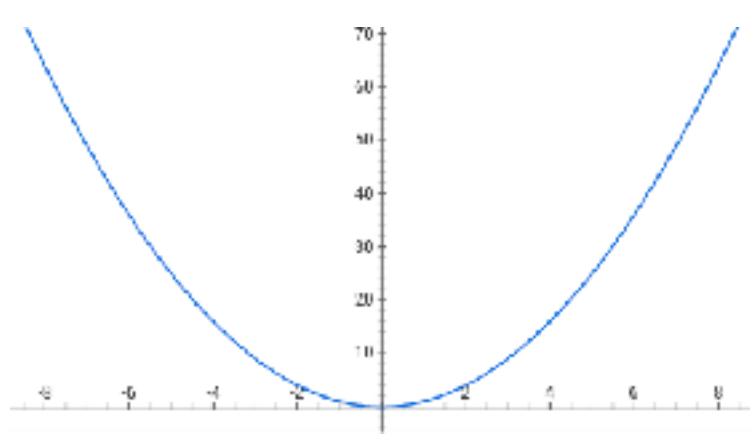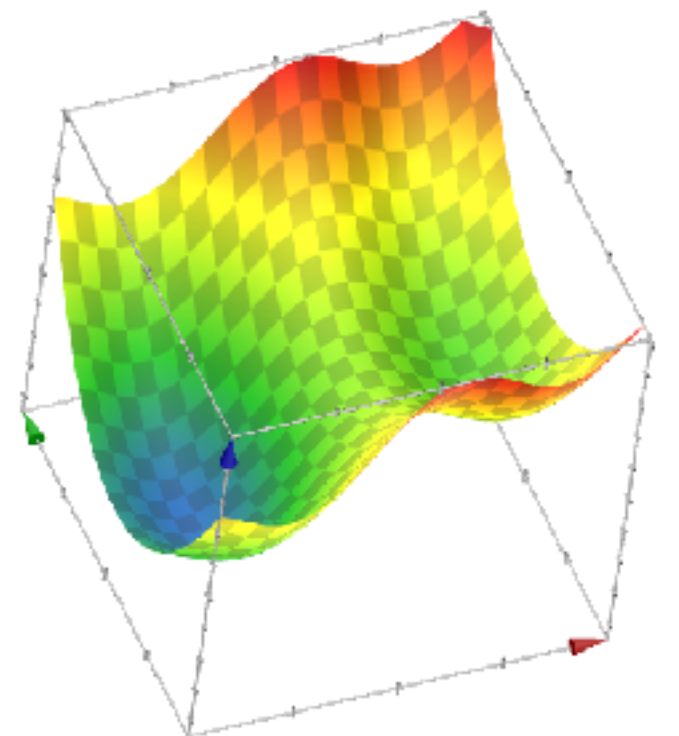
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
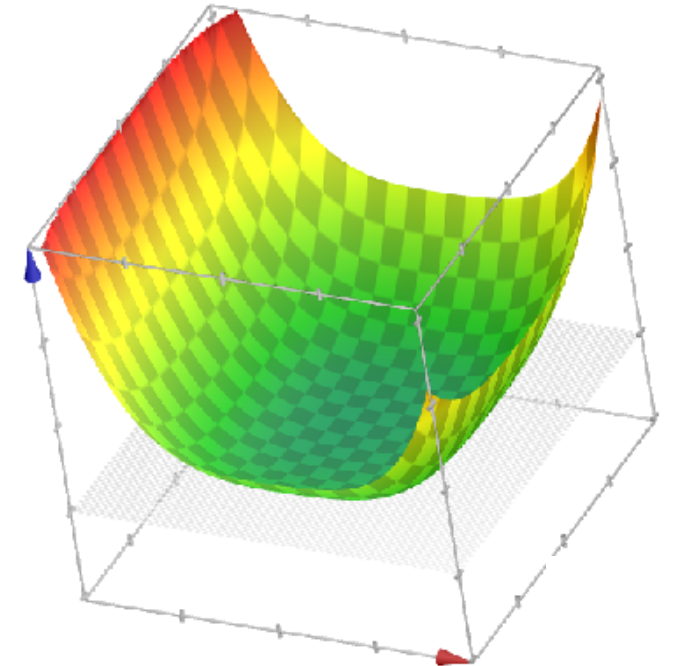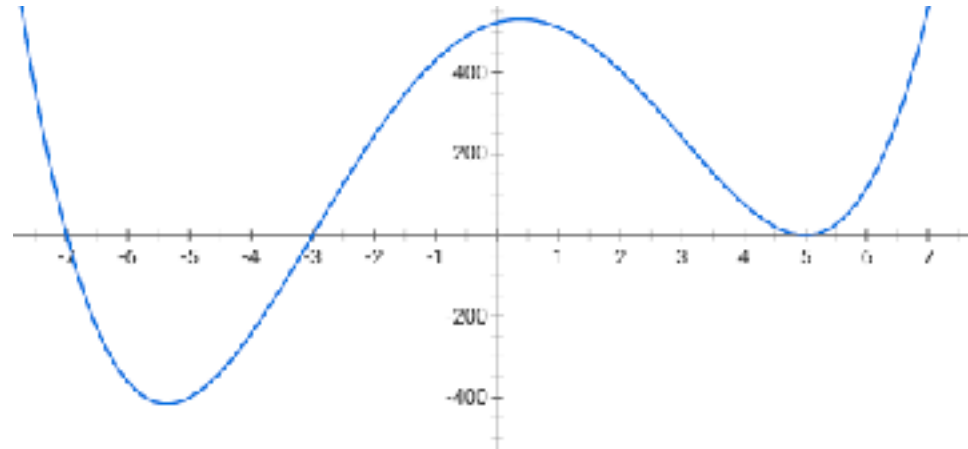
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
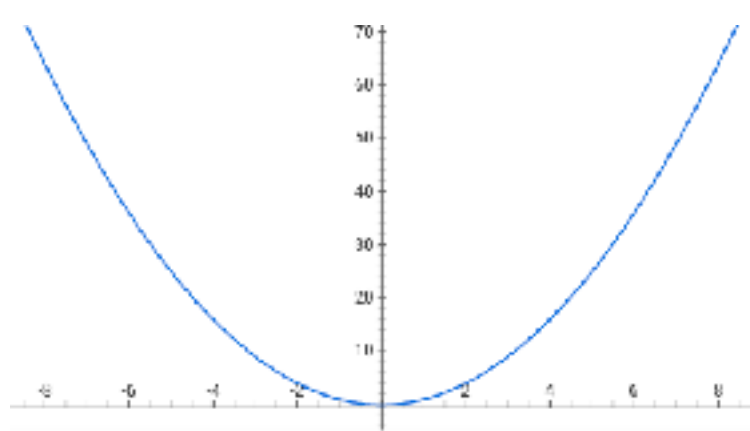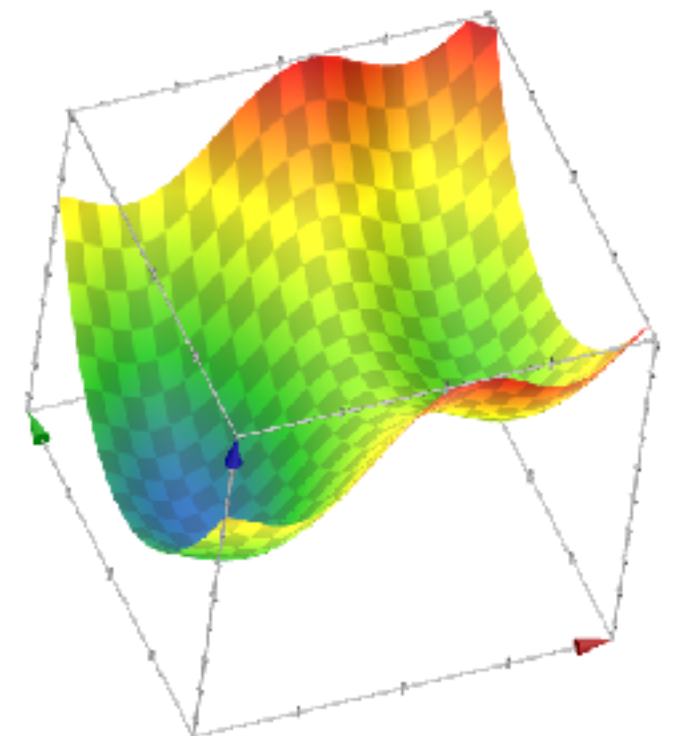
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
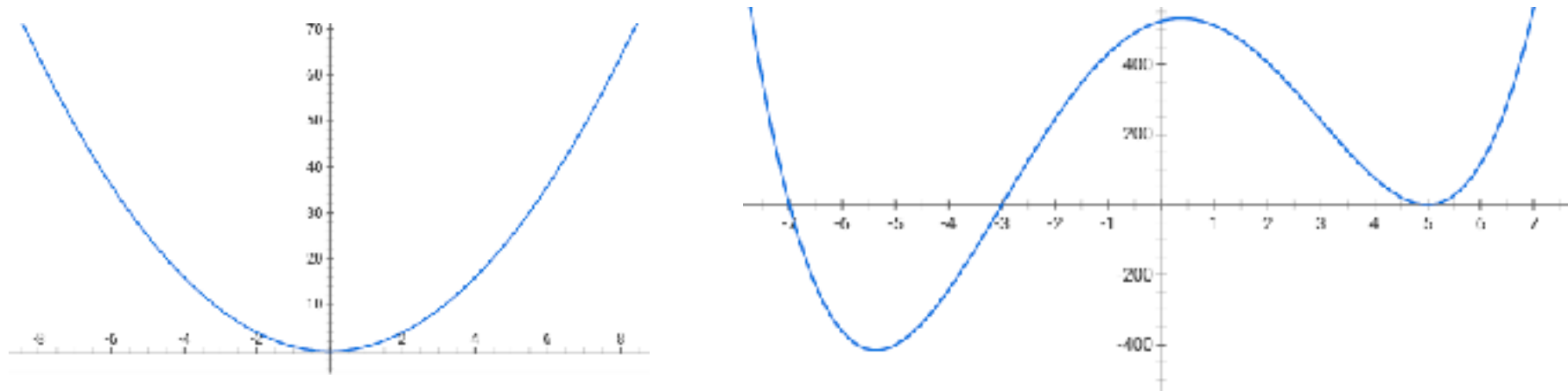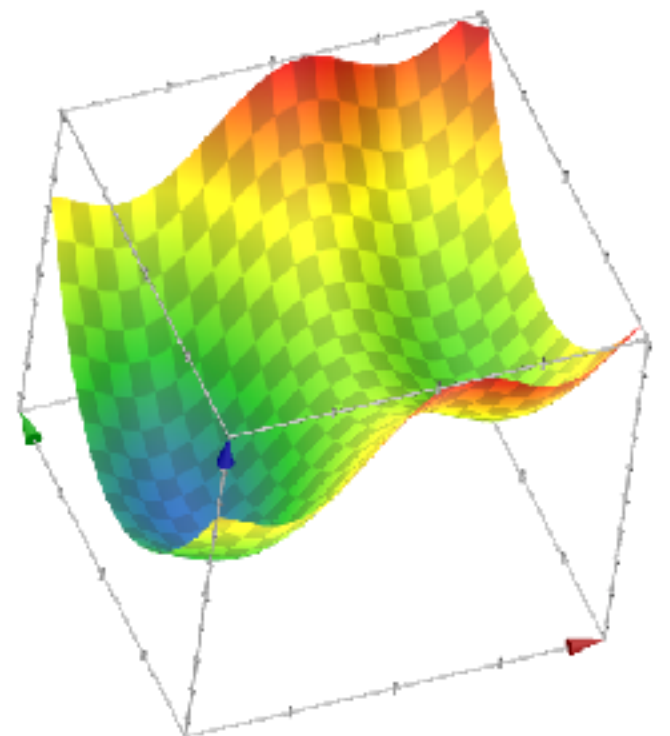
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph






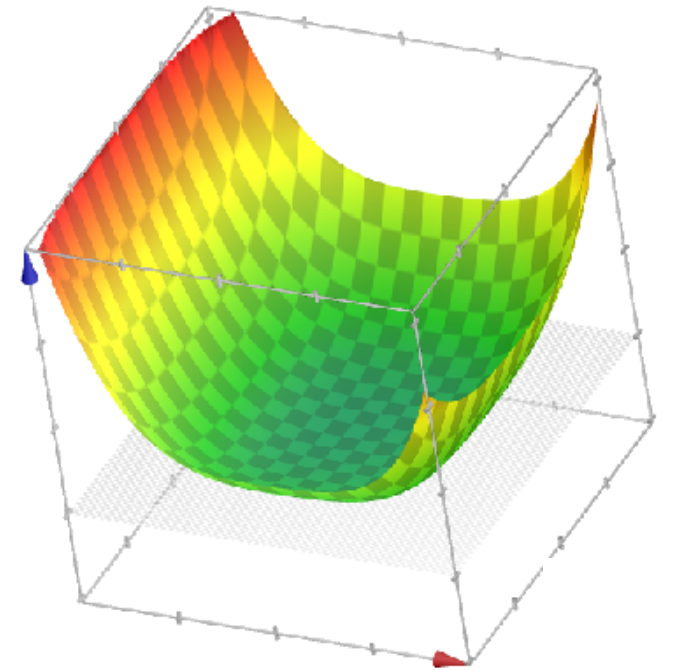
- **Theorem**: Gradient descent performance

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
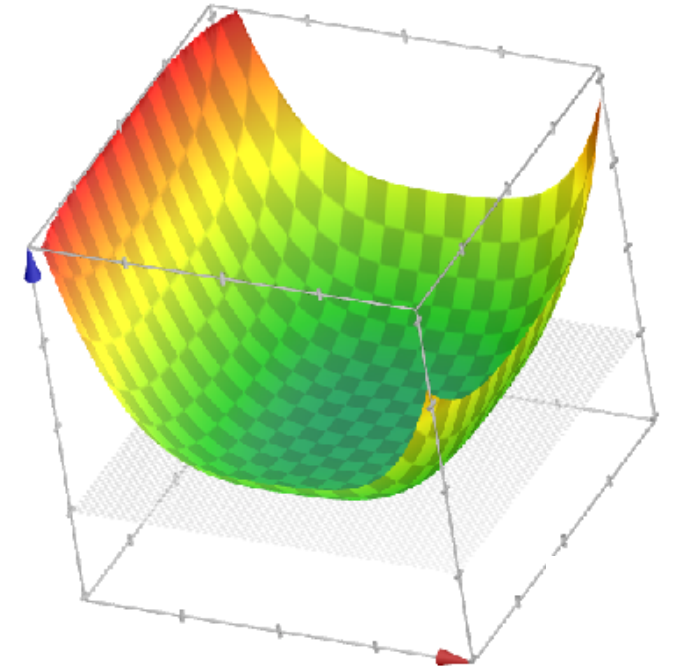
- **Theorem**: Gradient descent performance
  - **Assumptions**:

# Gradient descent properties

- A function *f* on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of *f* lies above or on the graph
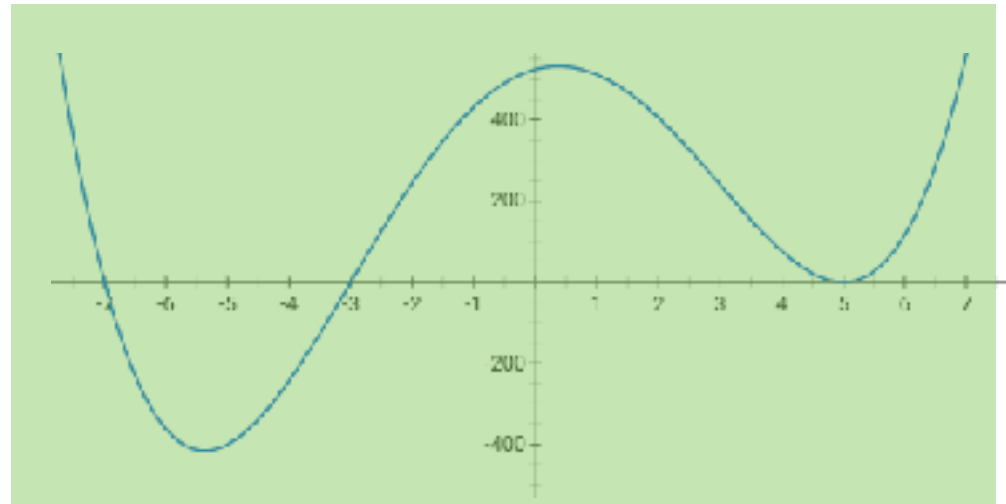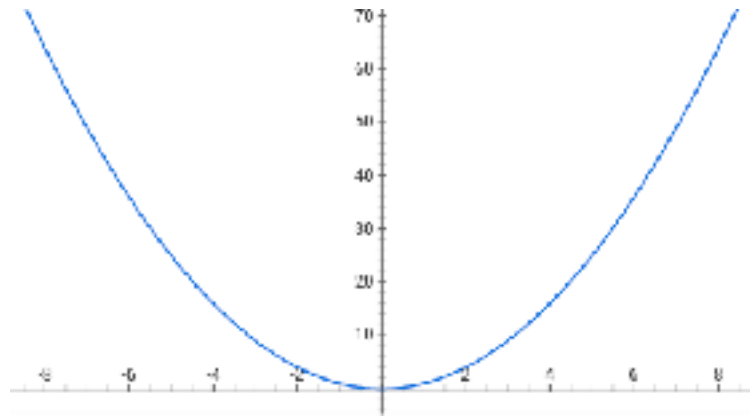
- **Theorem**: Gradient descent performance
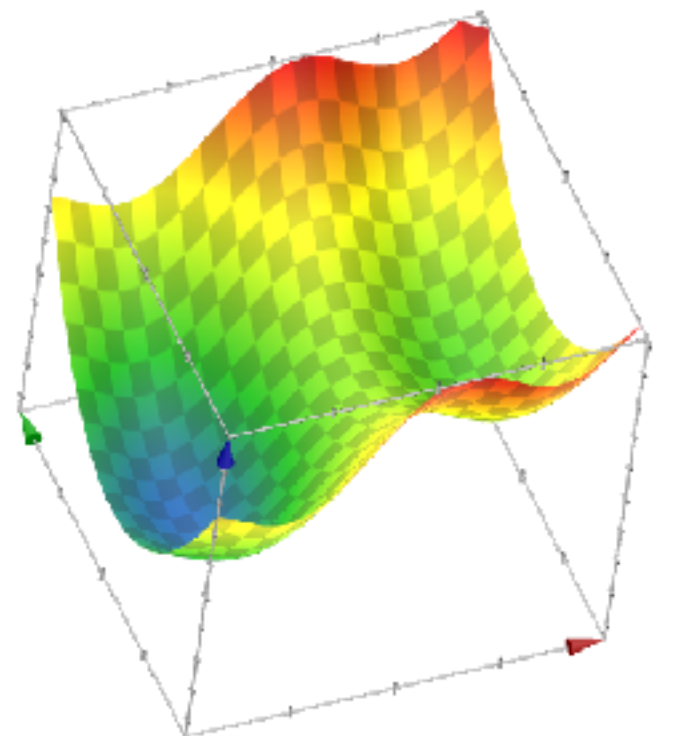  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
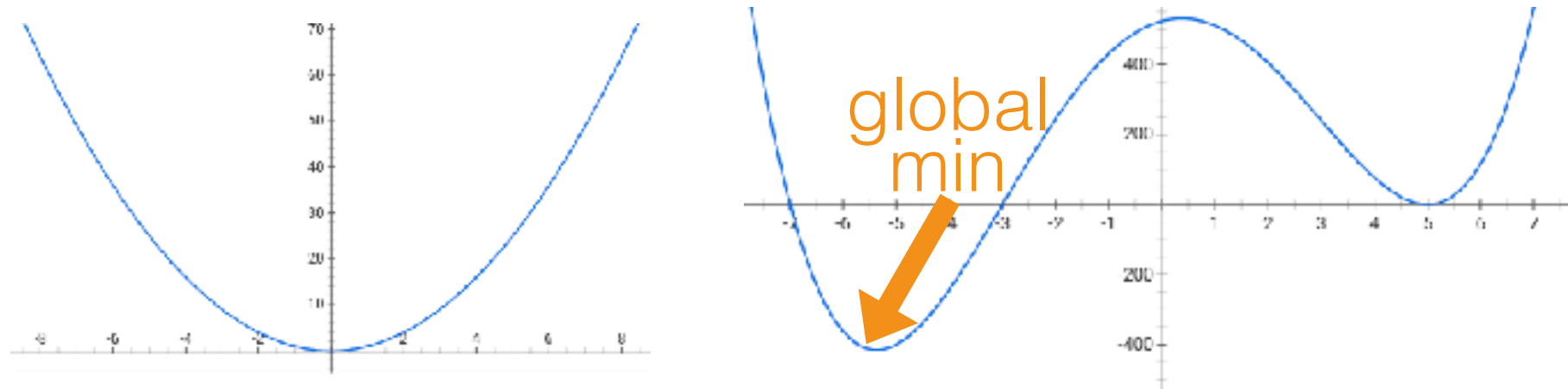
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
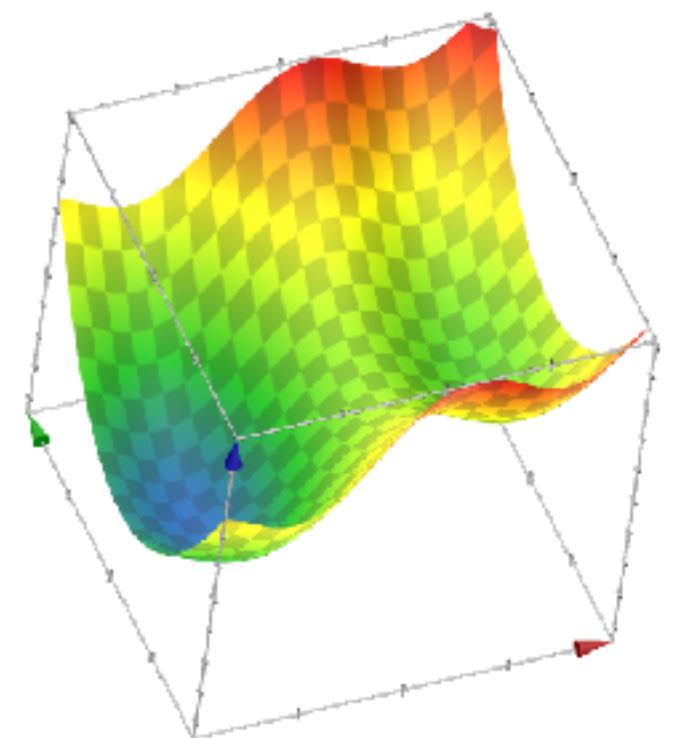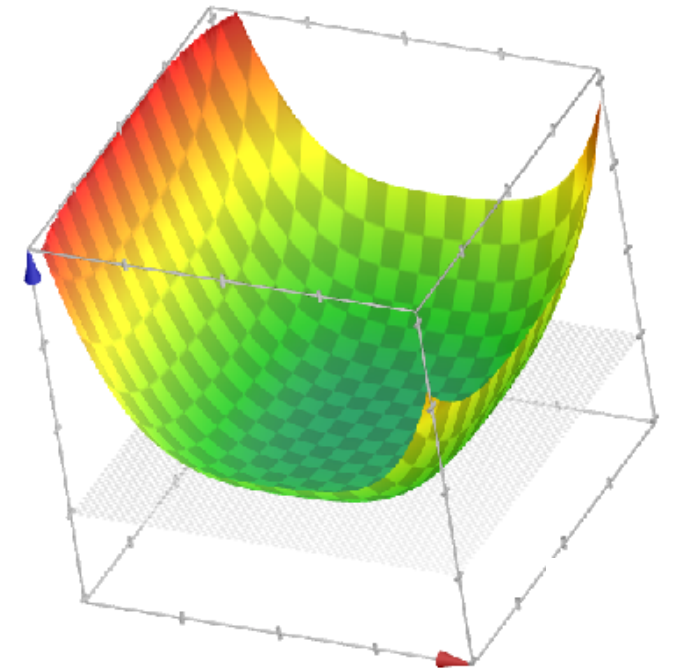    - $f$ is sufficiently "smooth" and convex

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph



- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum

# Gradient descent properties

- A function *f* on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of *f* lies above or on the graph
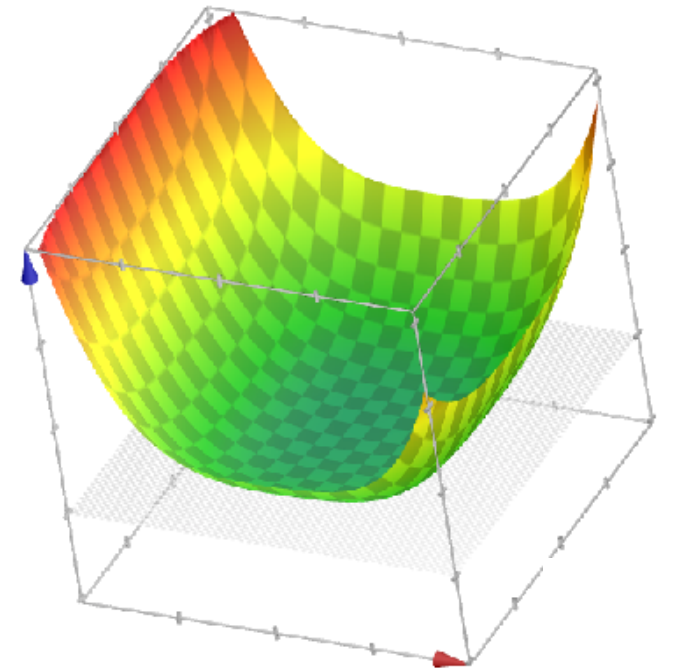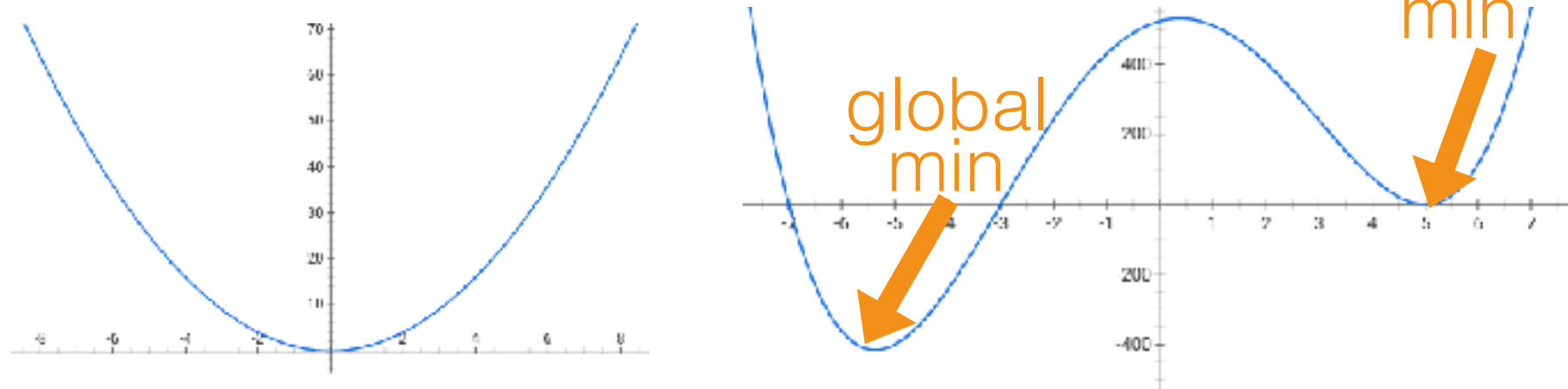
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - *f* is sufficiently "smooth" and convex
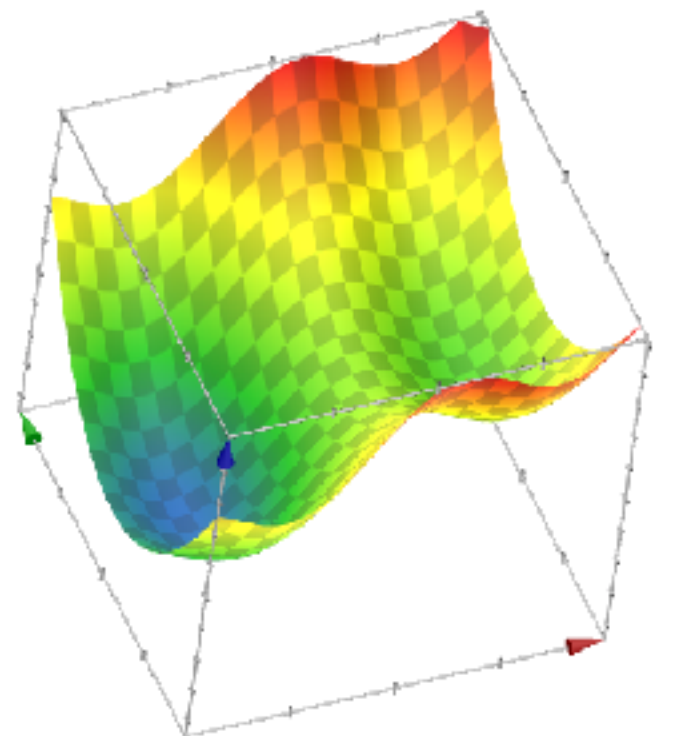    - *f* has at least one global optimum

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
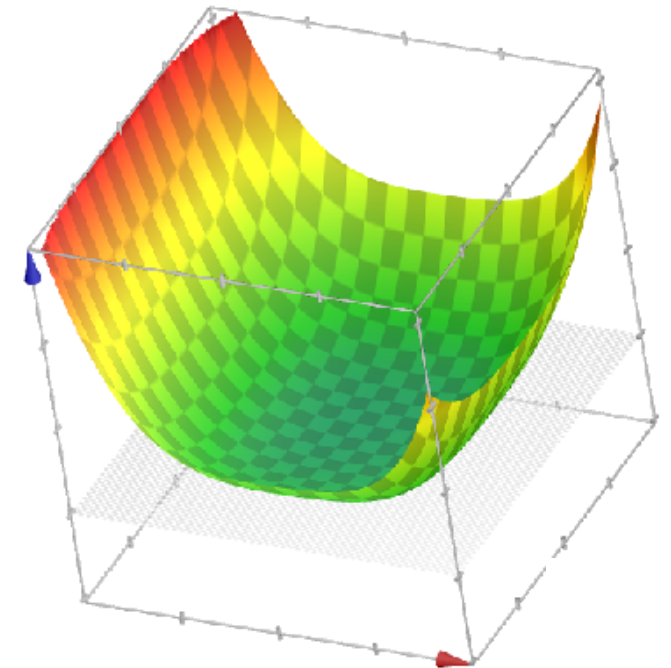


global min



- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
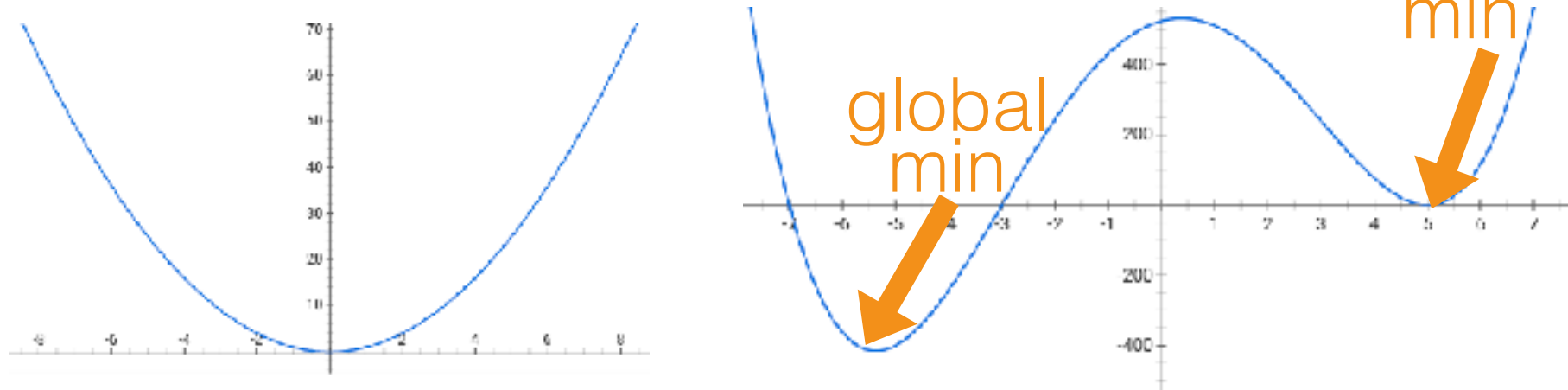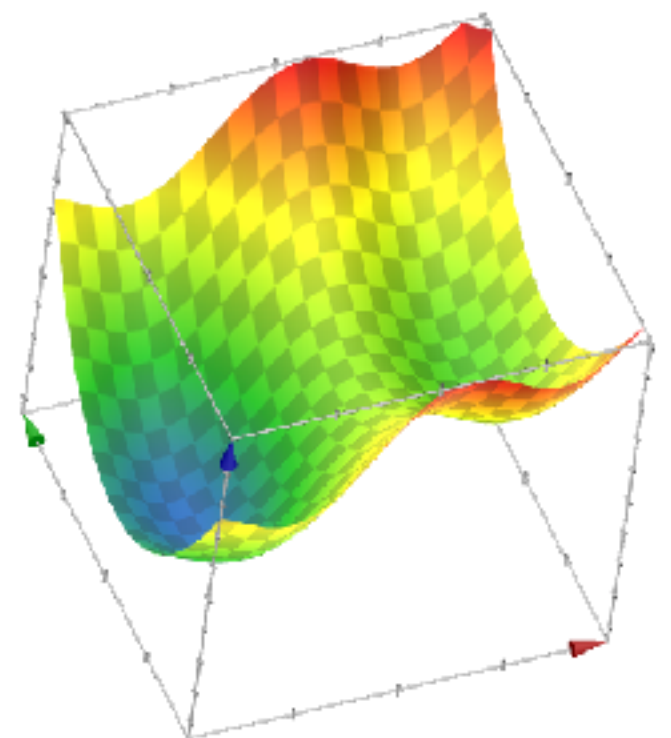


global min

local min

- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
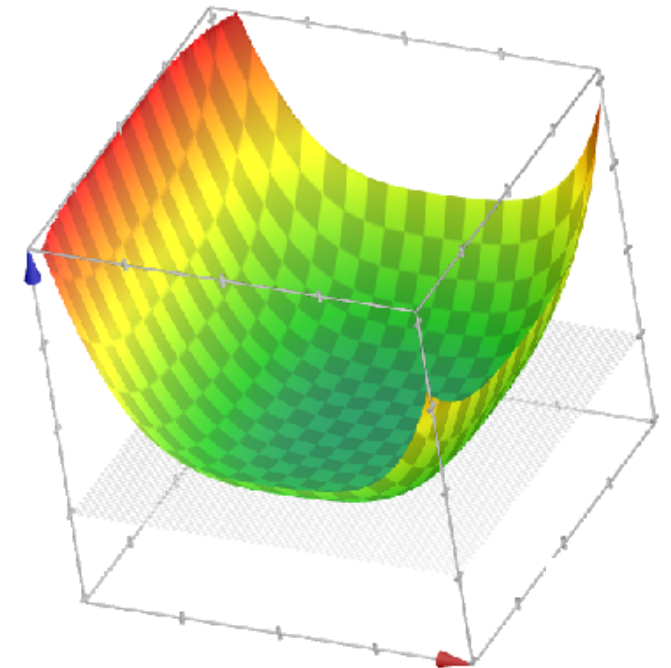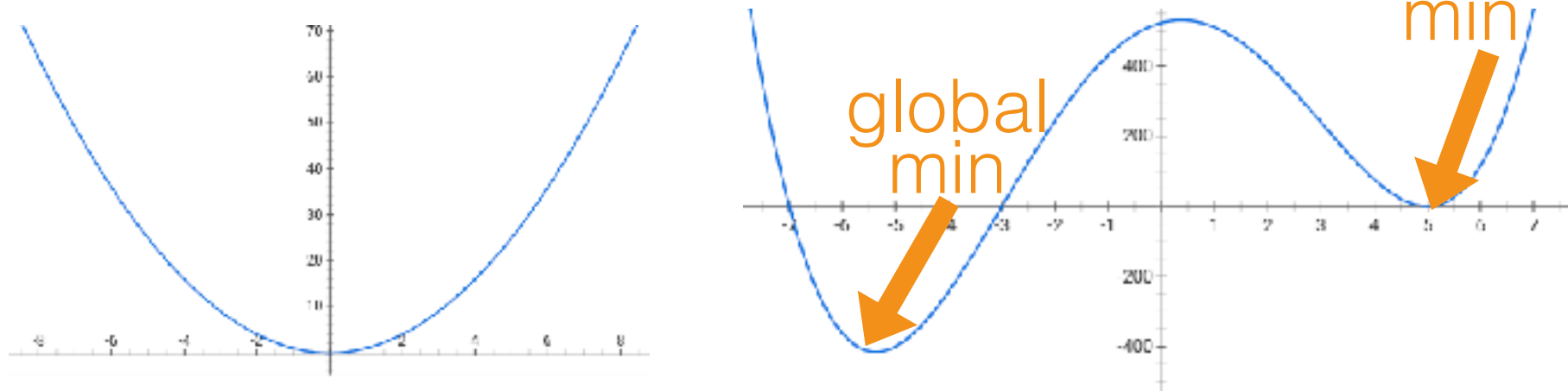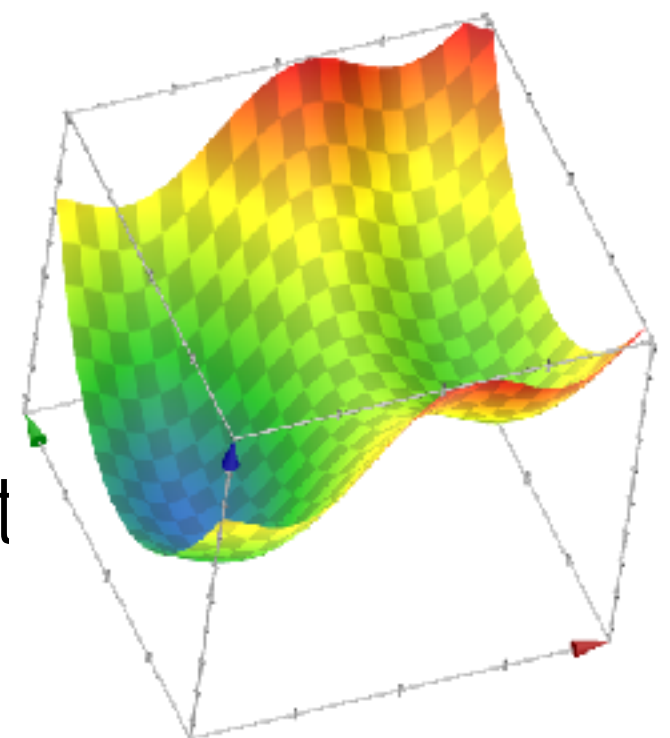


global min

local min

- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
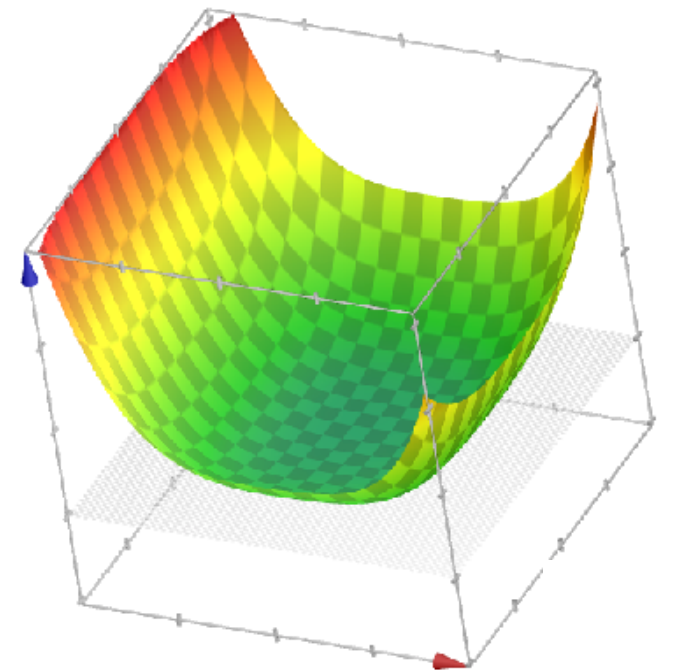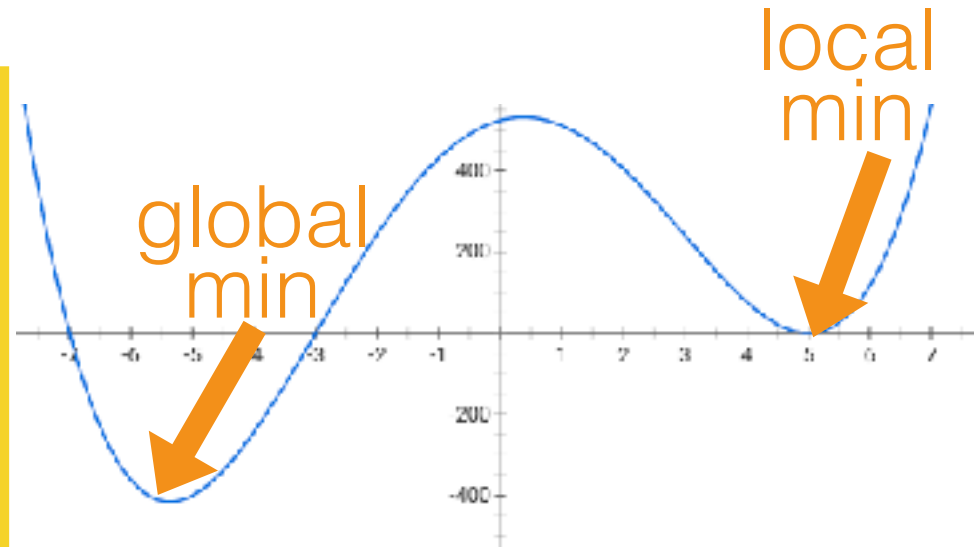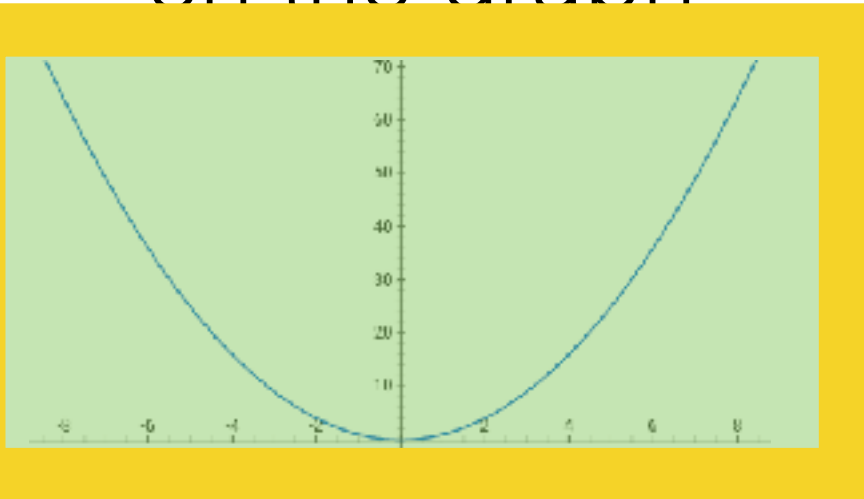


global min

local min



- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
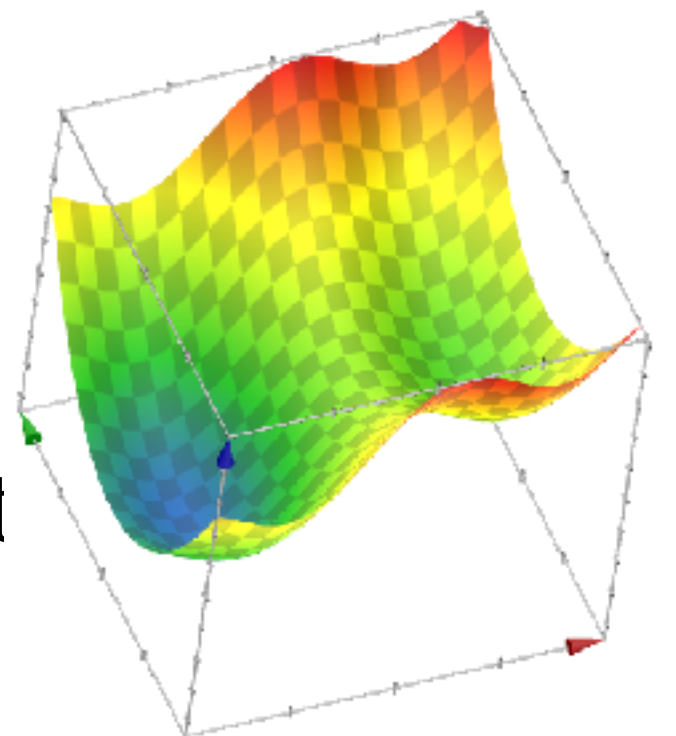  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
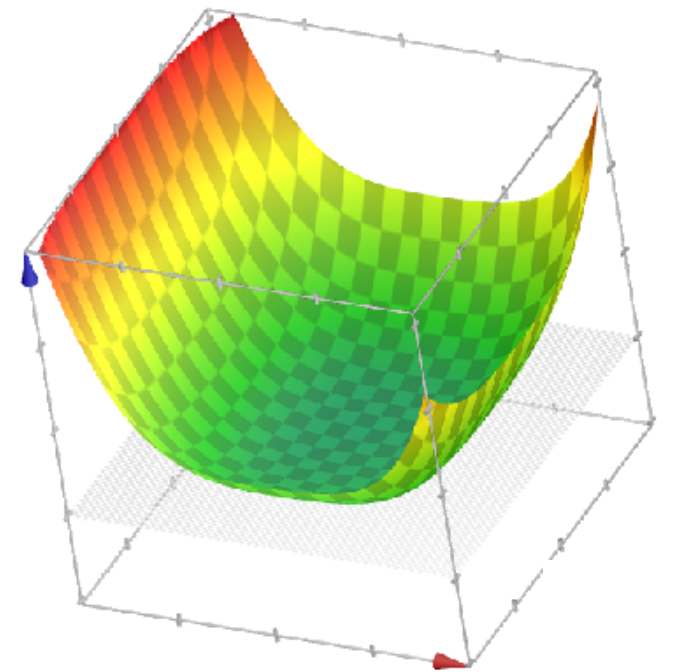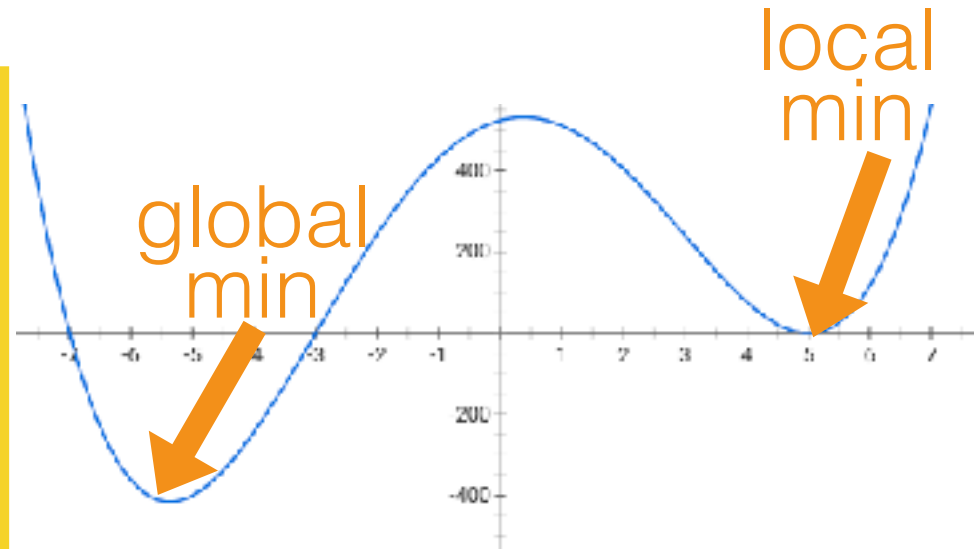


global min

local min

- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
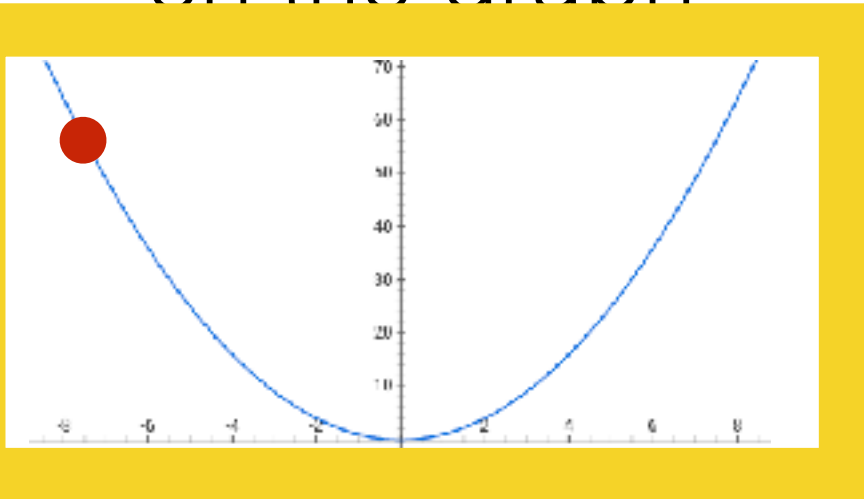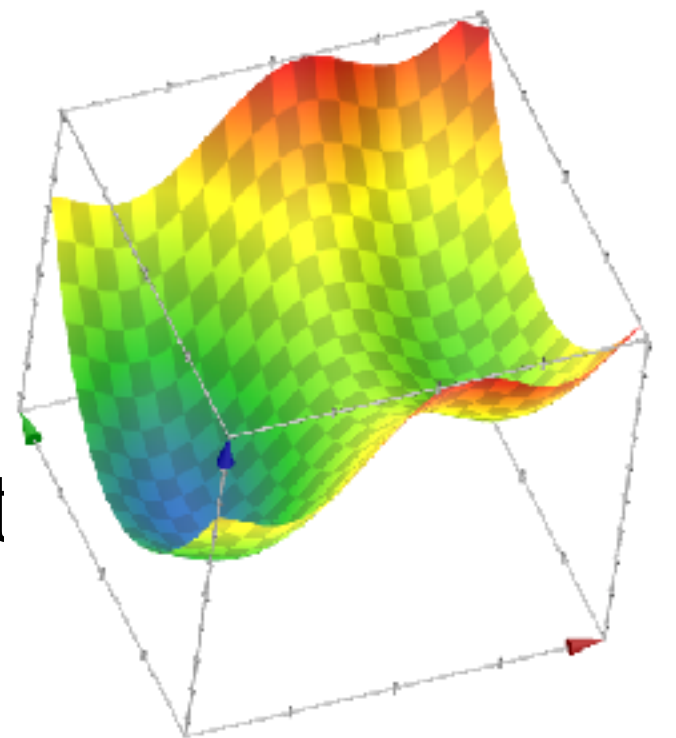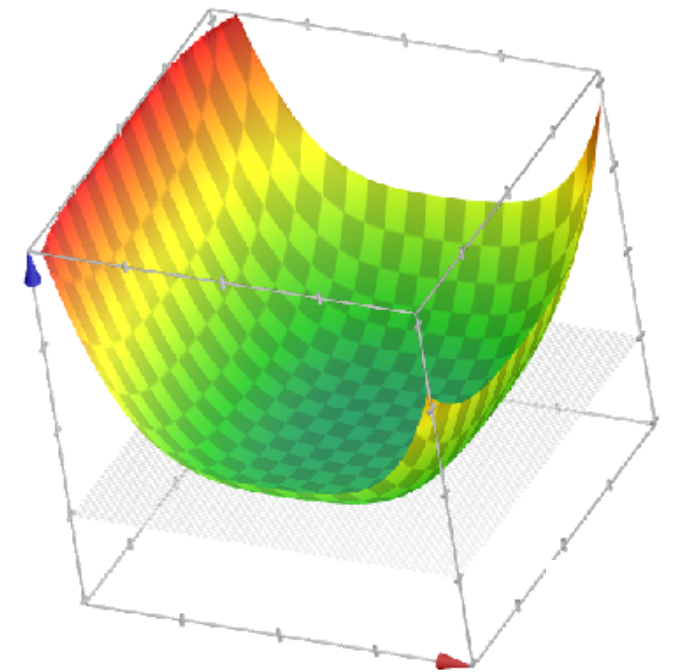
local
min

global
min

- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
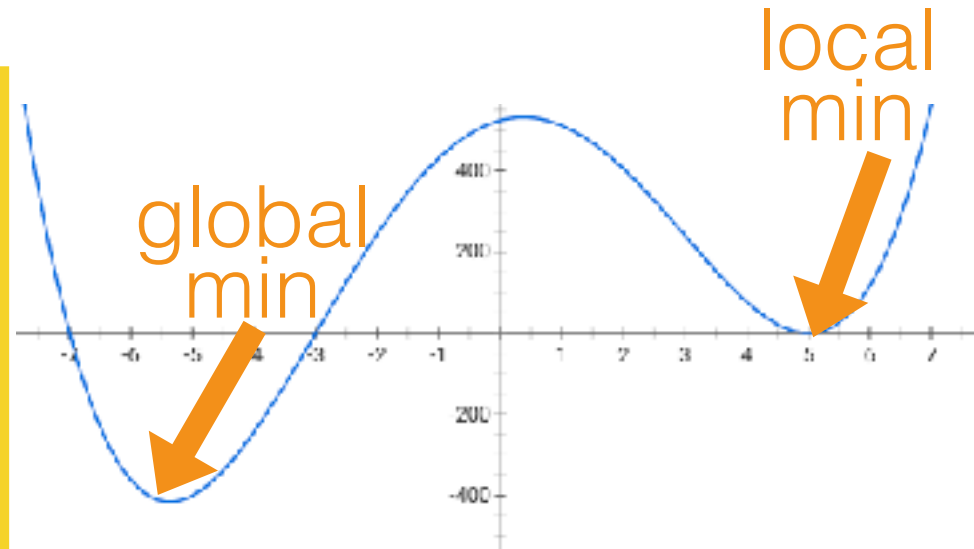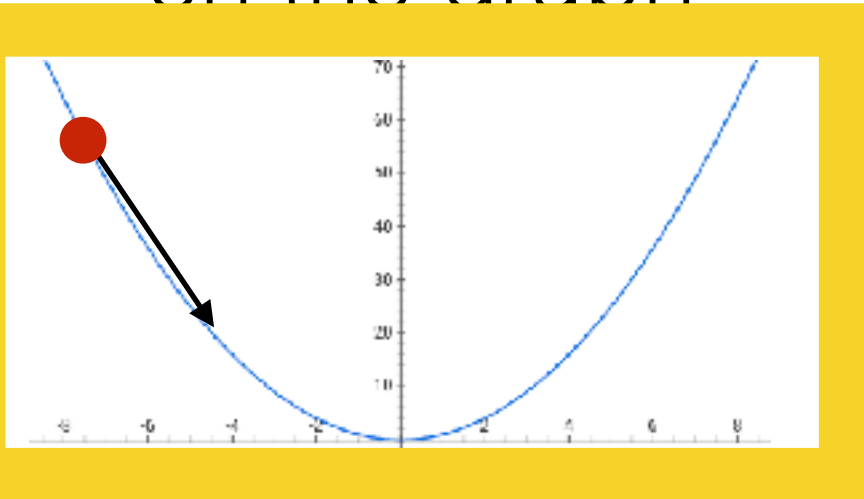


global min

local min

- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
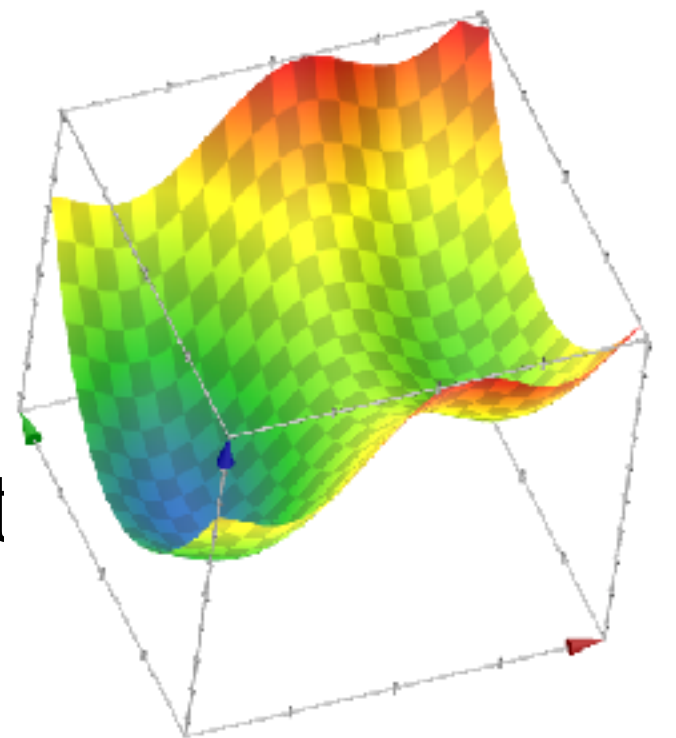  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

5

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
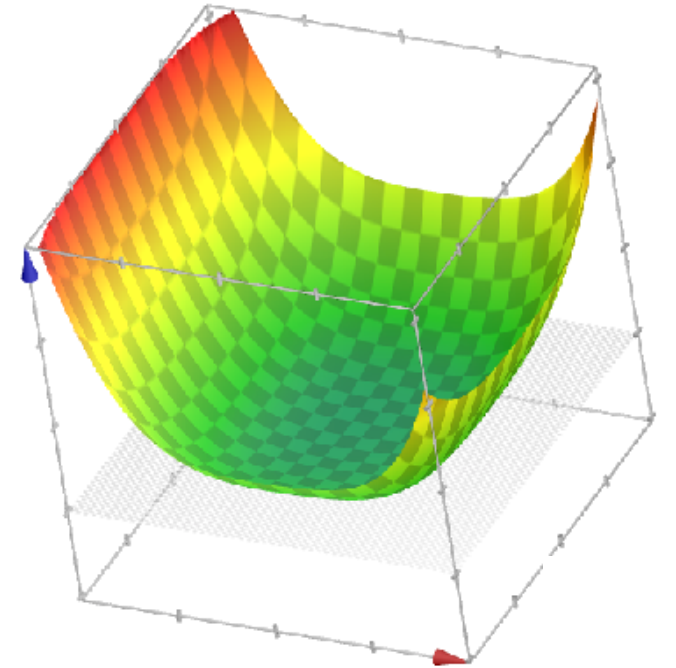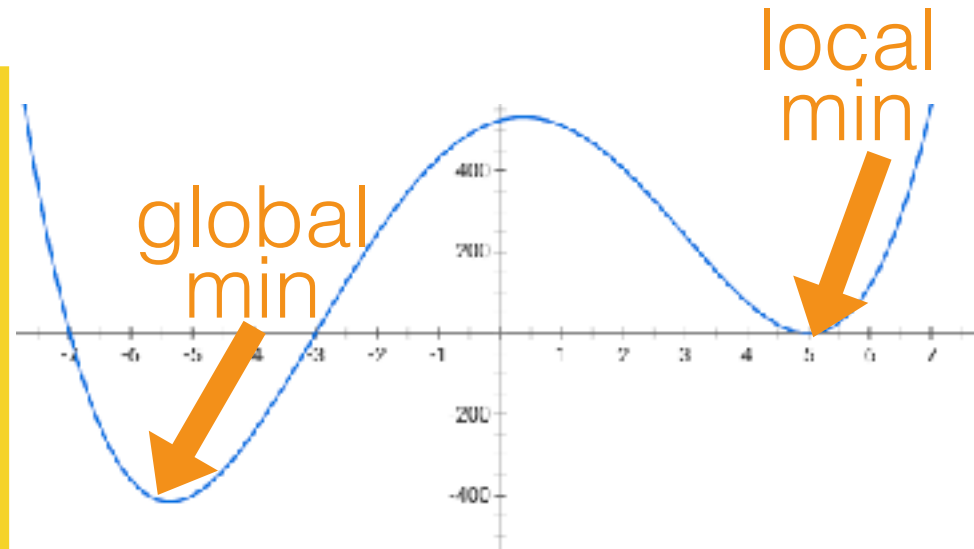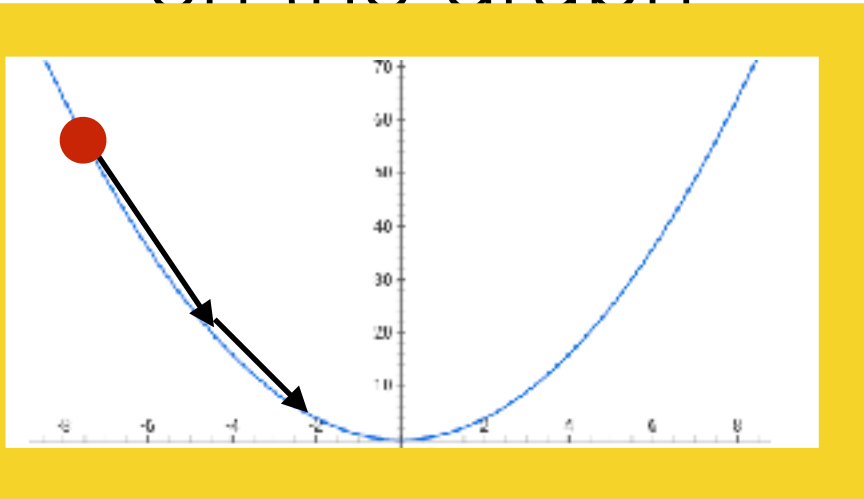


global min

local min

- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
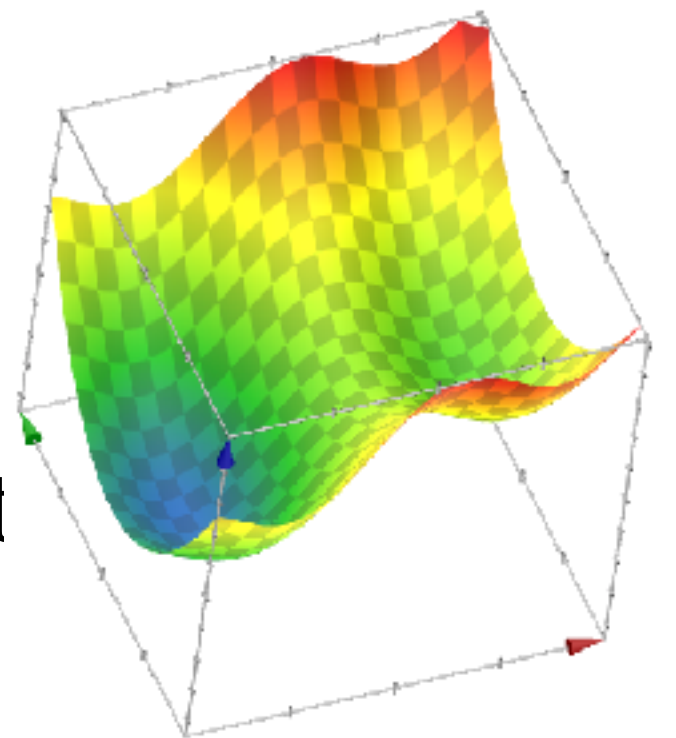  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
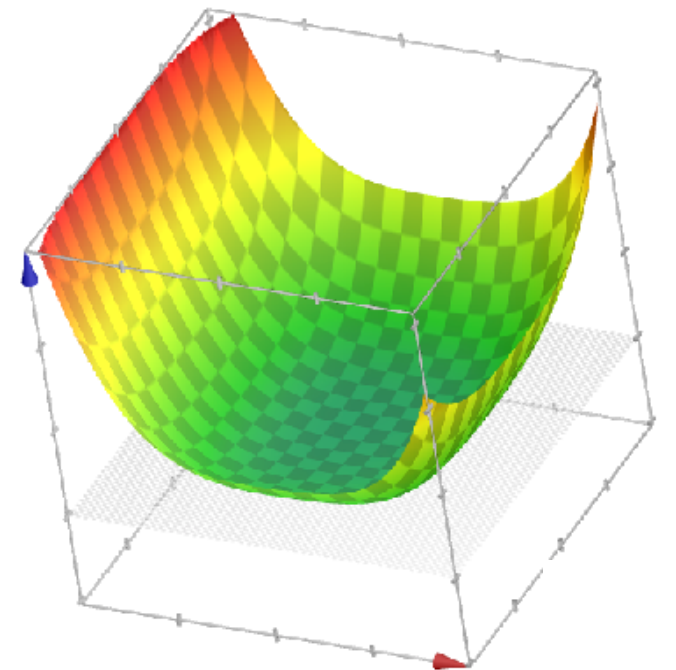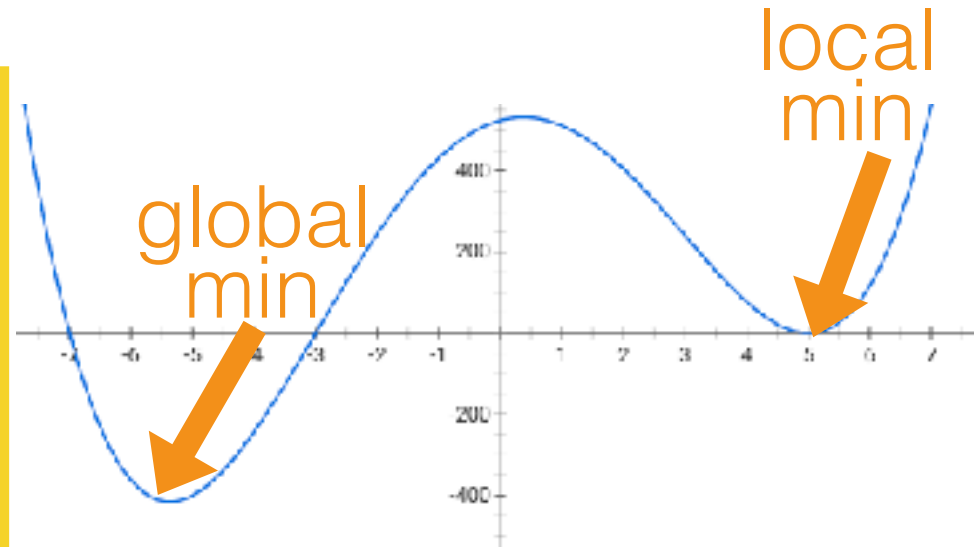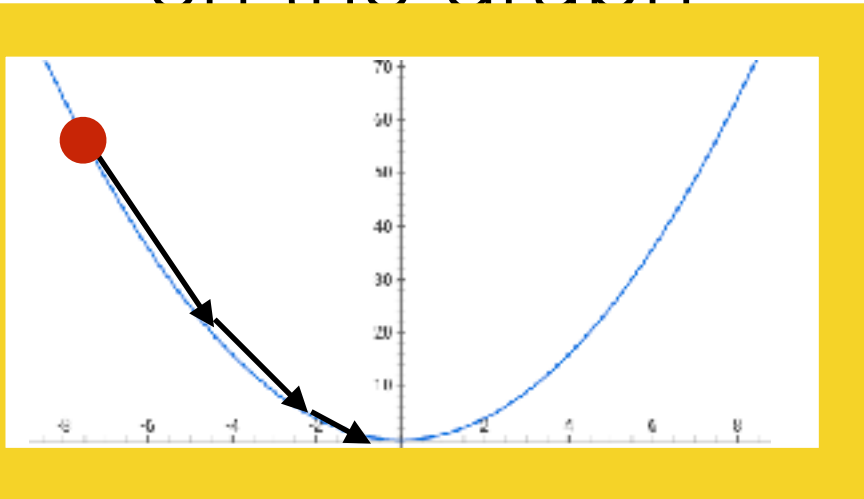
local min

global min

- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
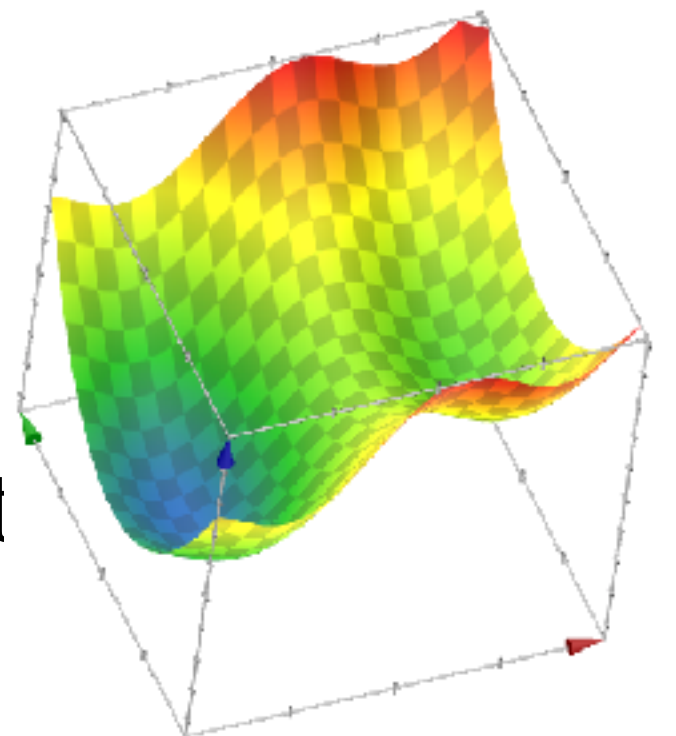  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
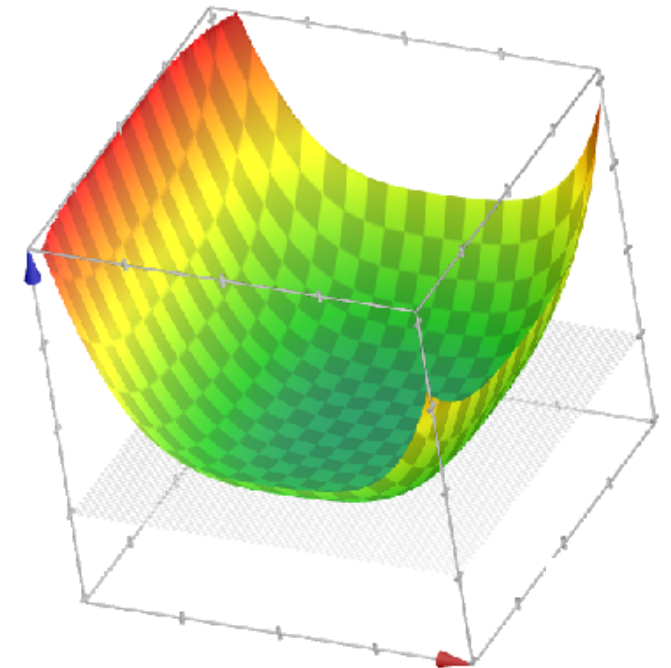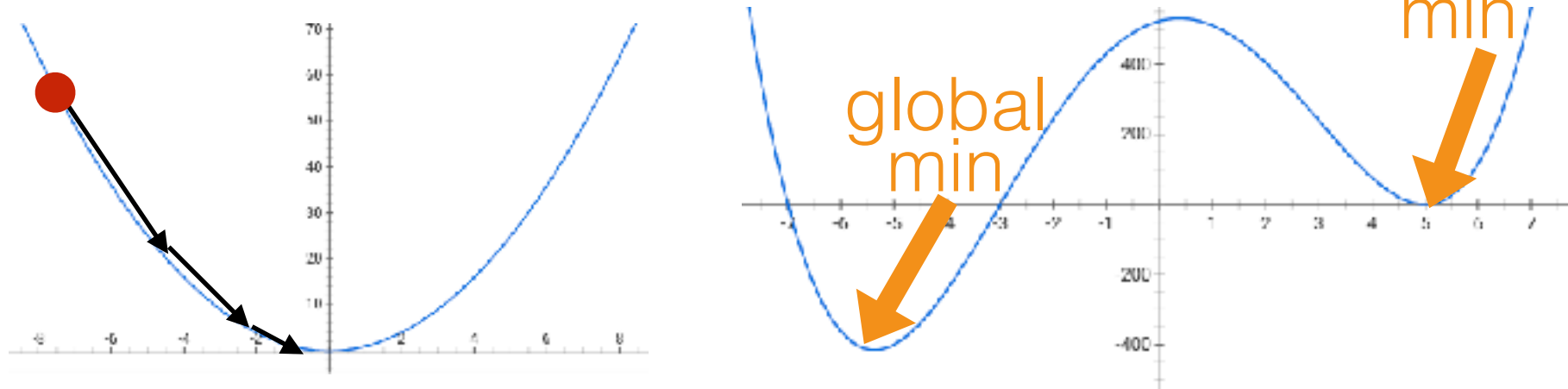
global min

local min

- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
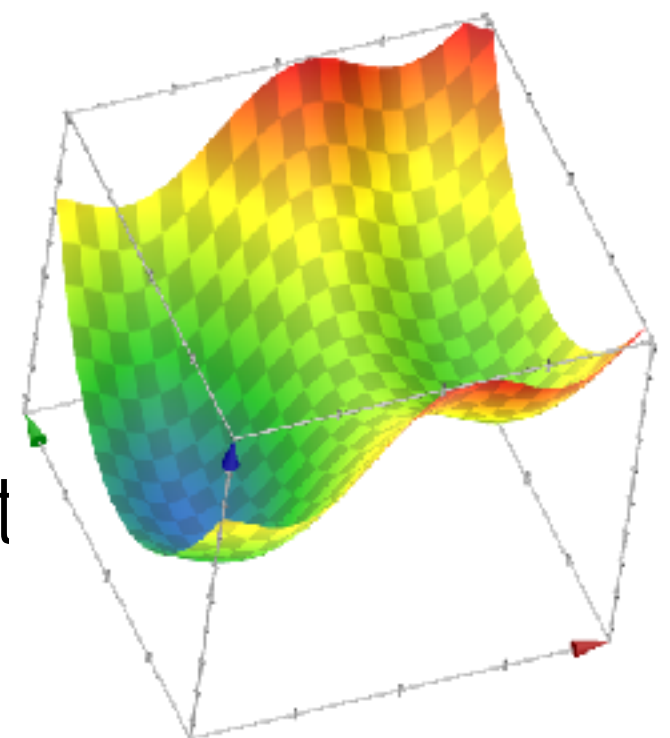  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
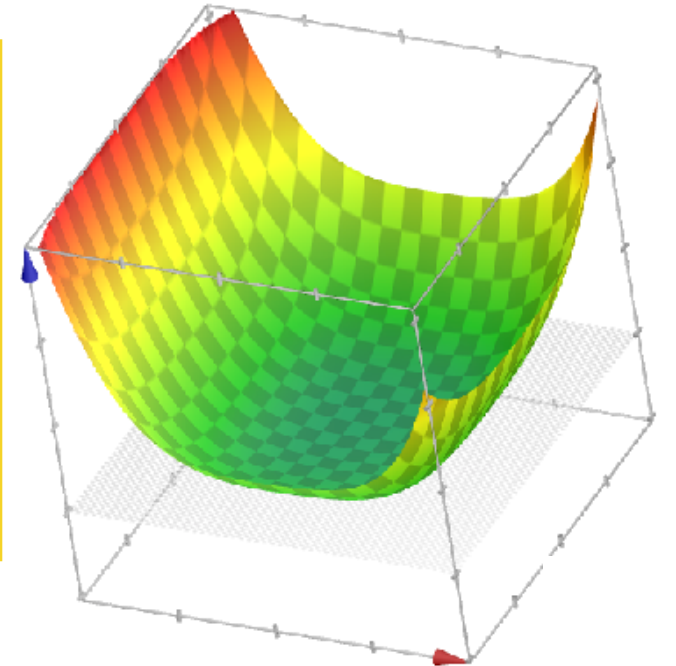






- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
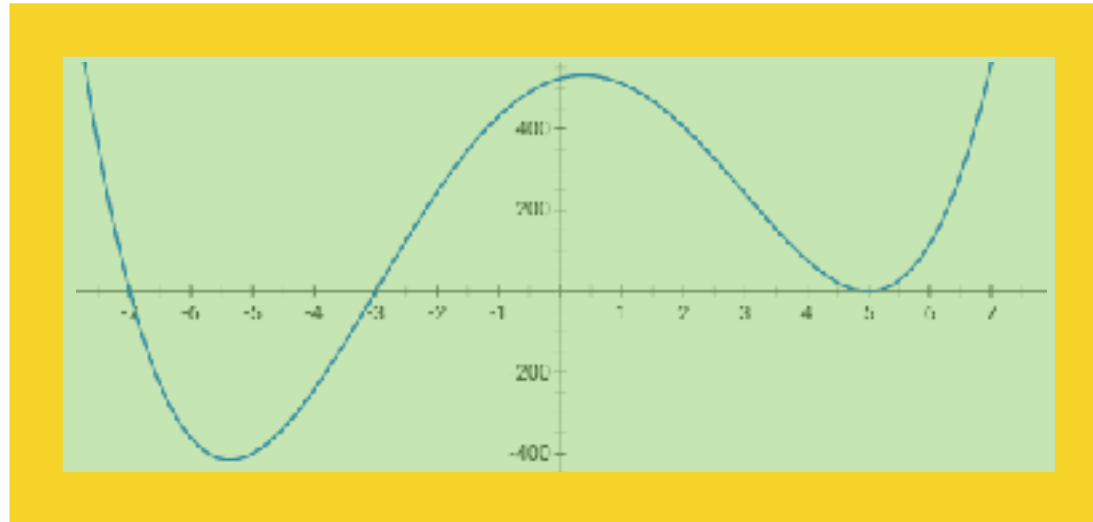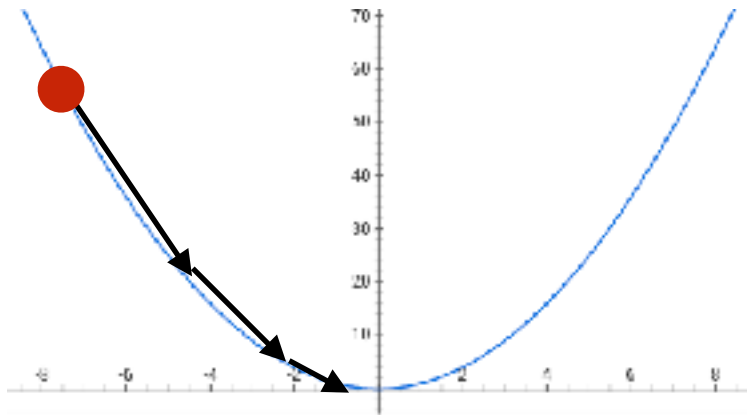
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
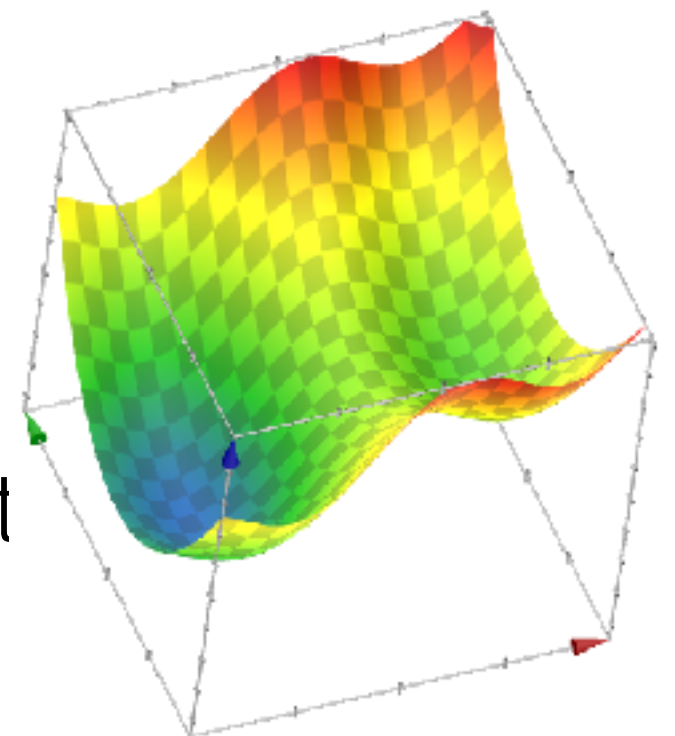  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
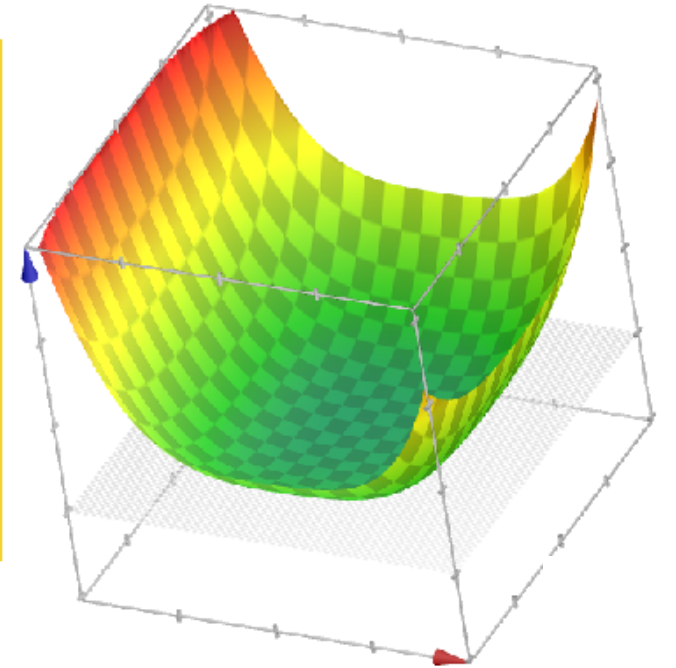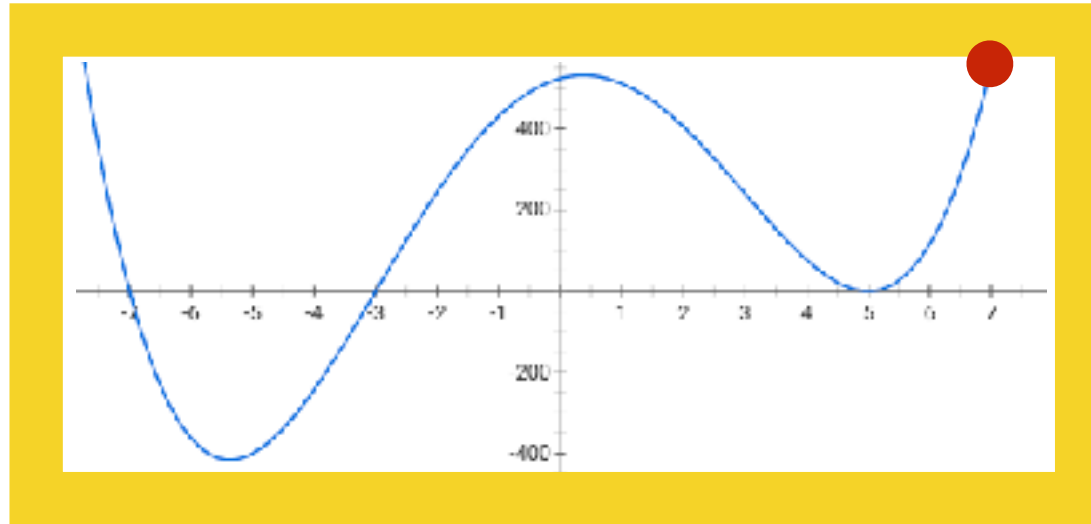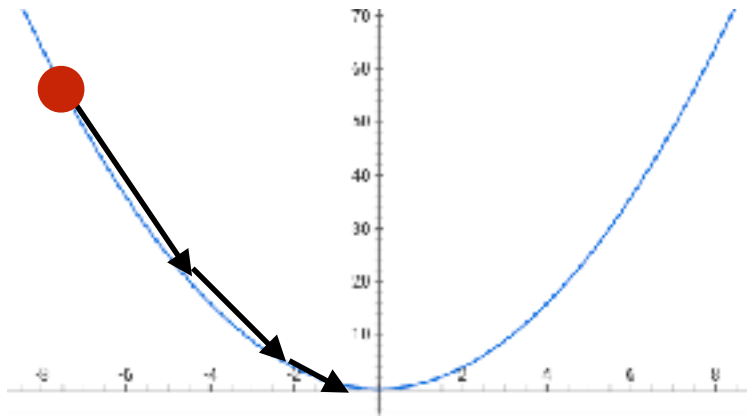


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
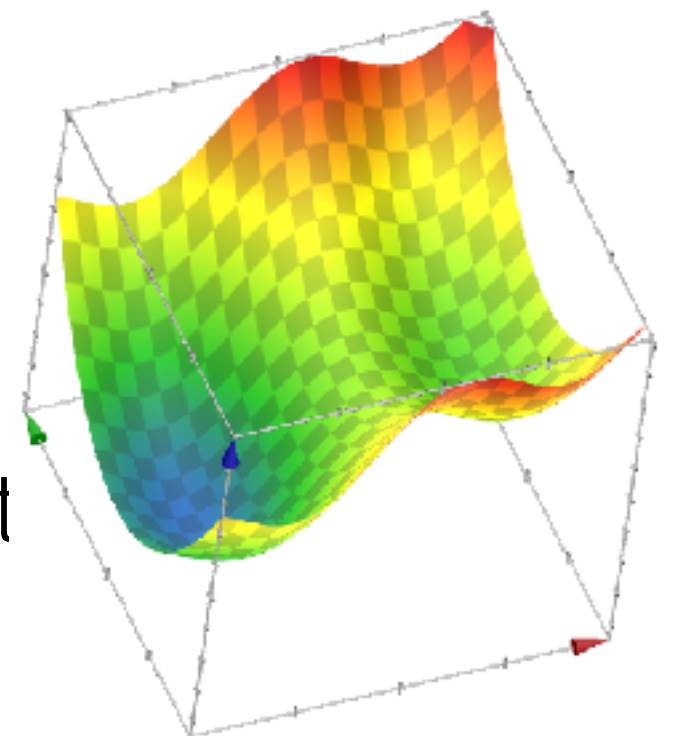  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
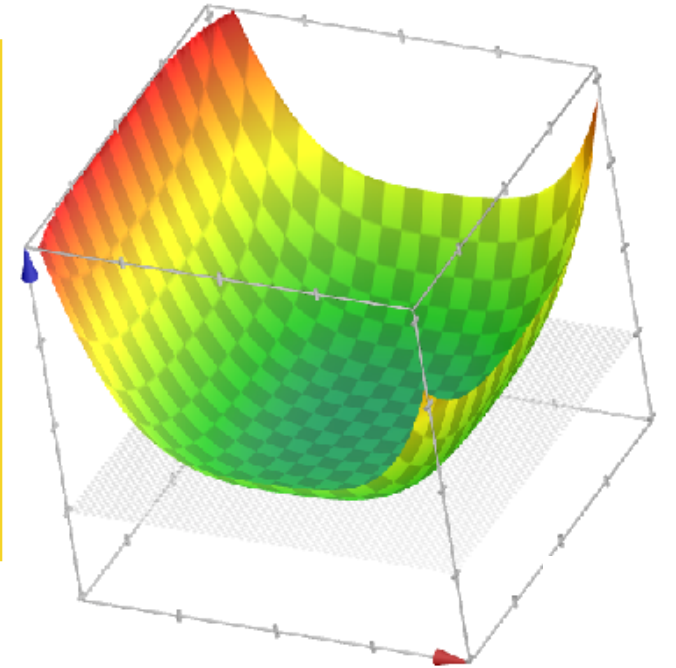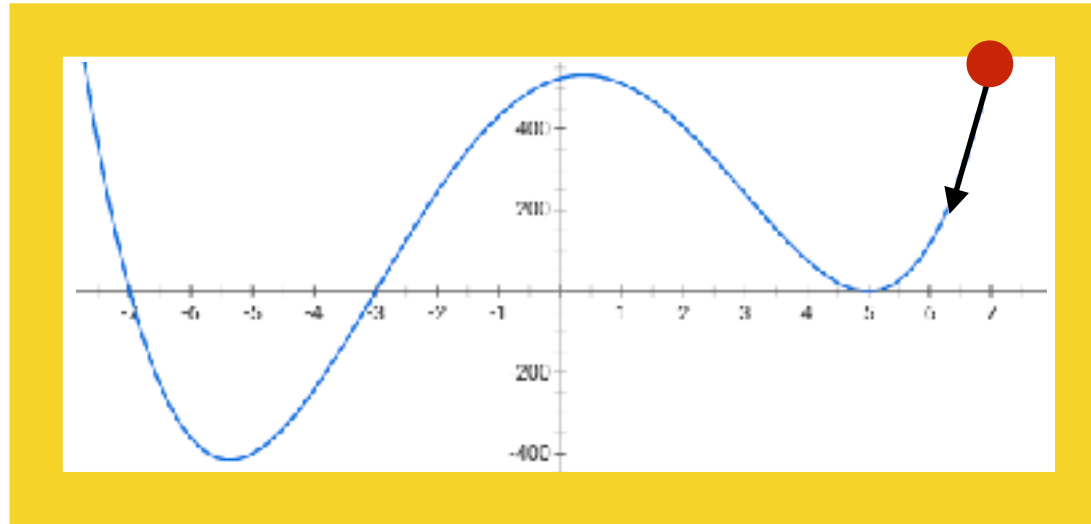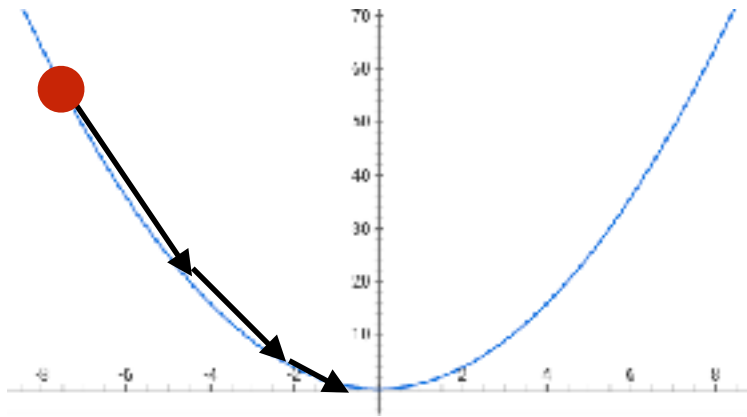


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
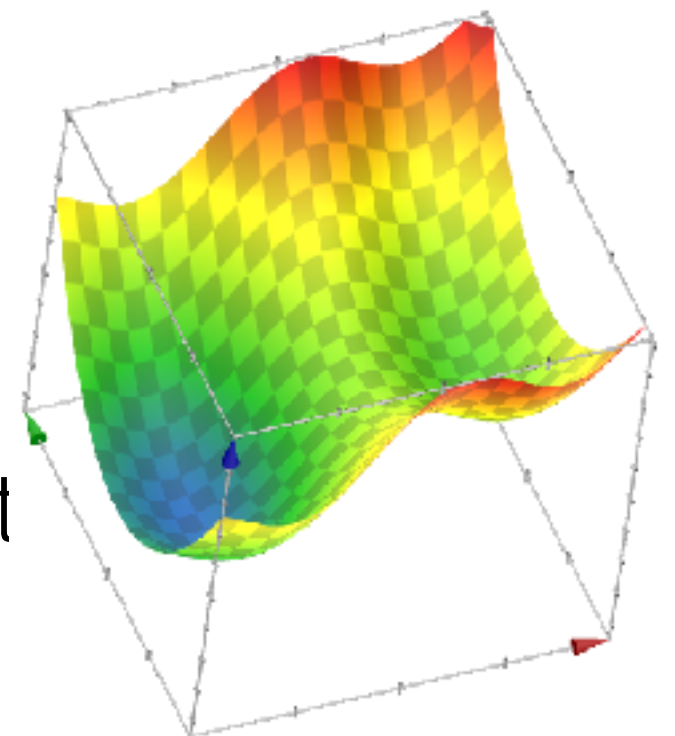  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
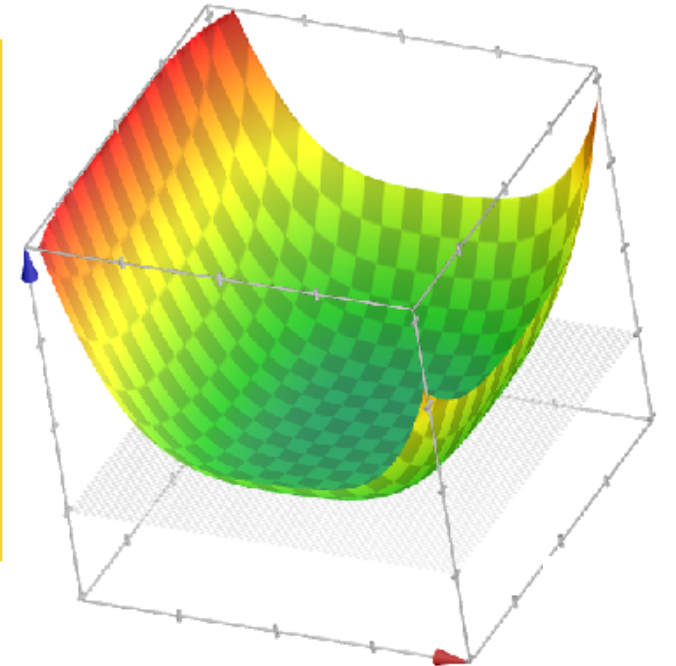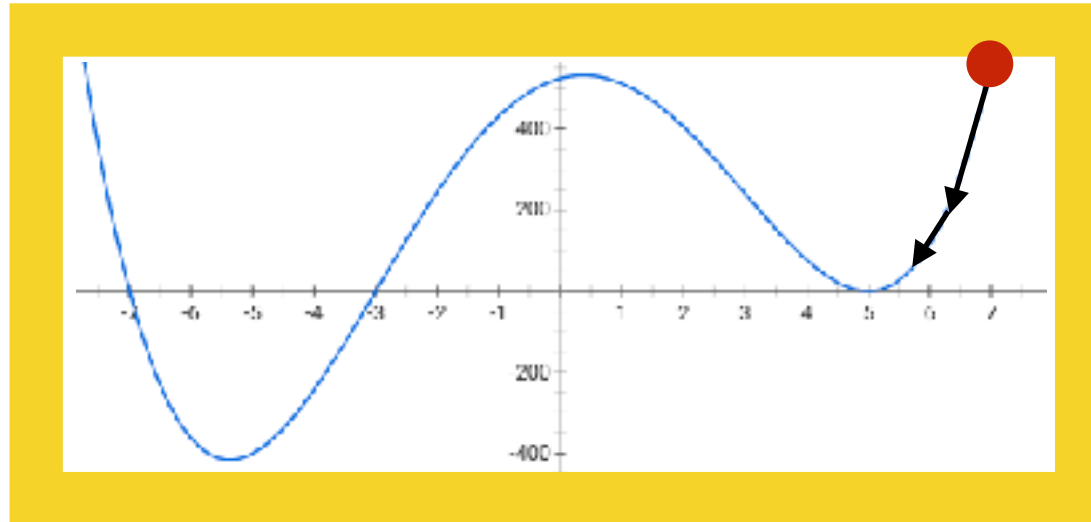


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
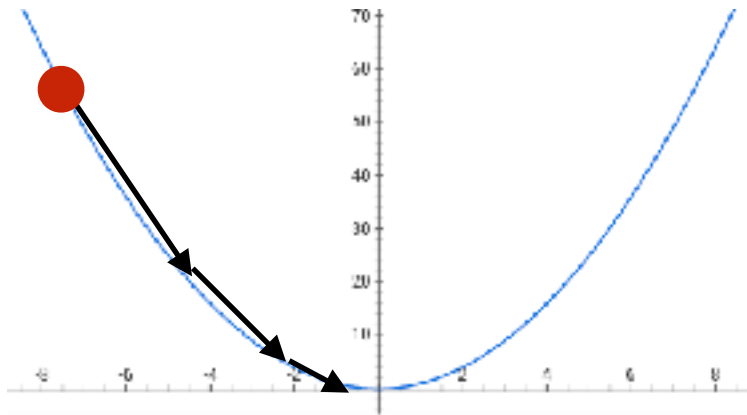


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
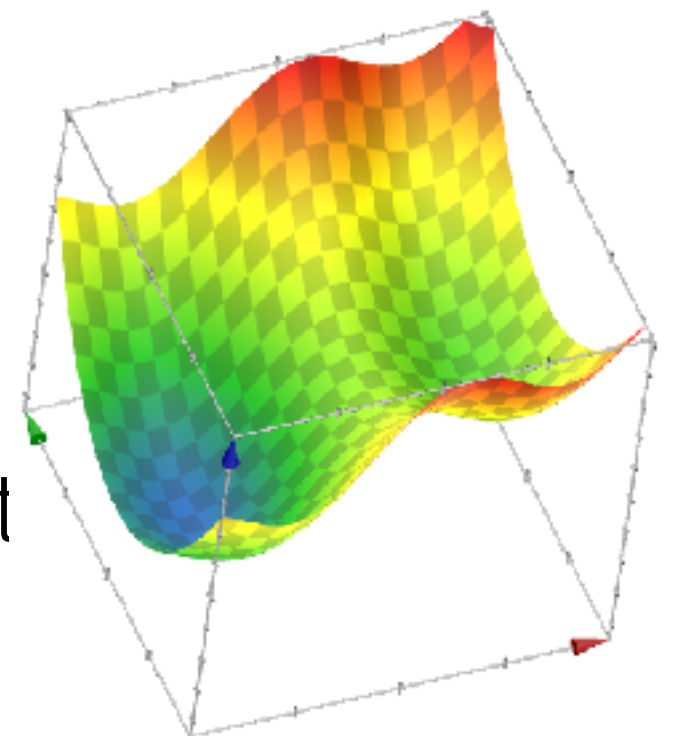  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
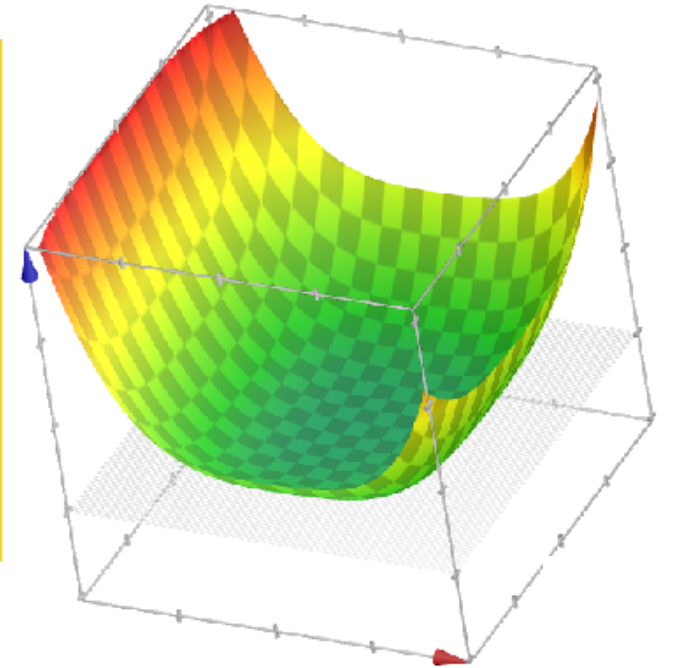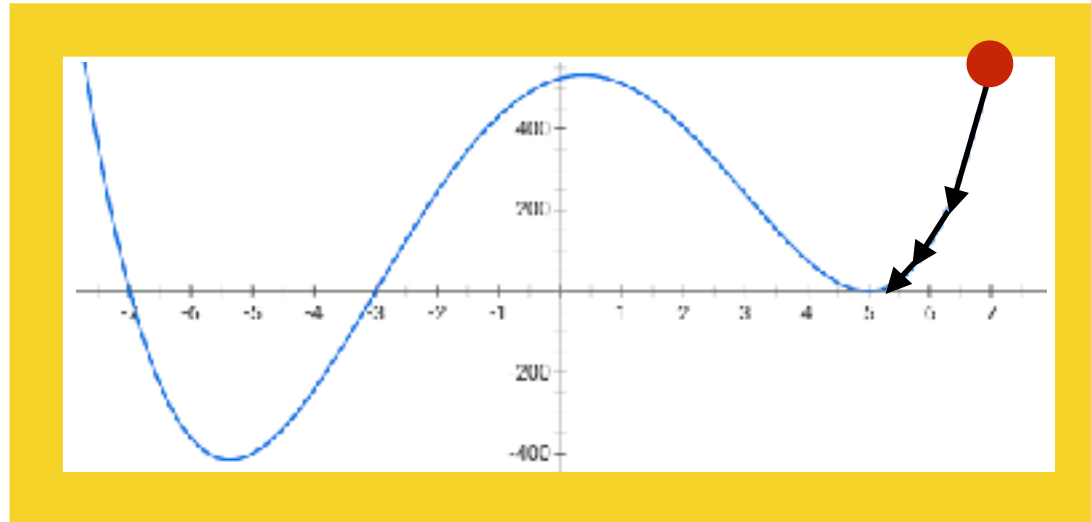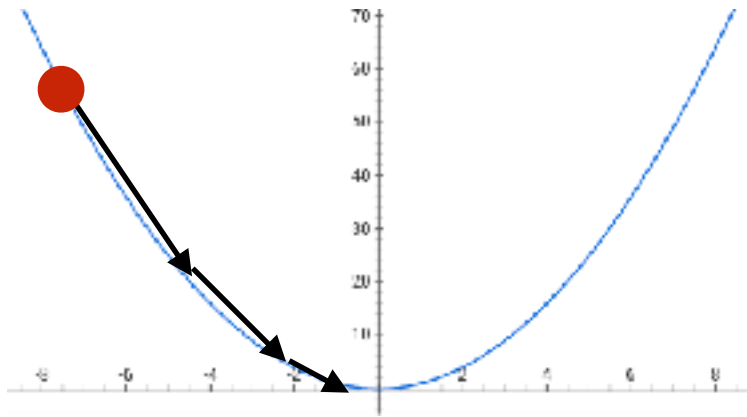


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
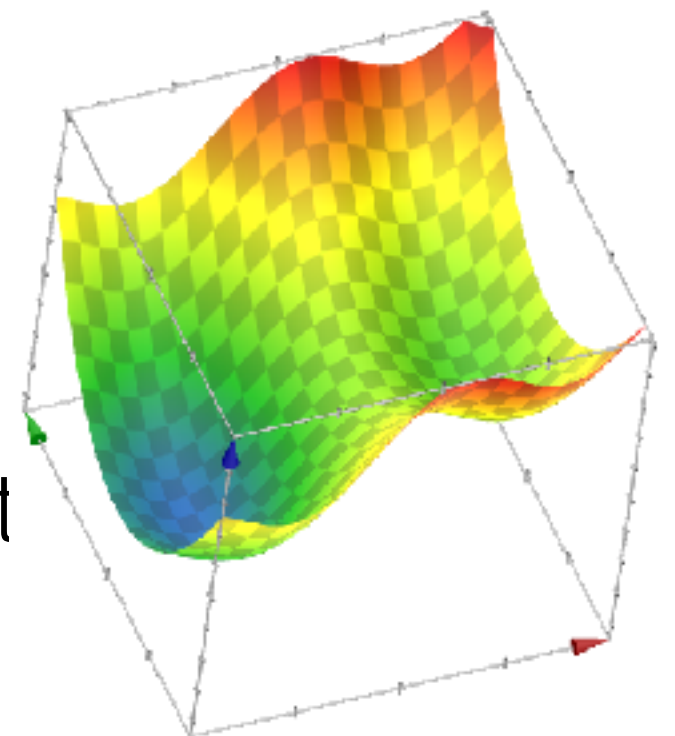  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function *f* on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of *f* lies above or on the graph
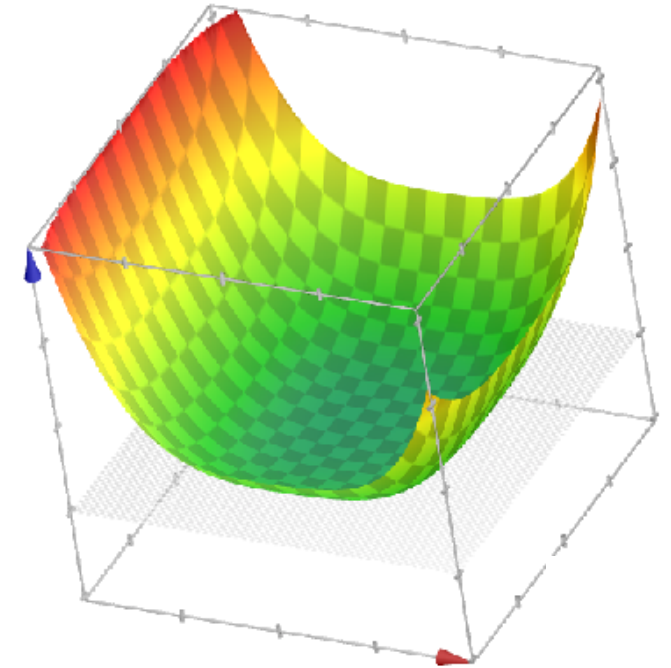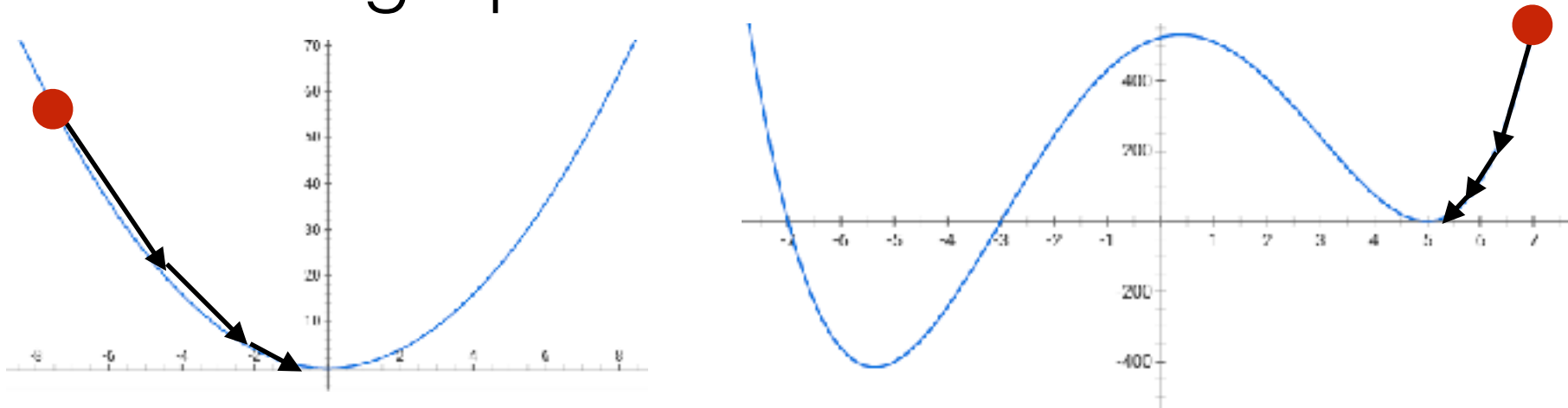




- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - *f* is sufficiently "smooth" and convex
    - *f* has at least one global optimum
    - $\eta$ is sufficiently small
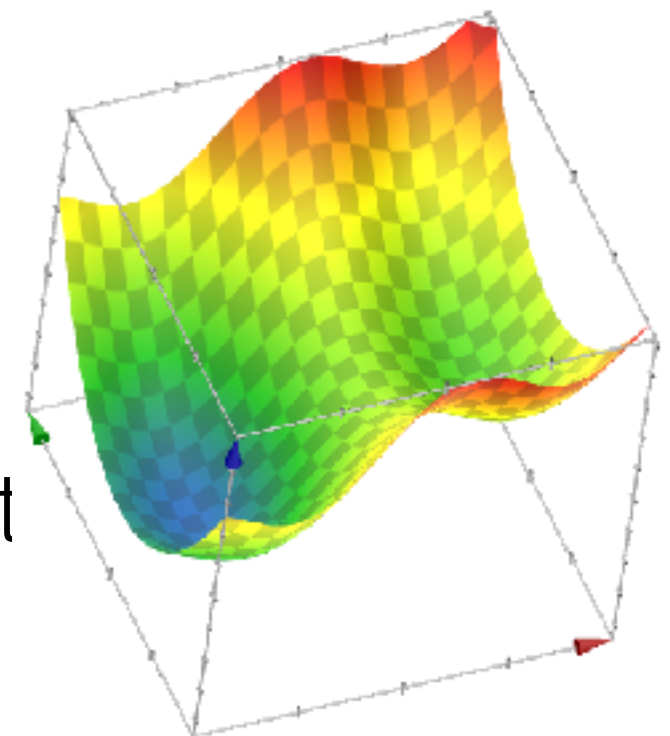  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
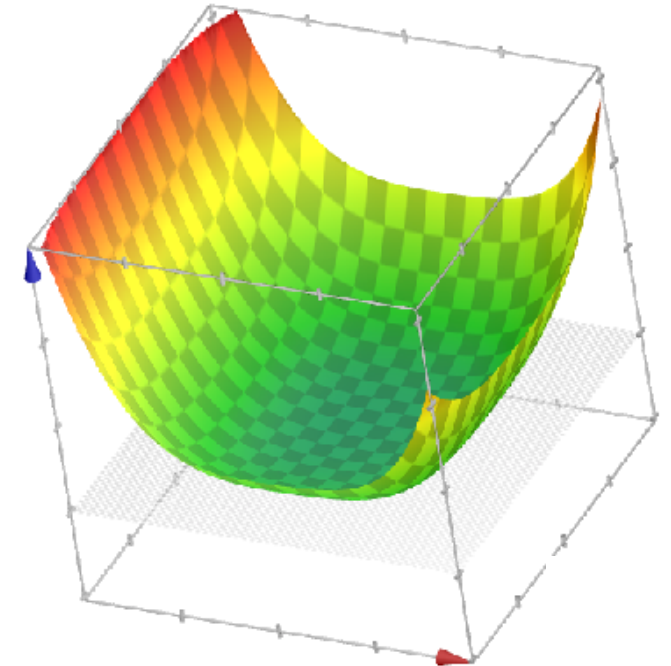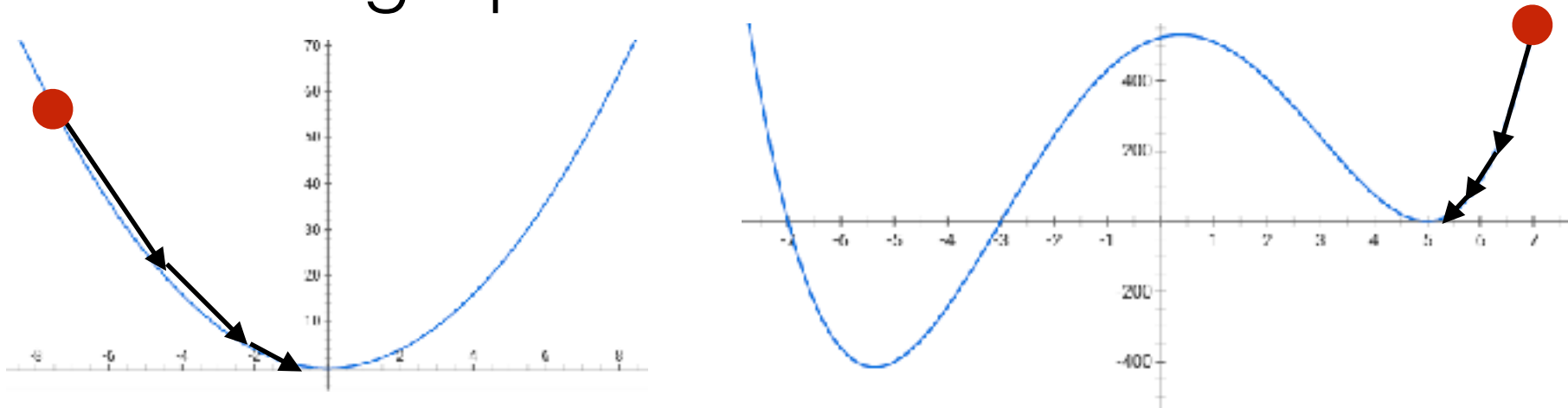
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
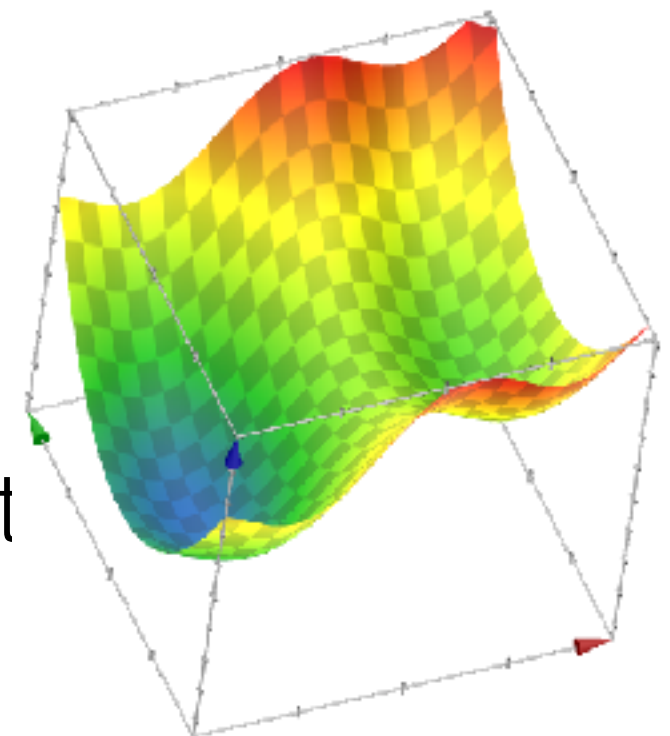  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
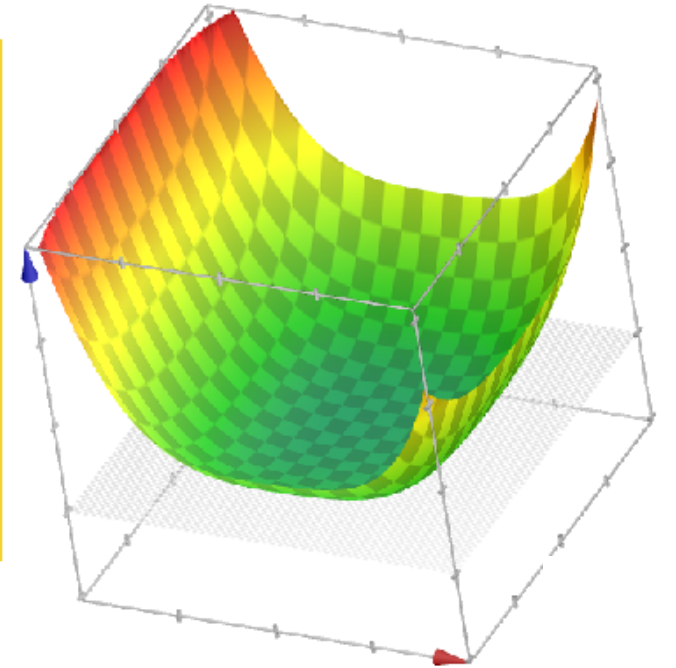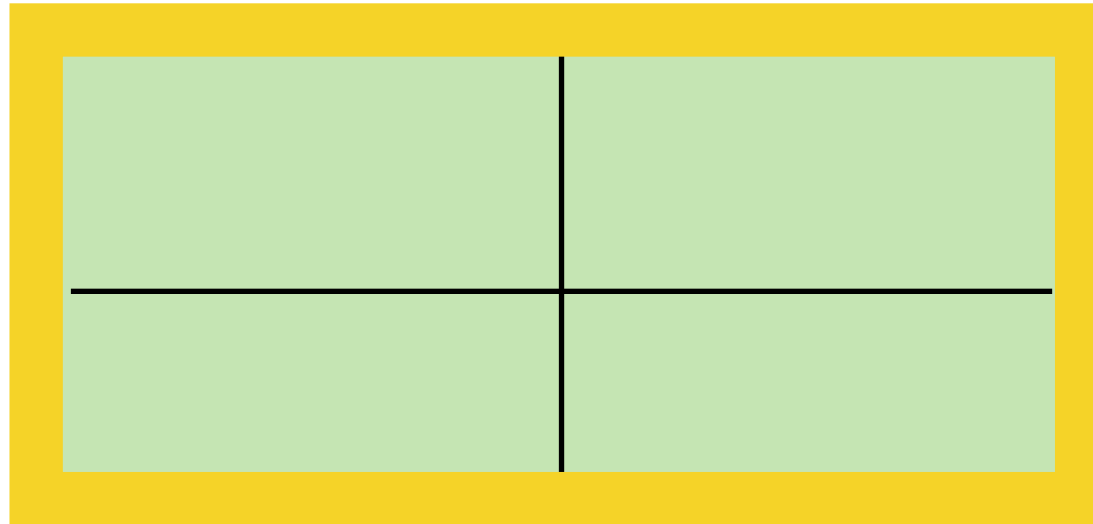


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
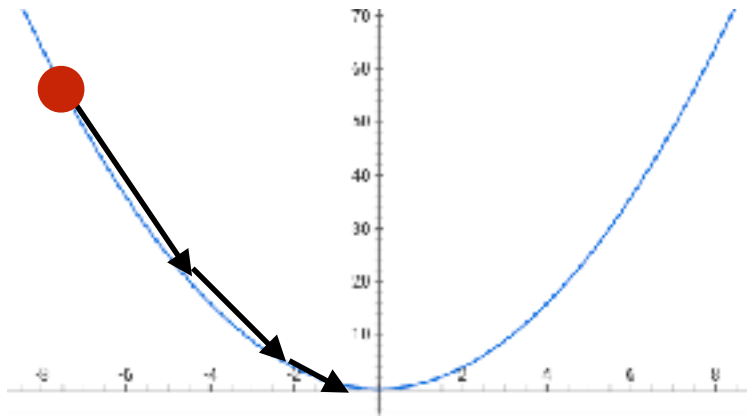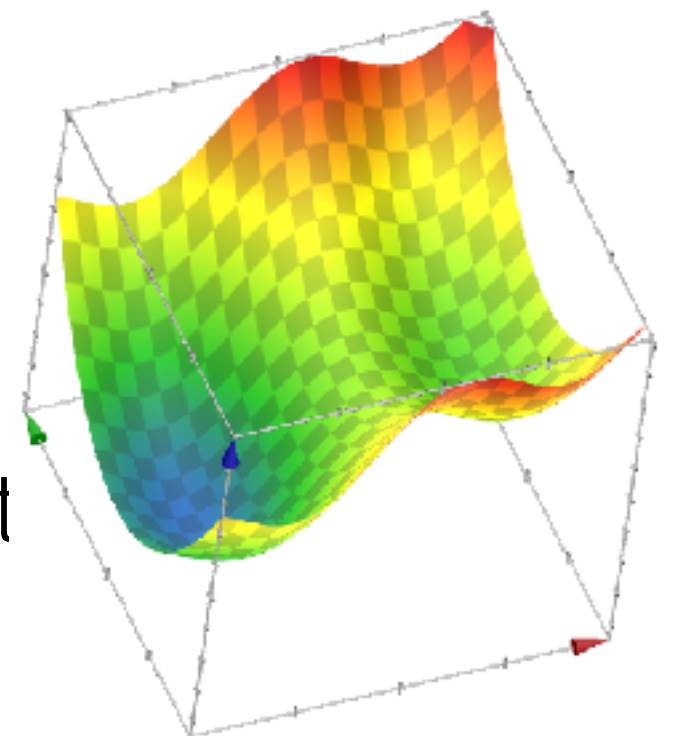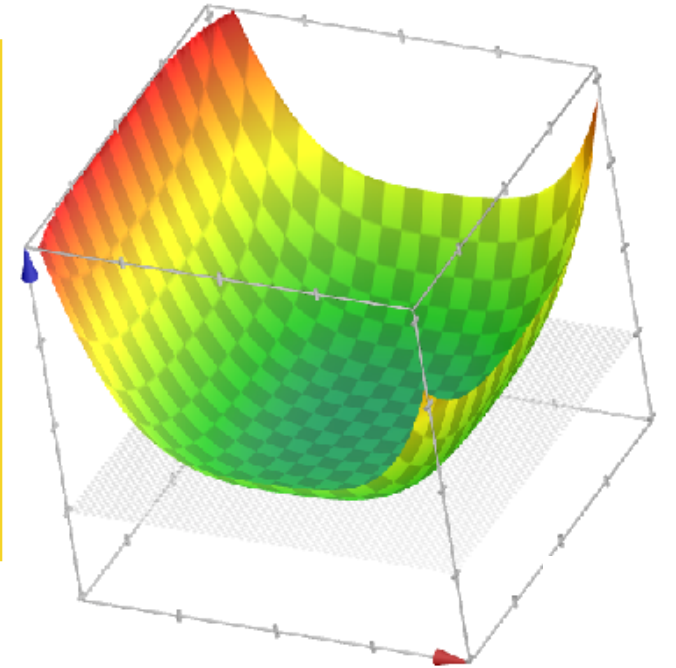
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
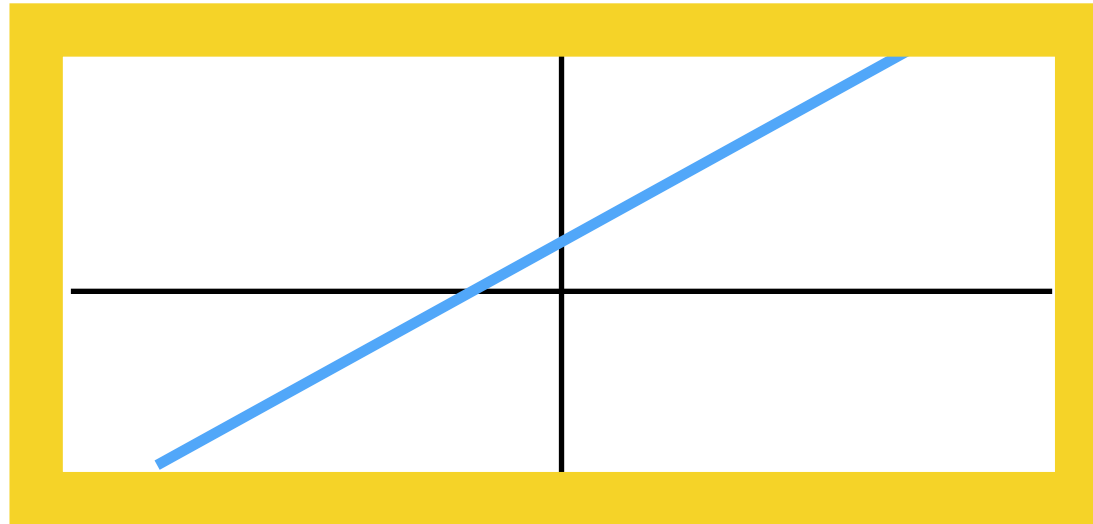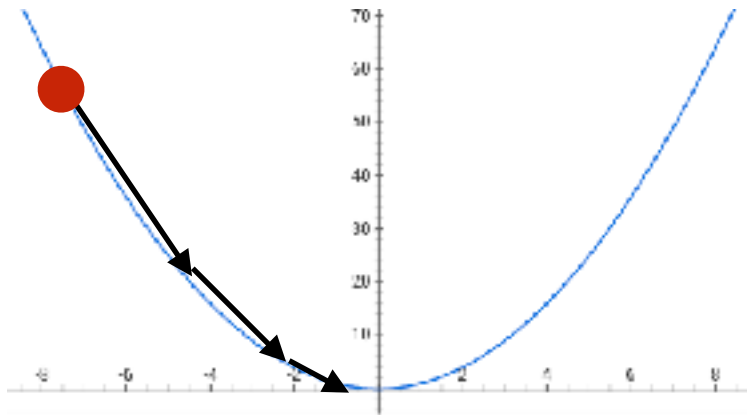


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
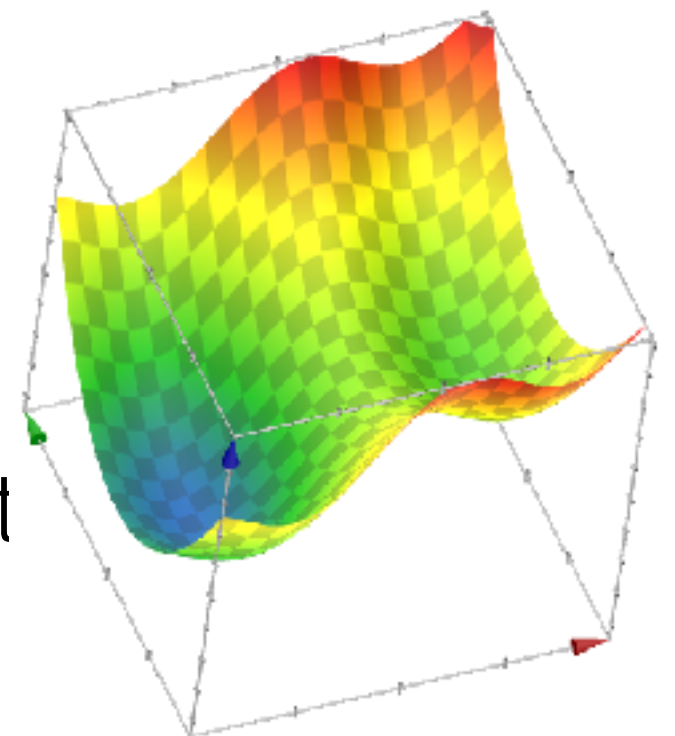  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

6

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
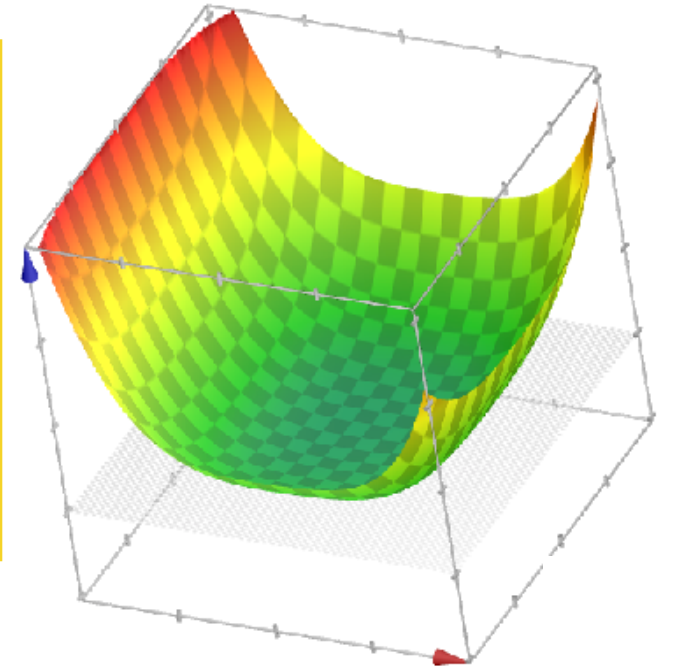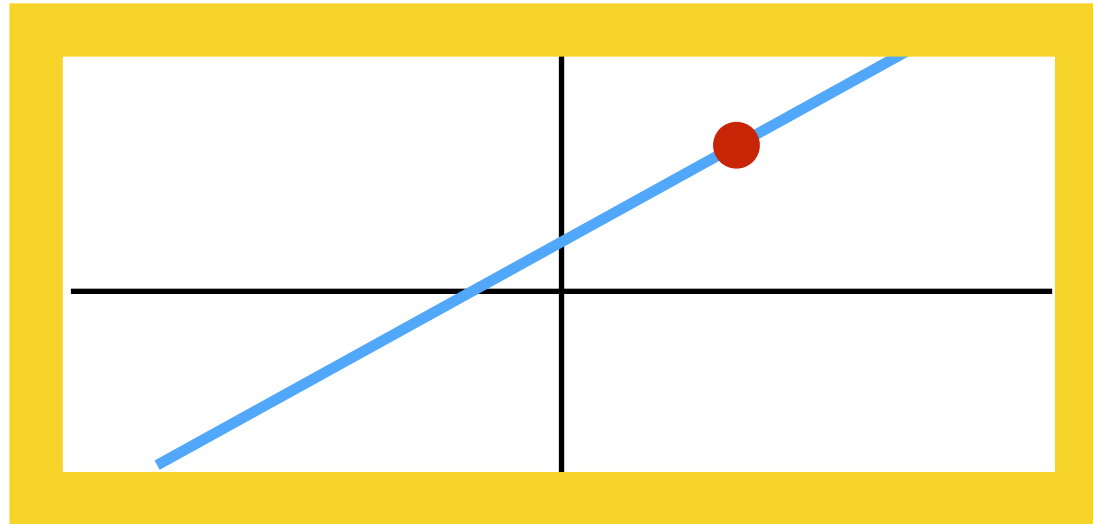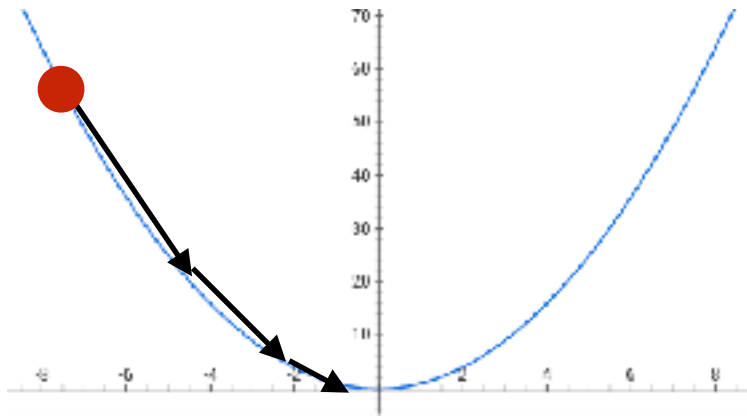




- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
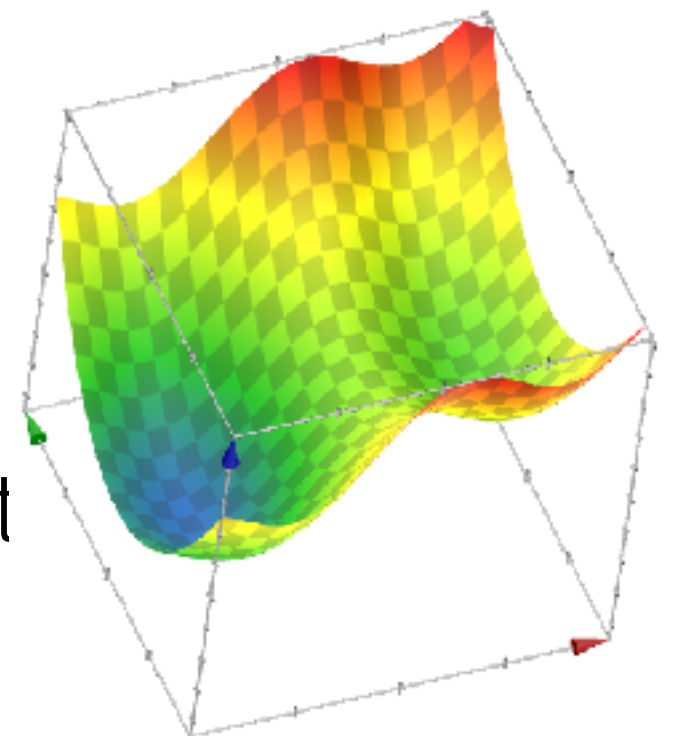  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
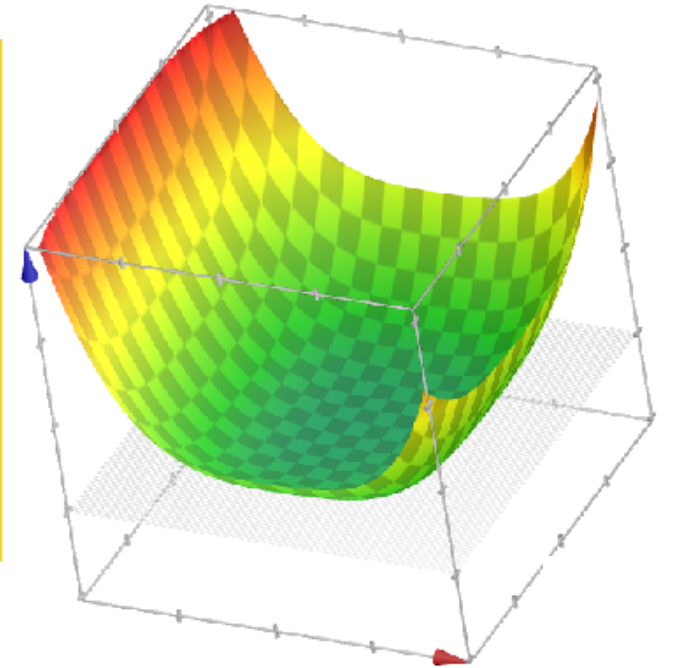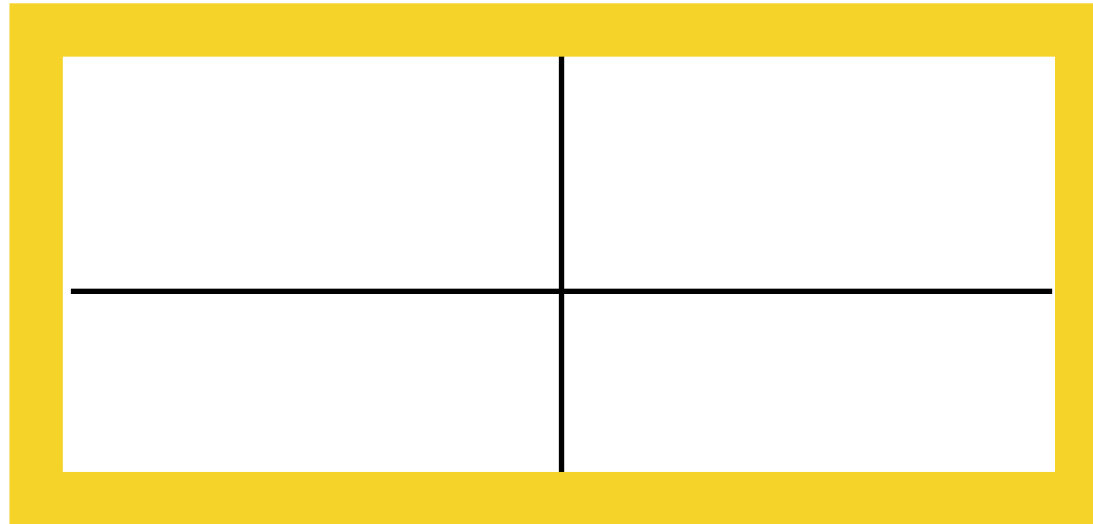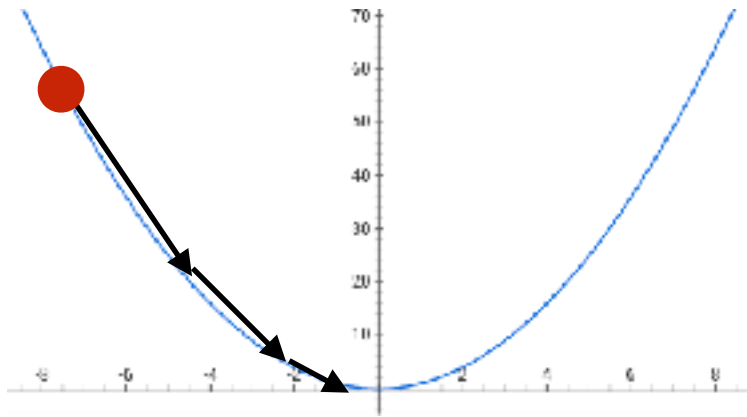


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
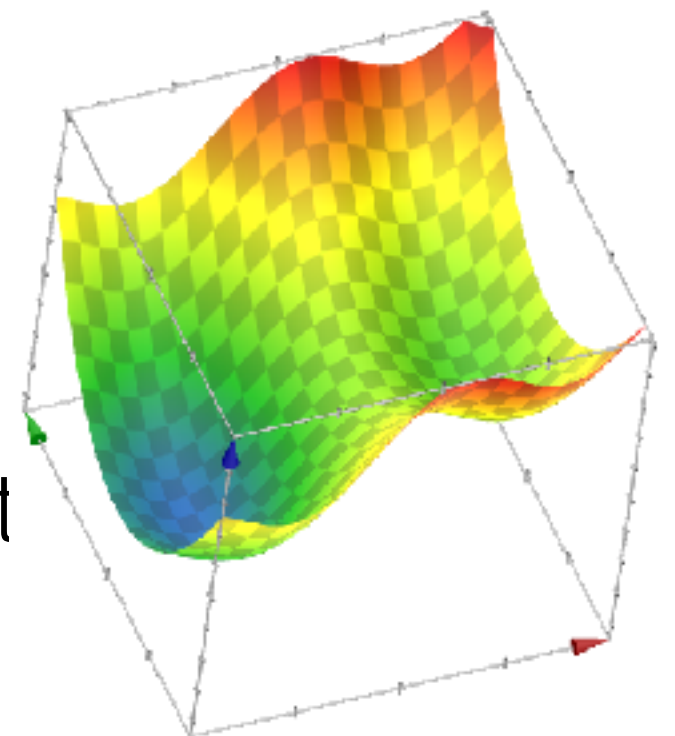  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
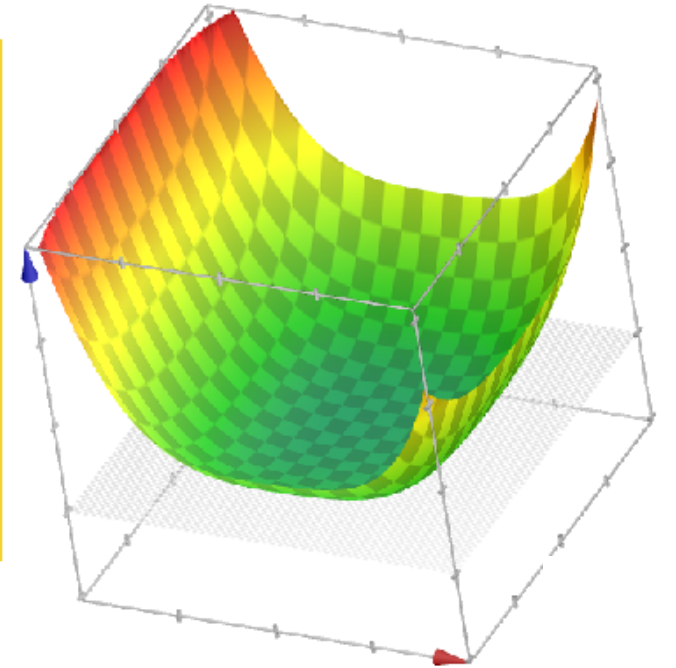
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

6

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
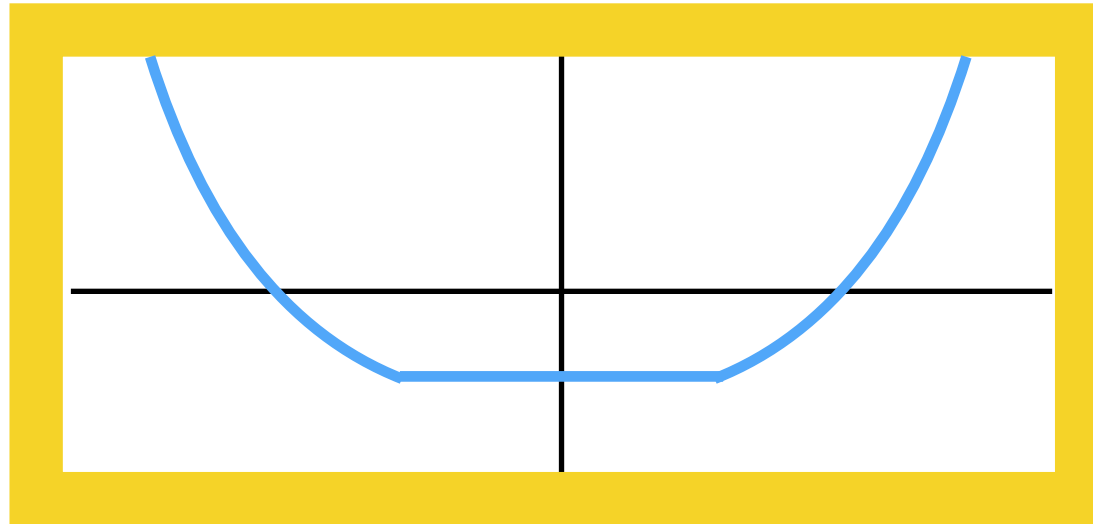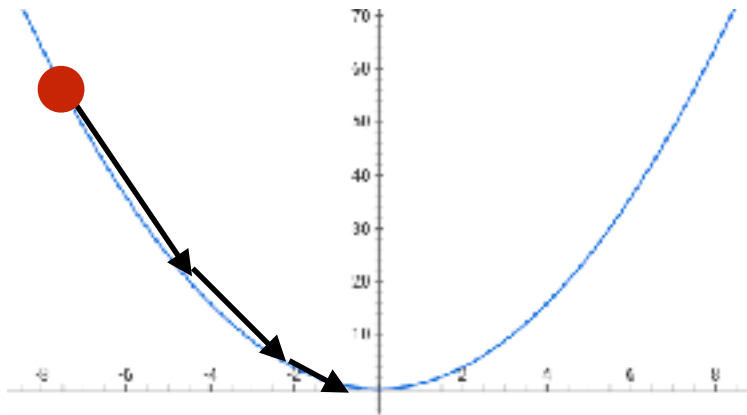


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
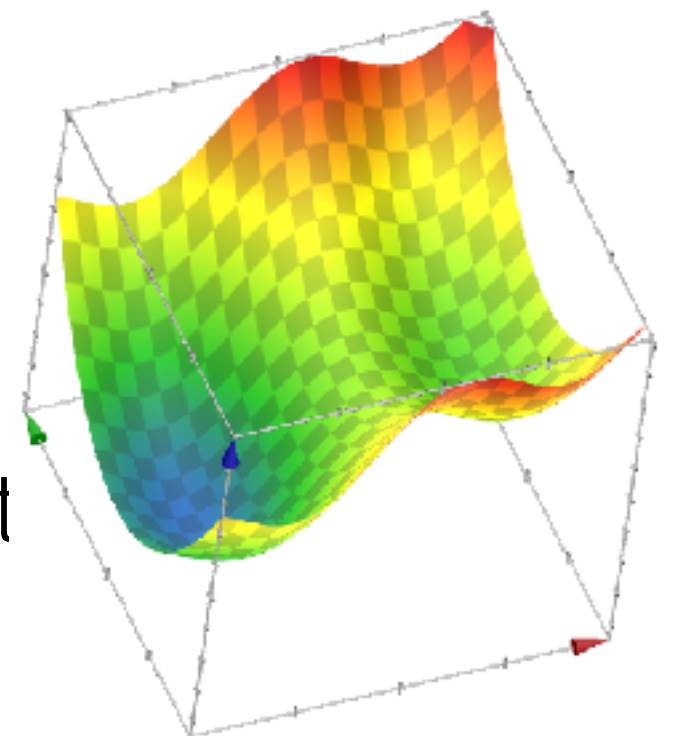  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
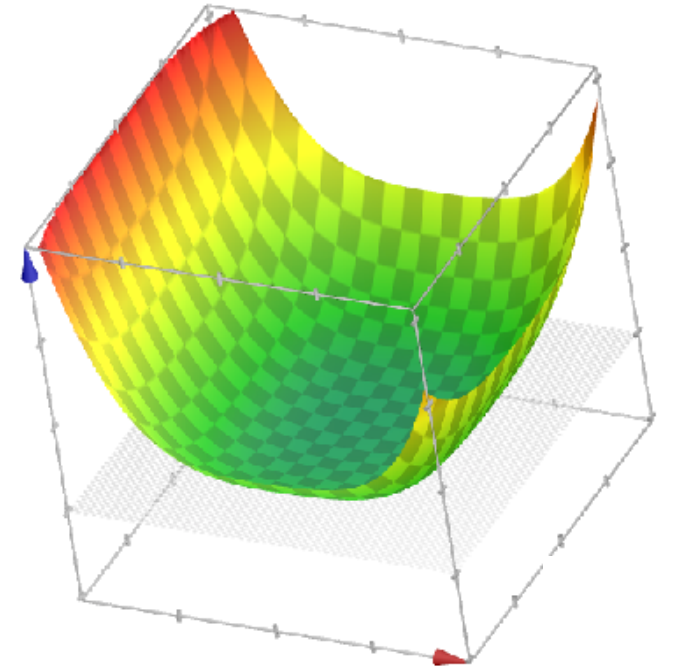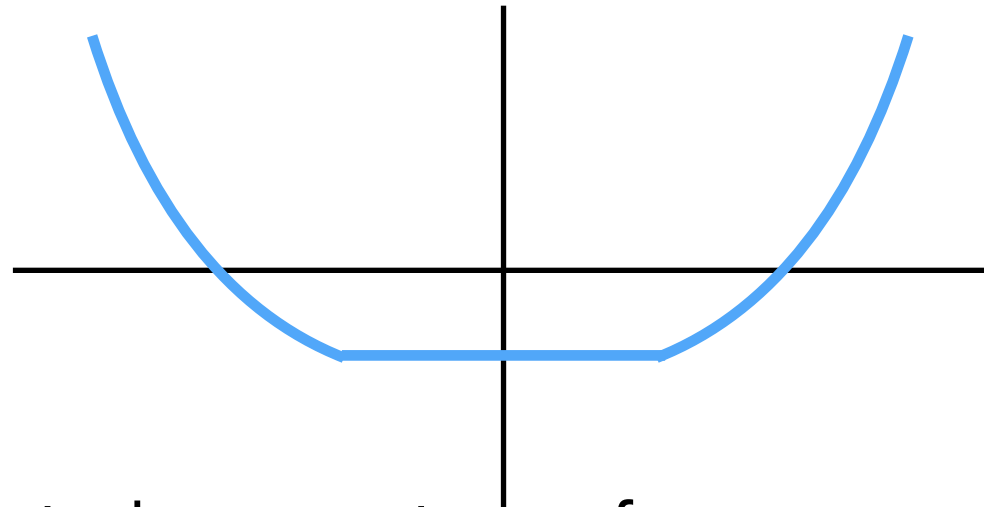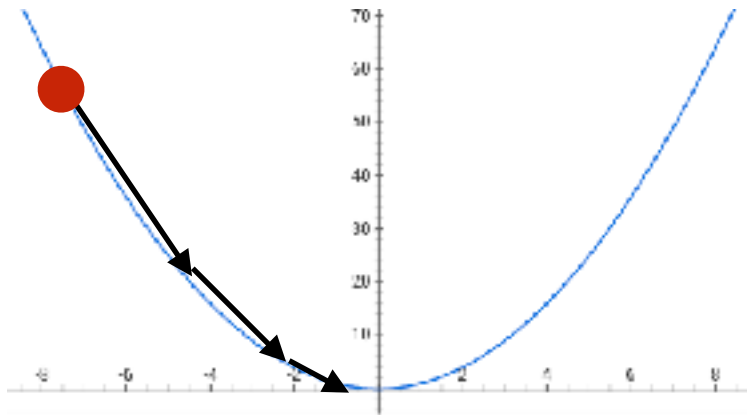


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
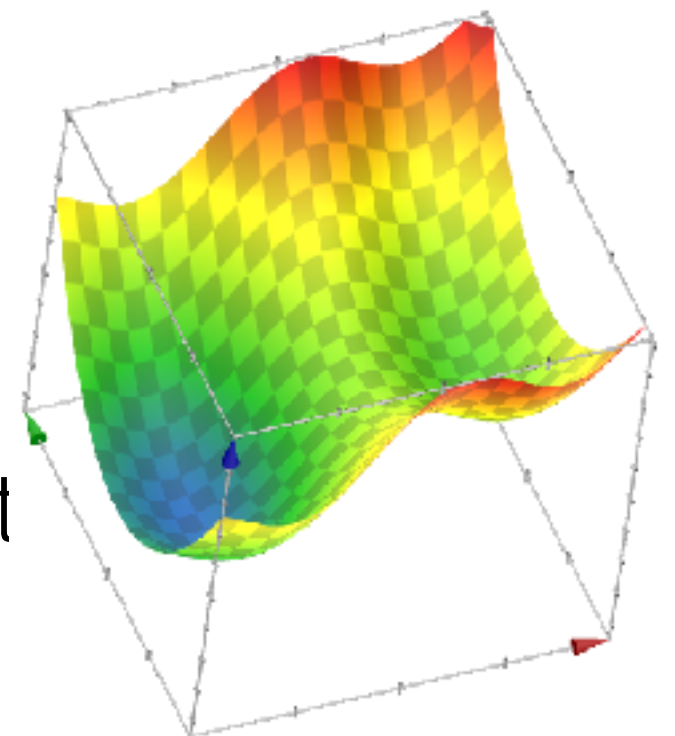  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
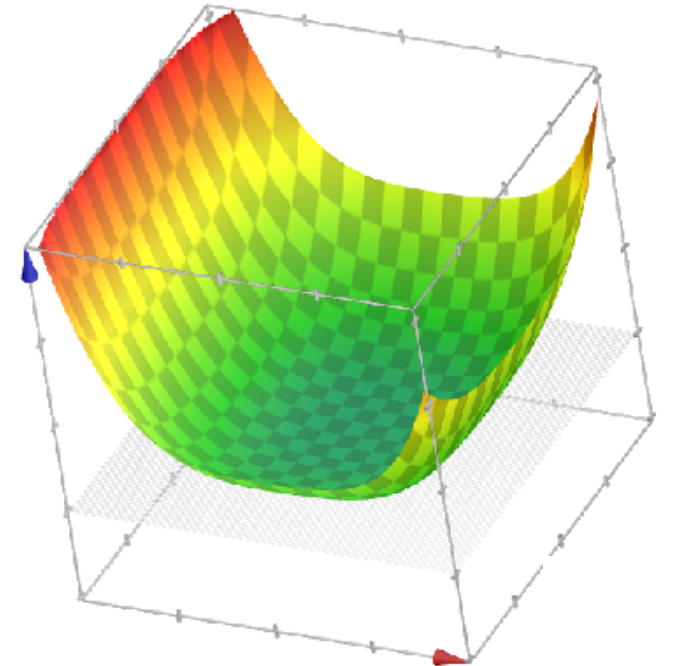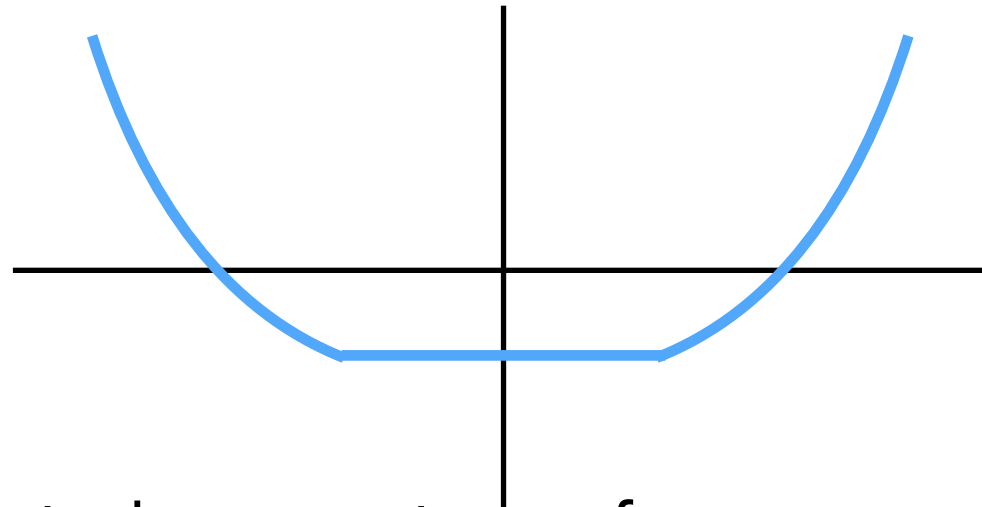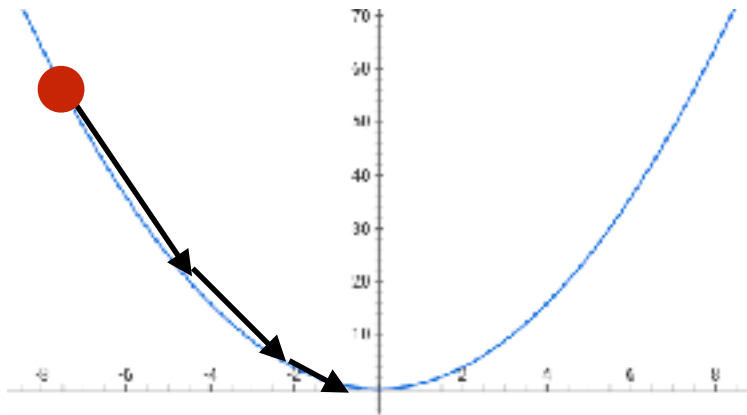
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
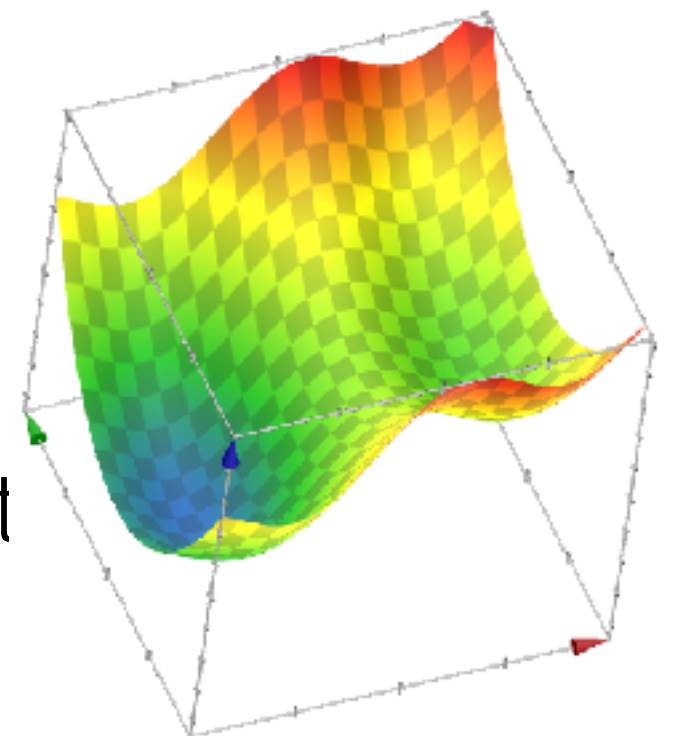  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
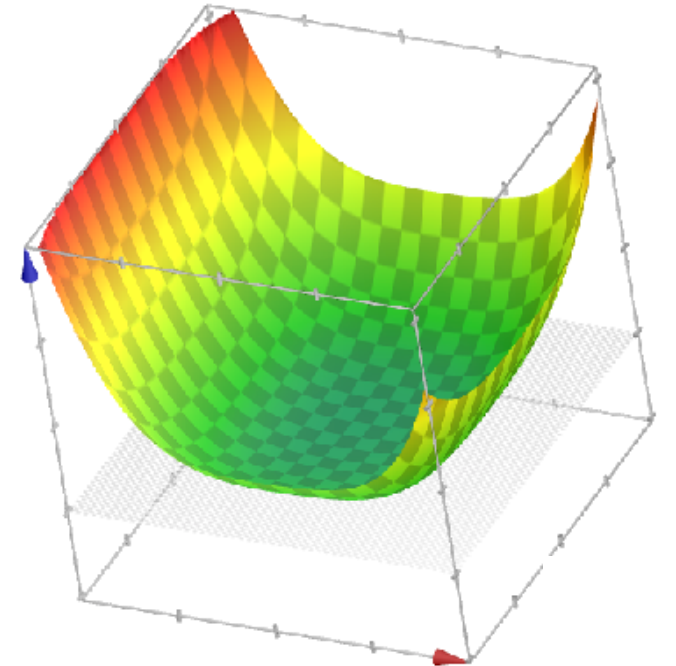
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
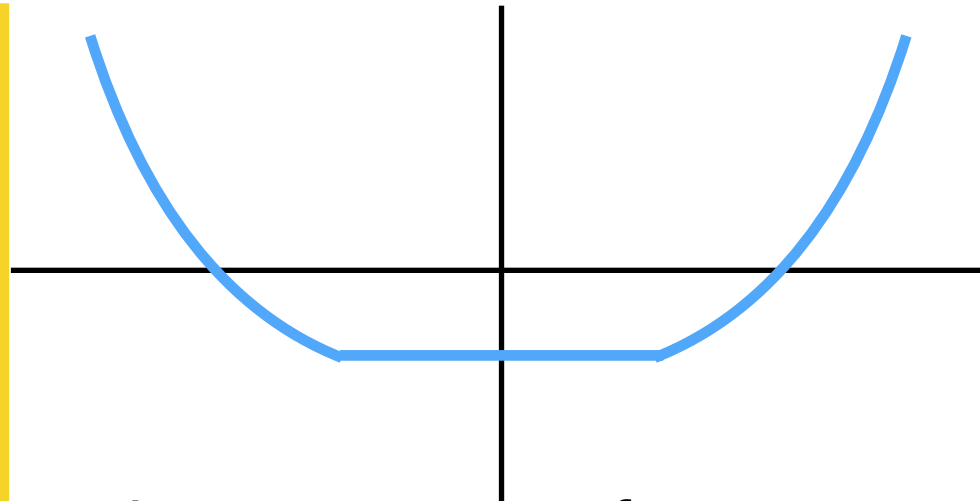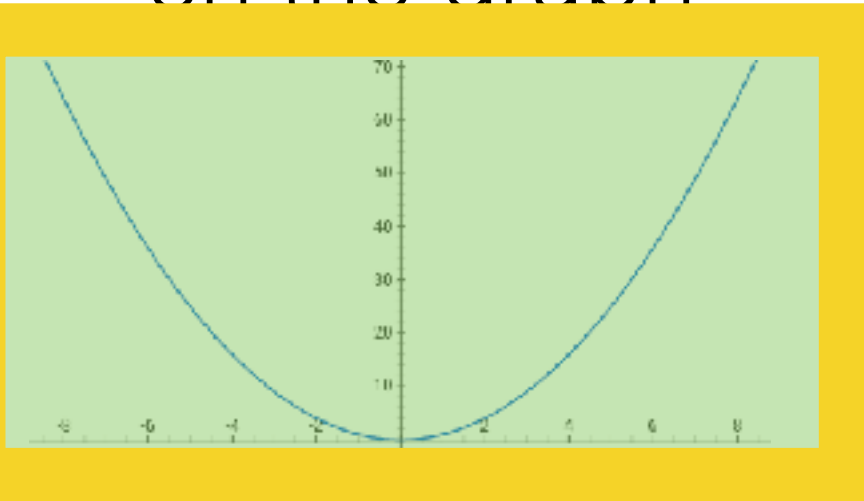


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
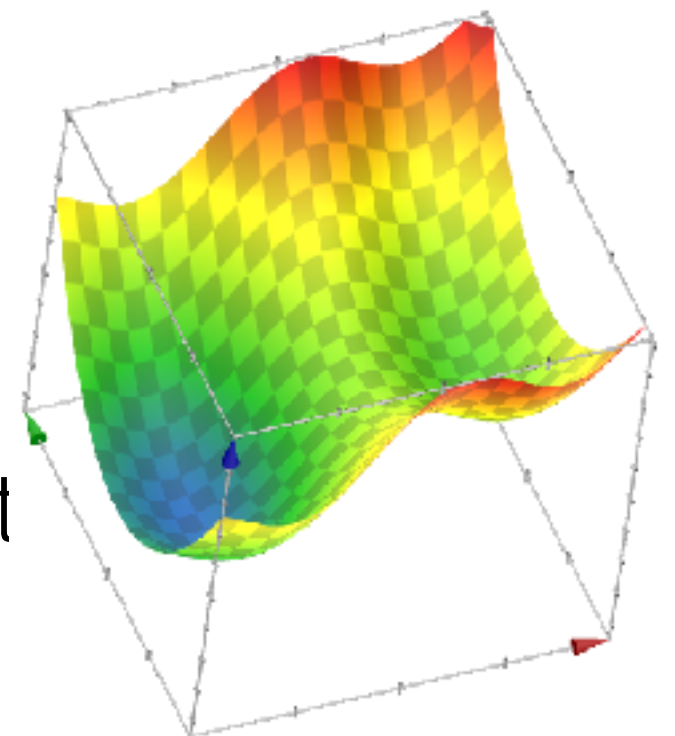  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
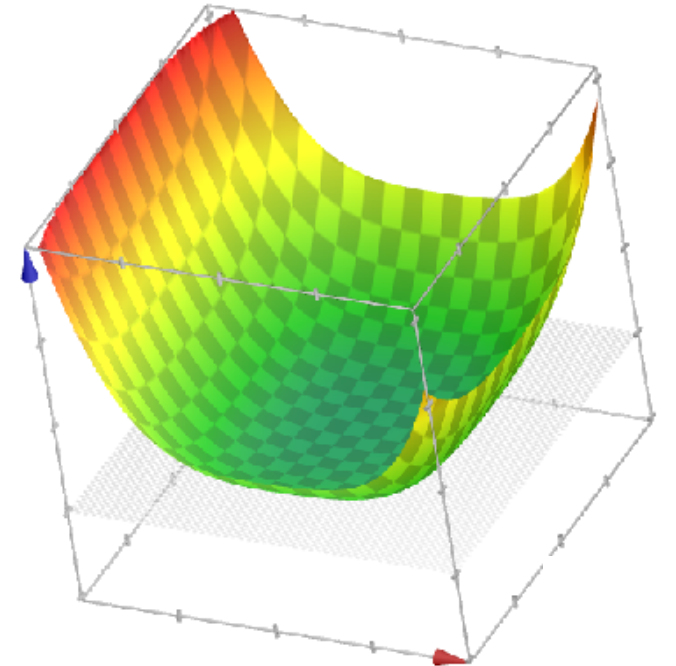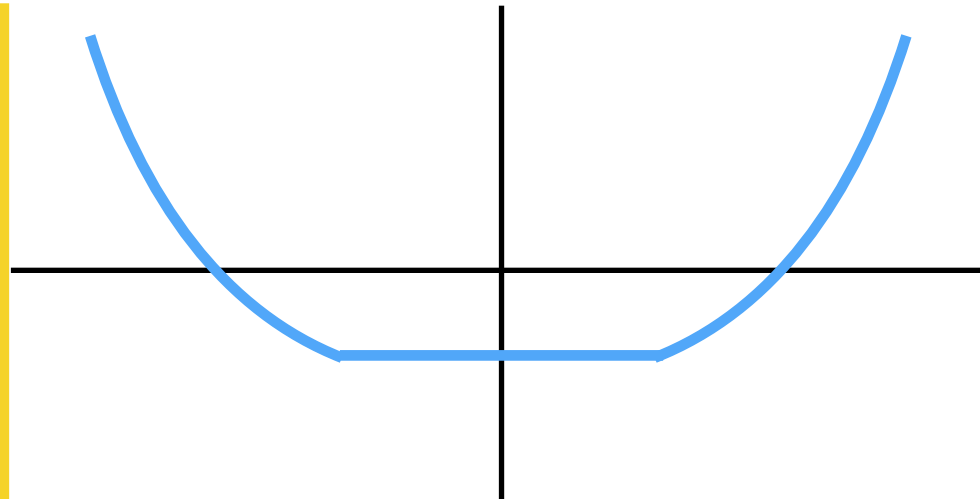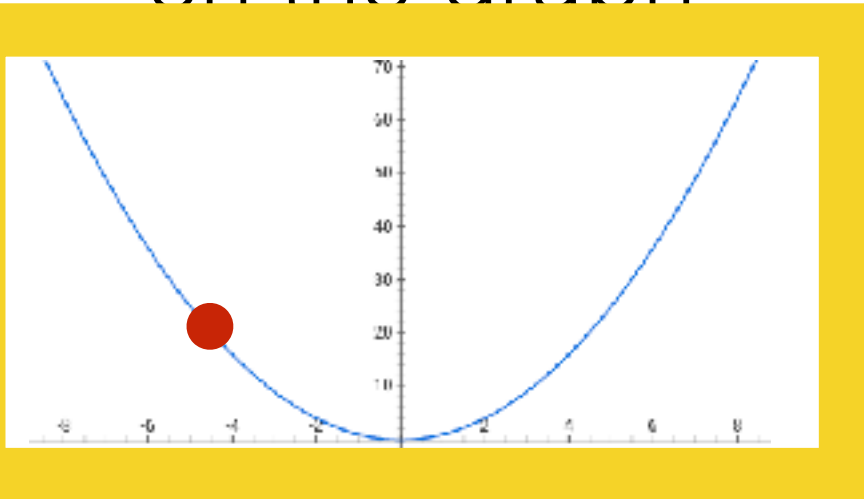


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
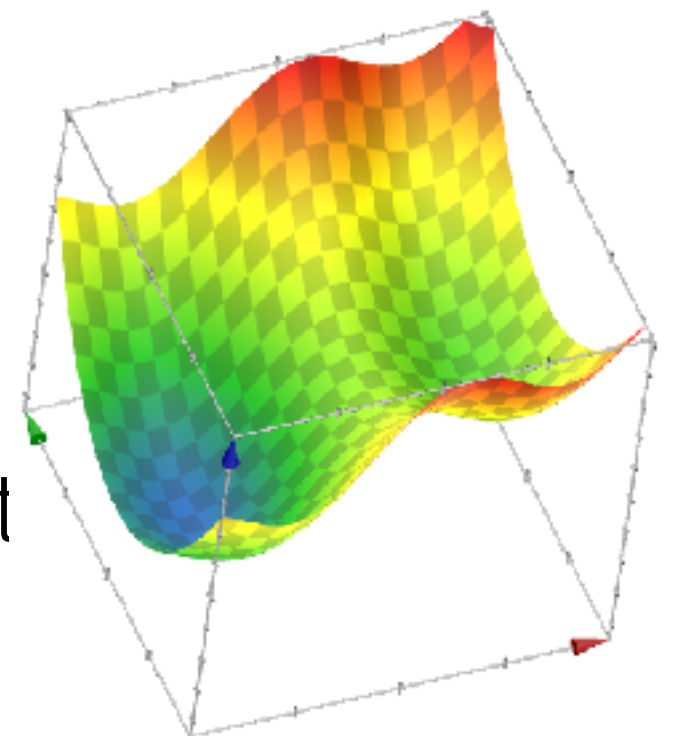  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
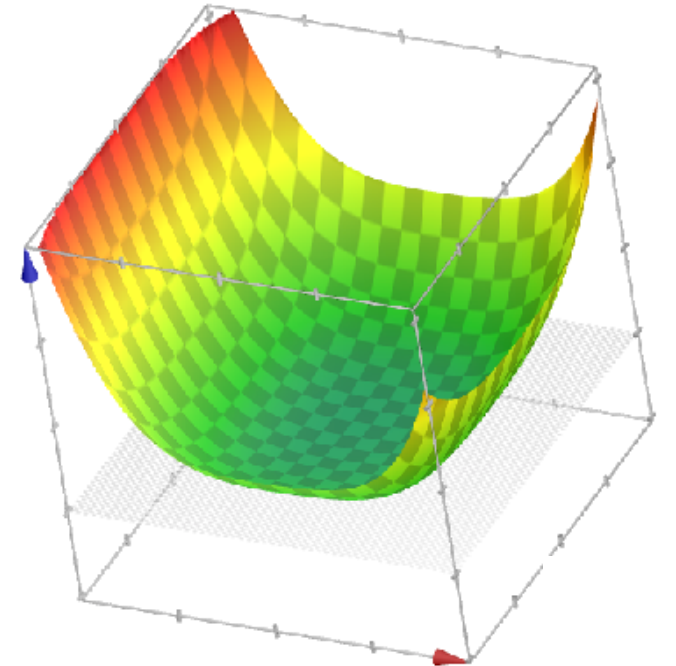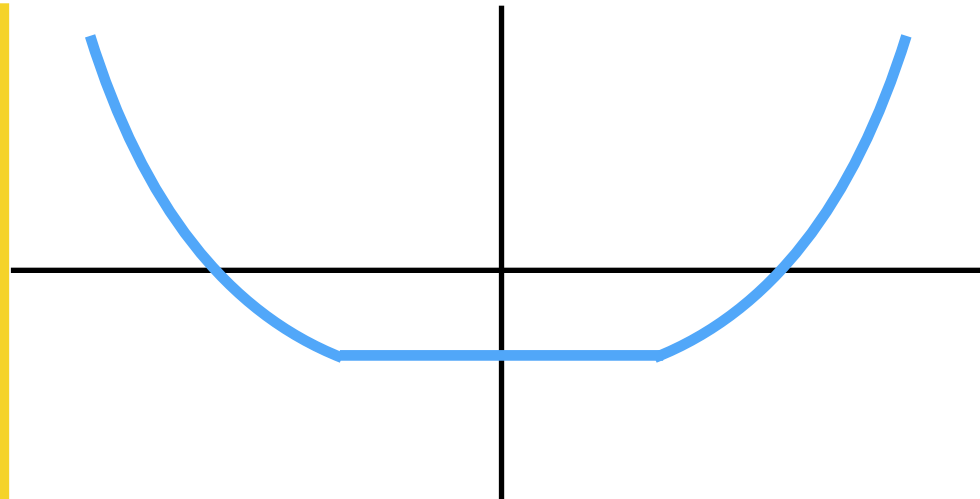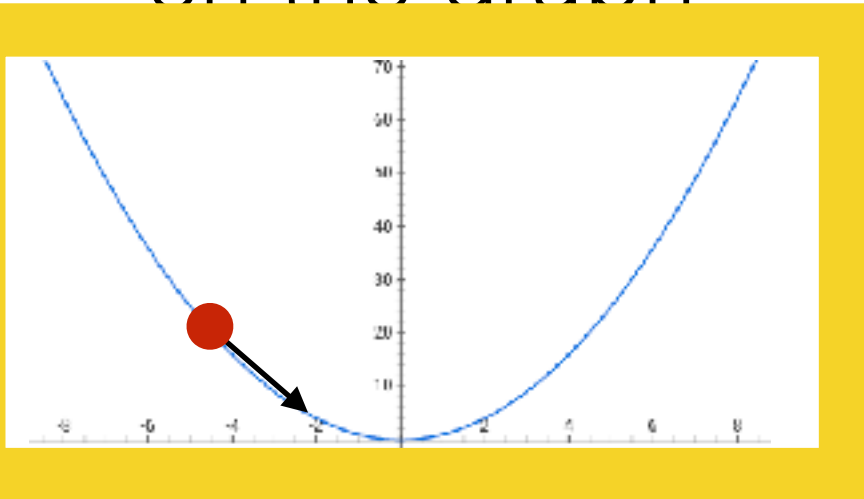


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
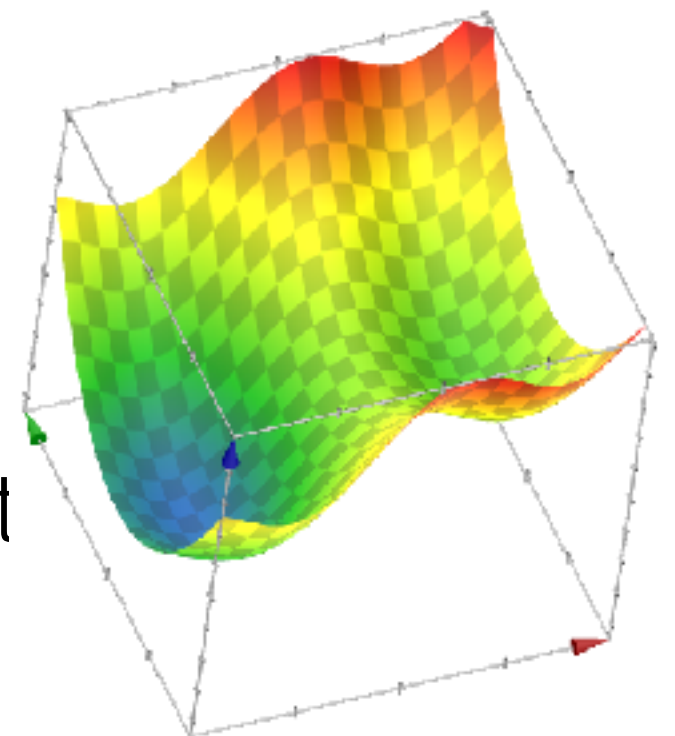  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
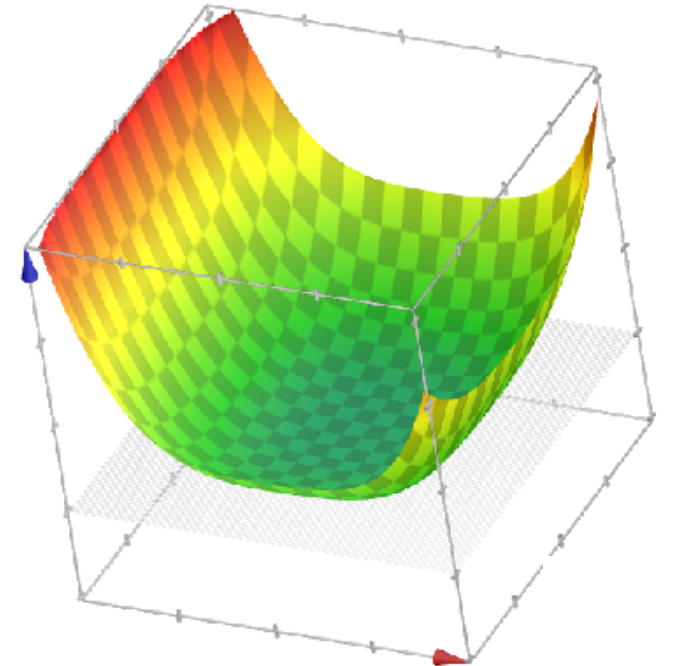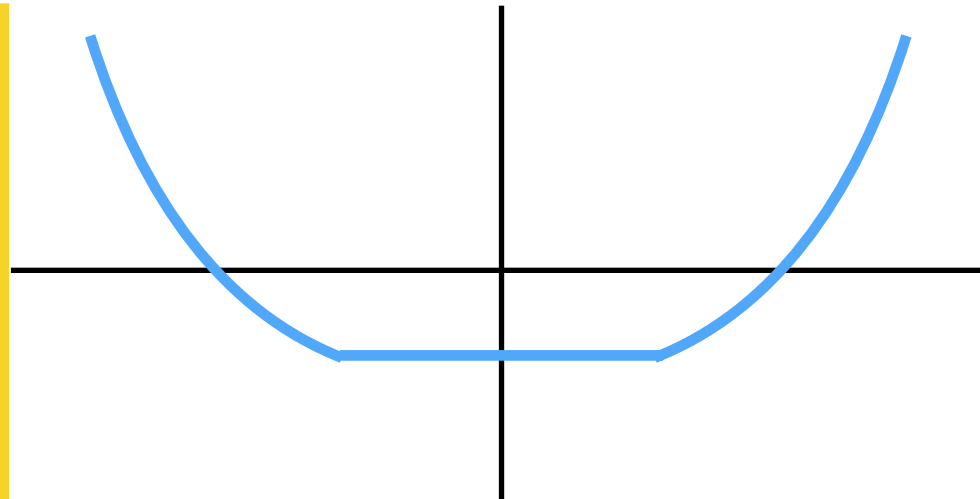
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
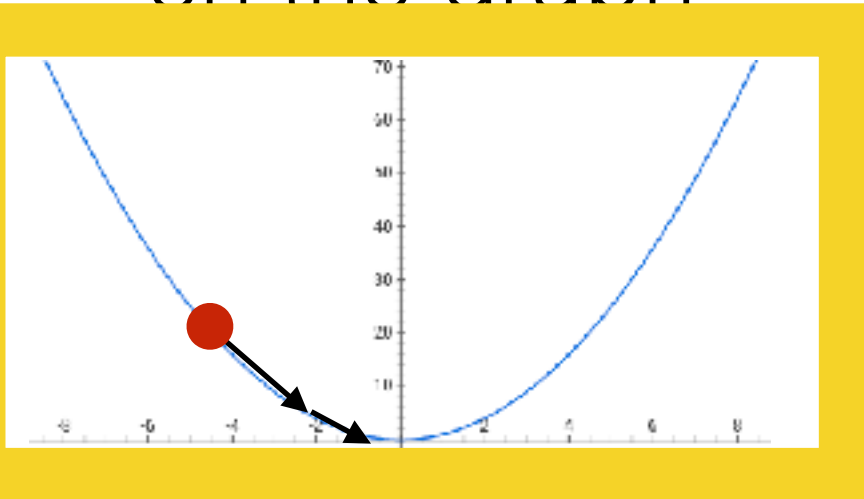


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
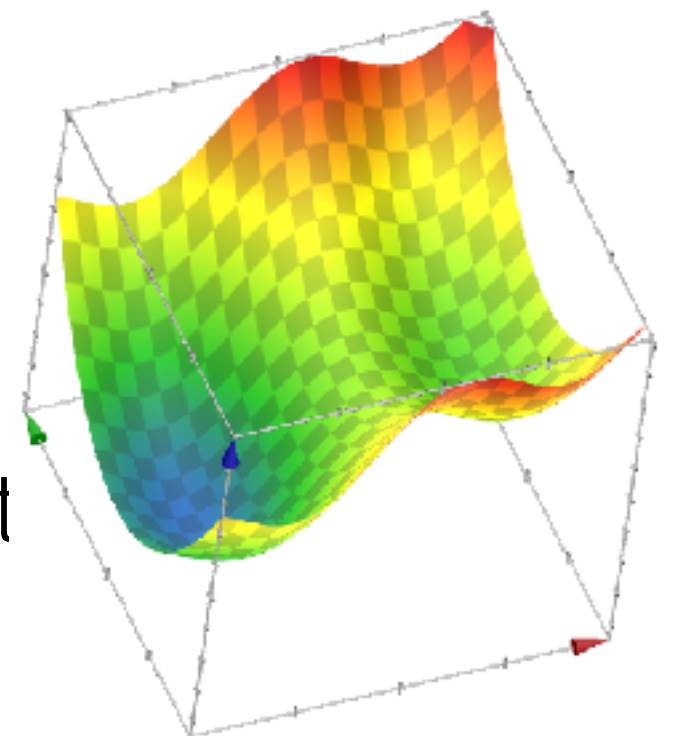  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
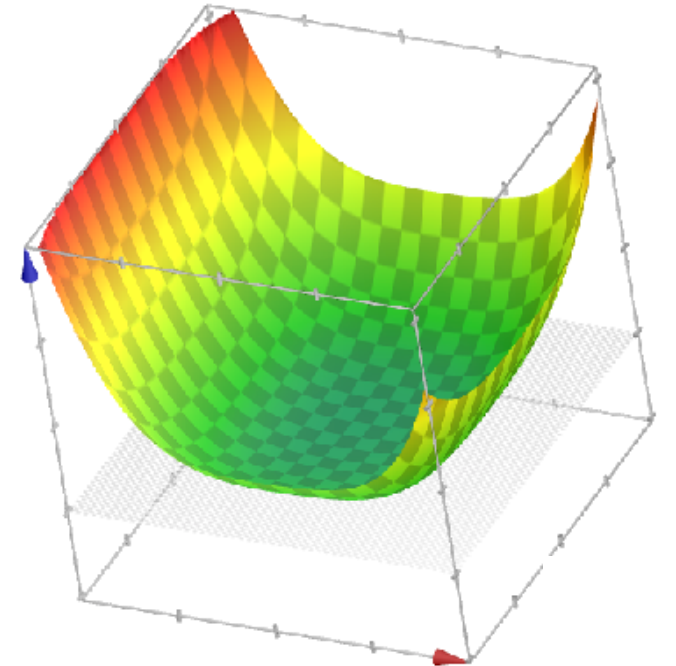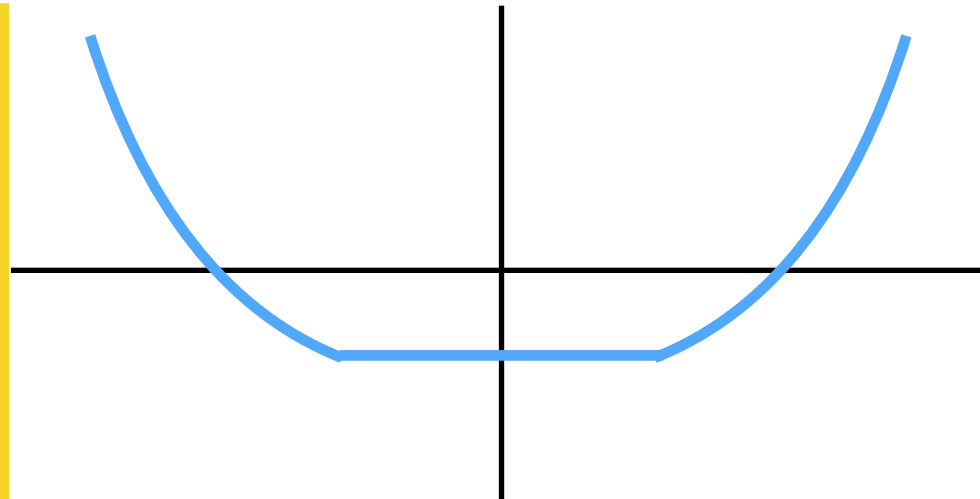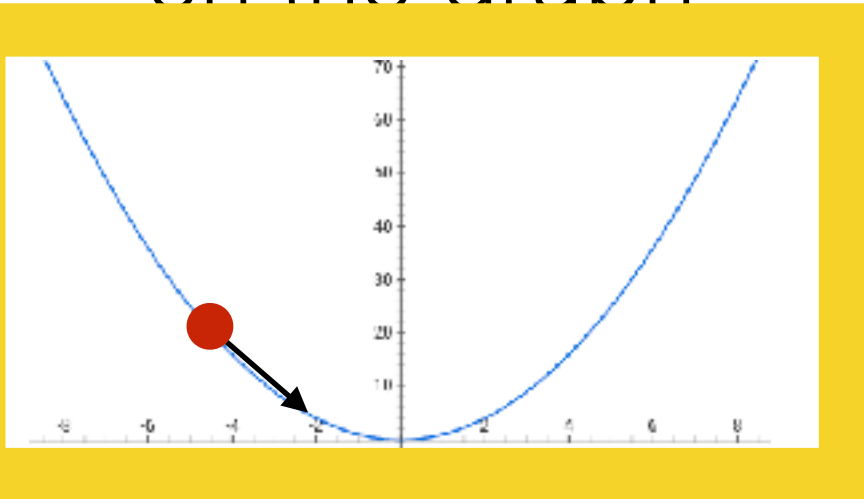
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
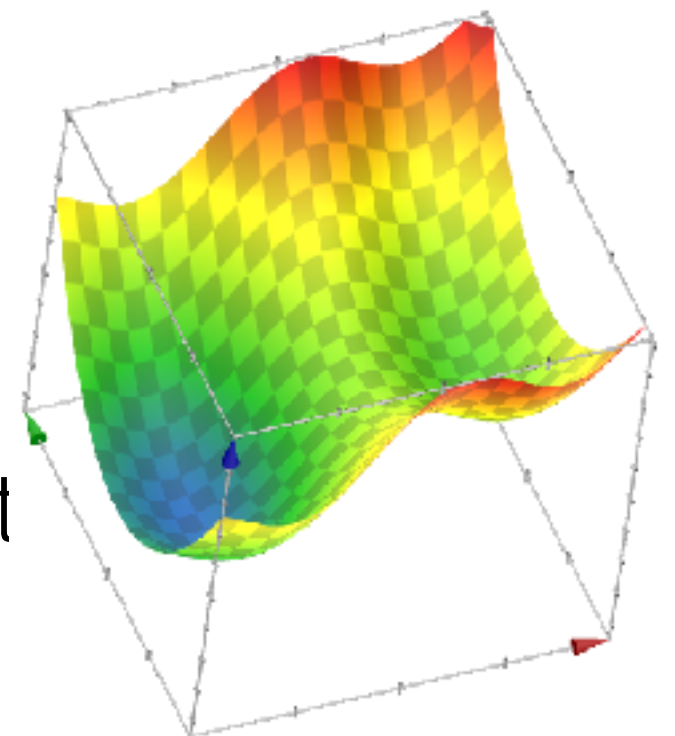  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
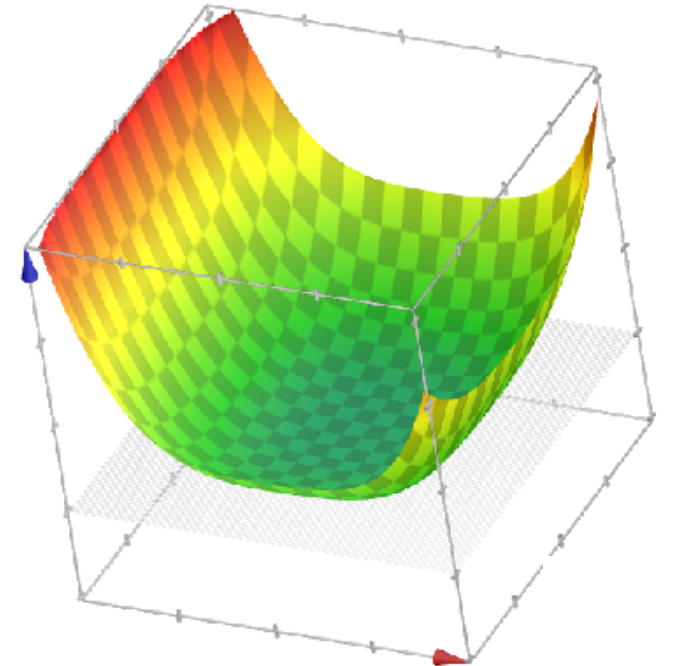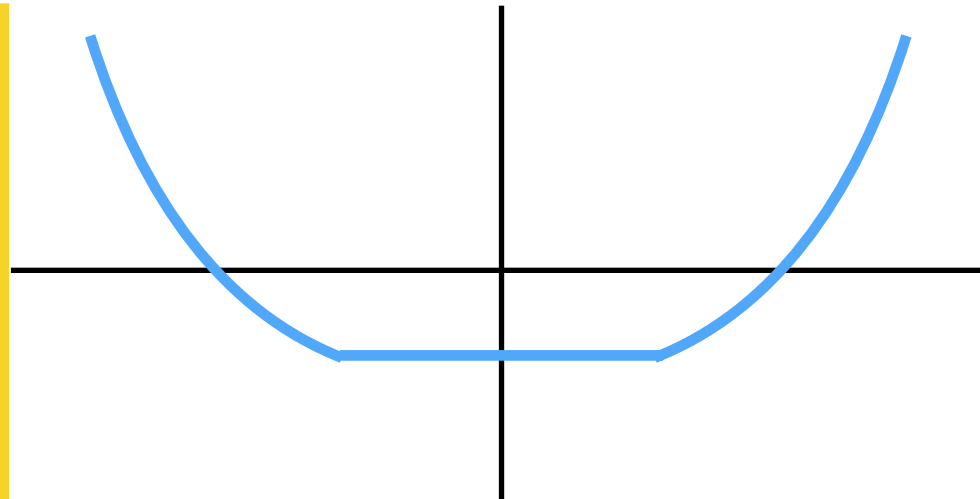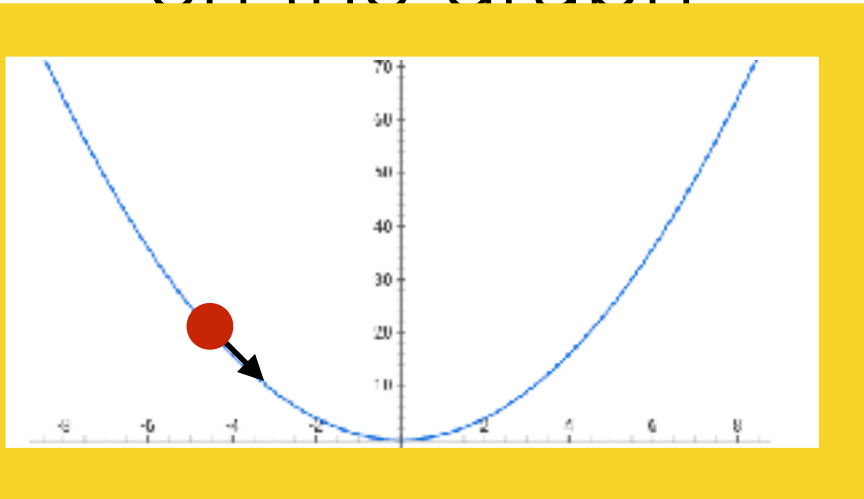


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
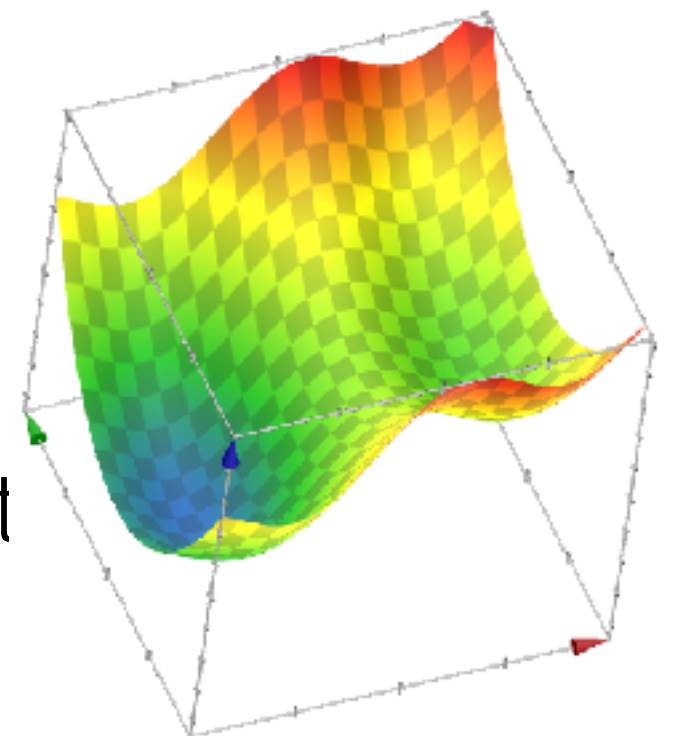  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

6

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph.
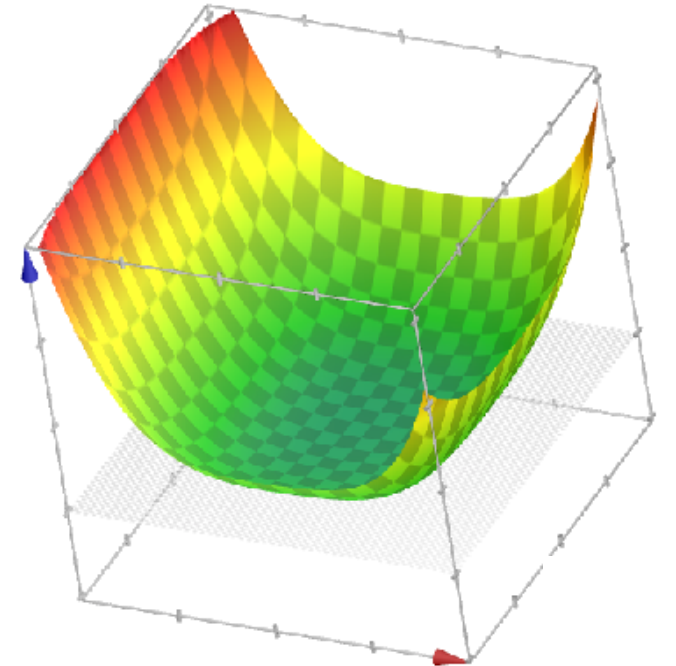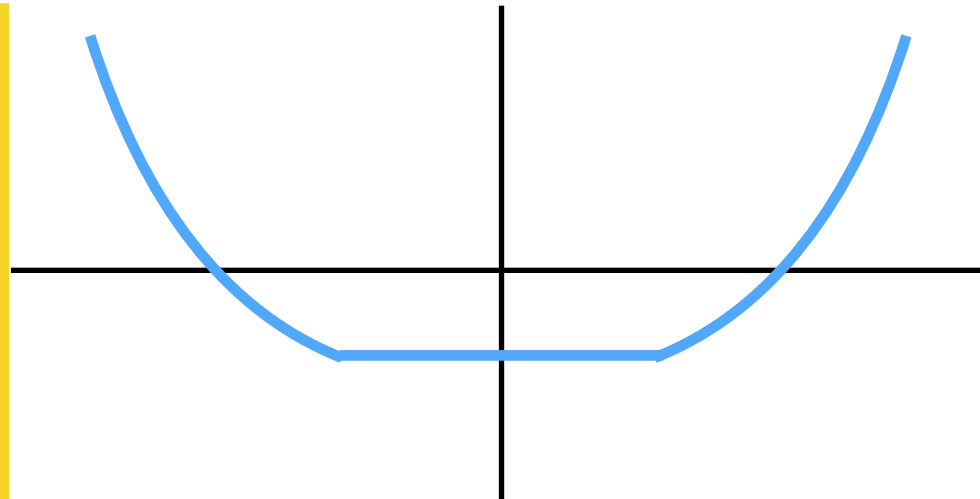
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
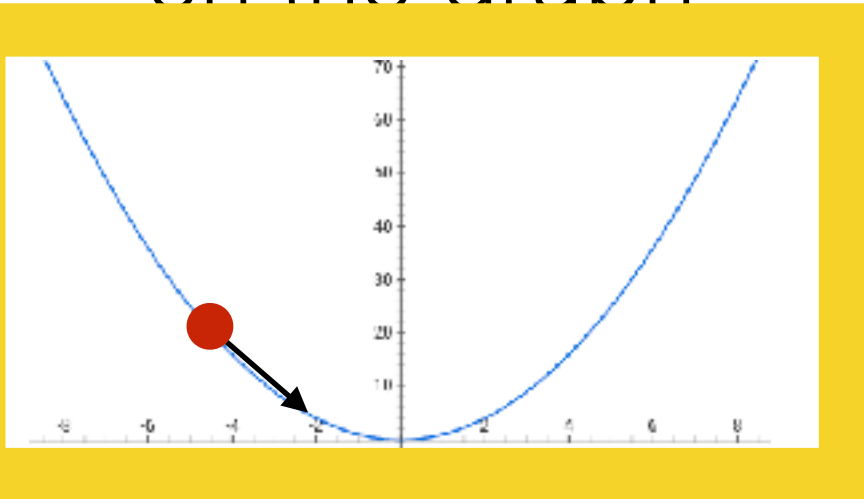
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
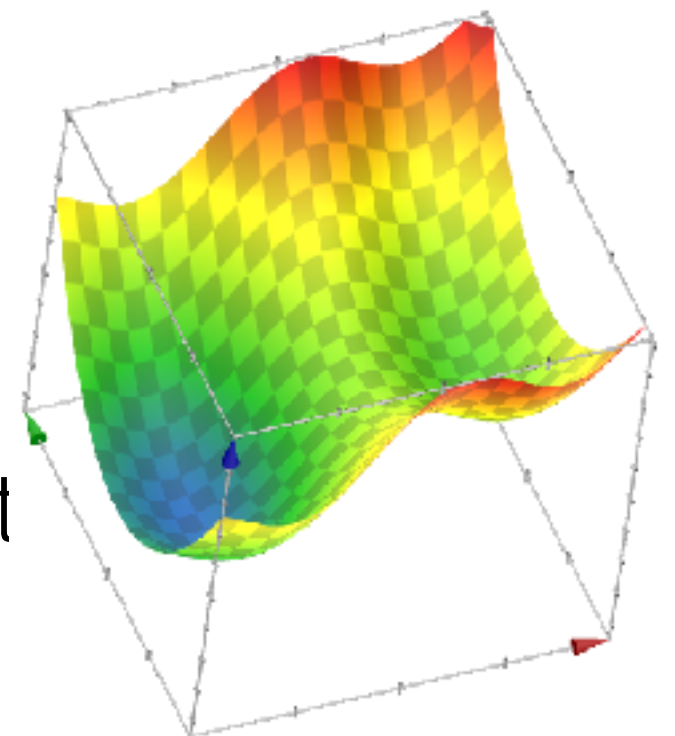  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
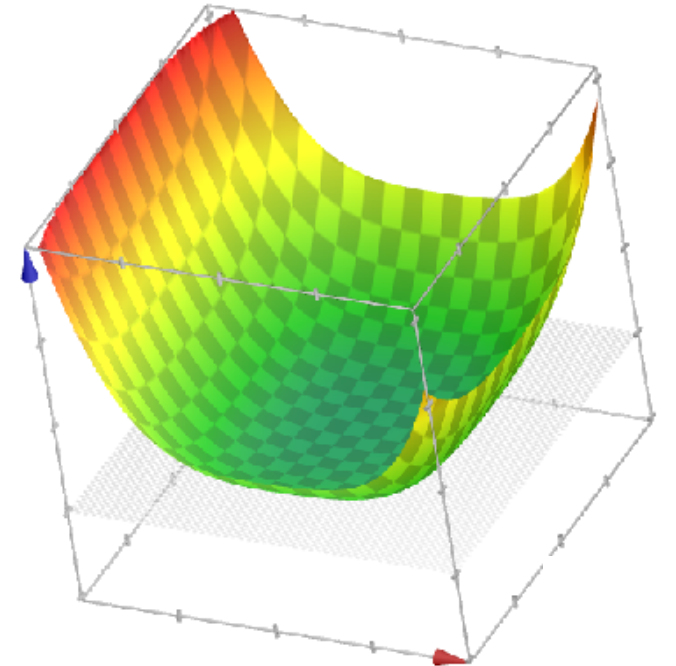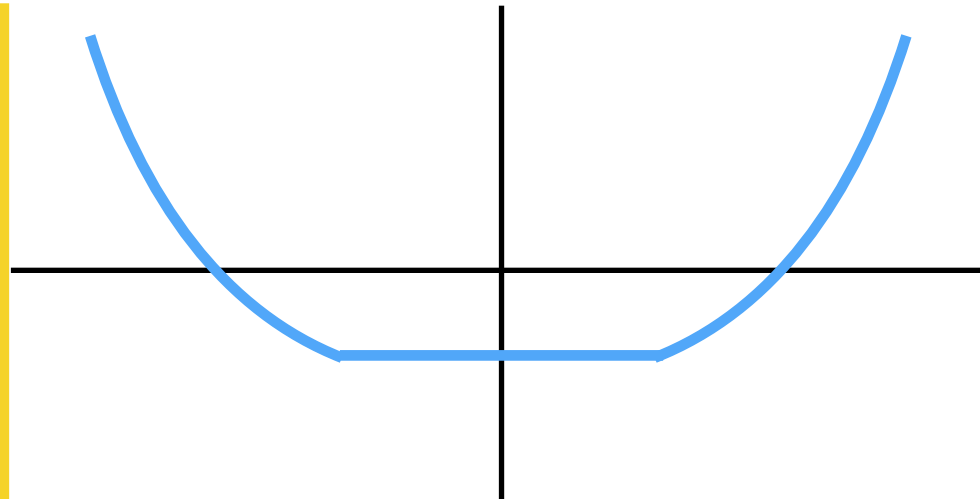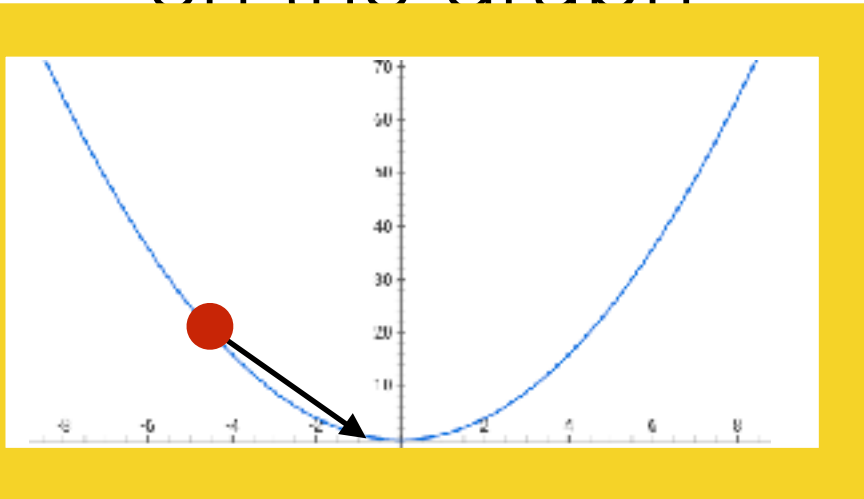


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
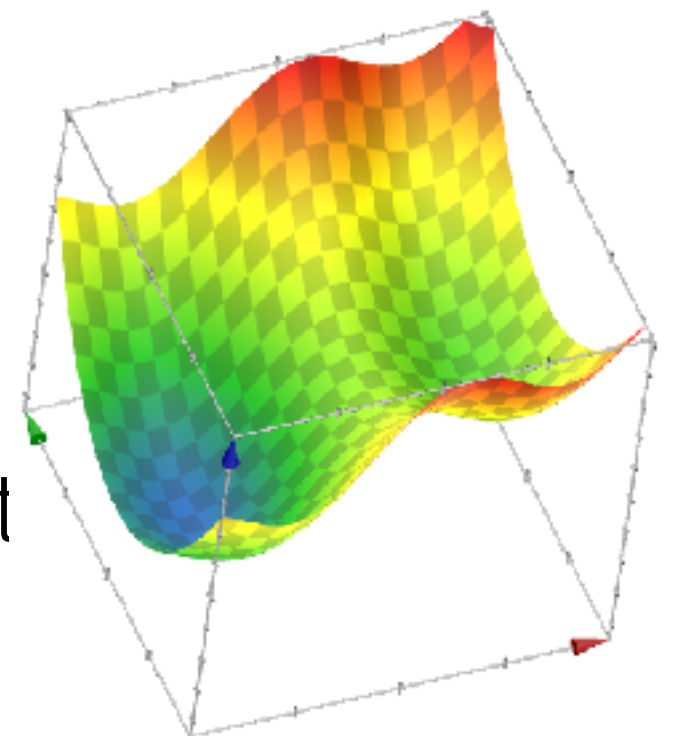  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
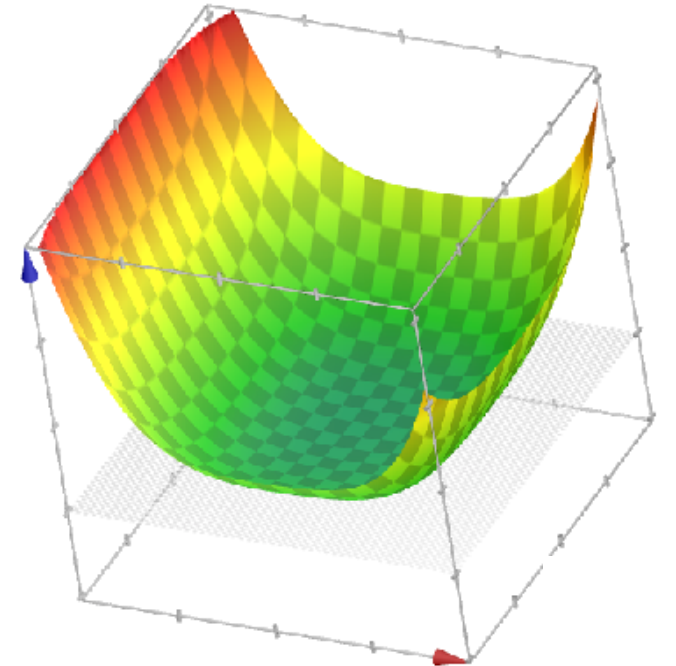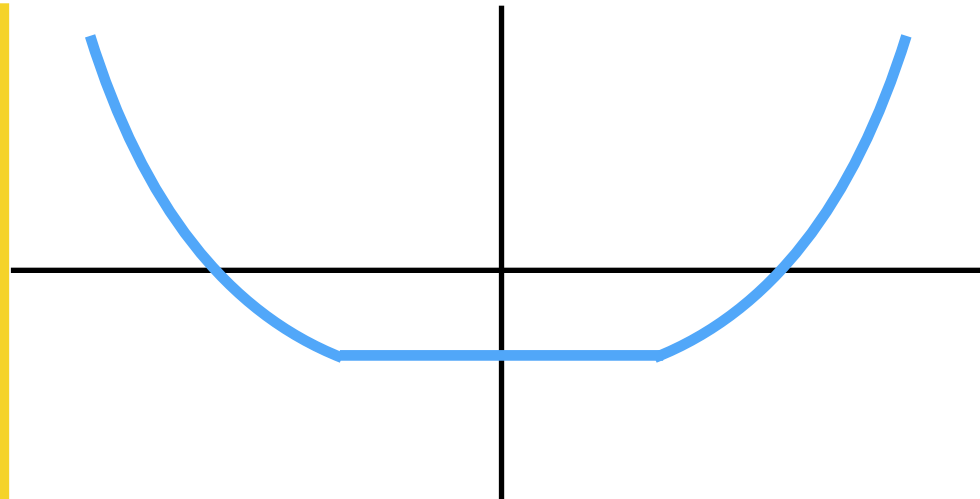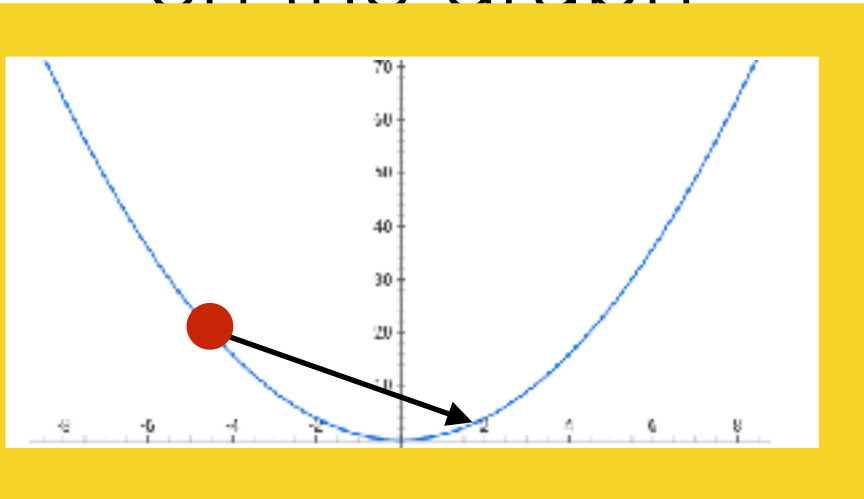


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: <mark>If run long enough</mark>, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$
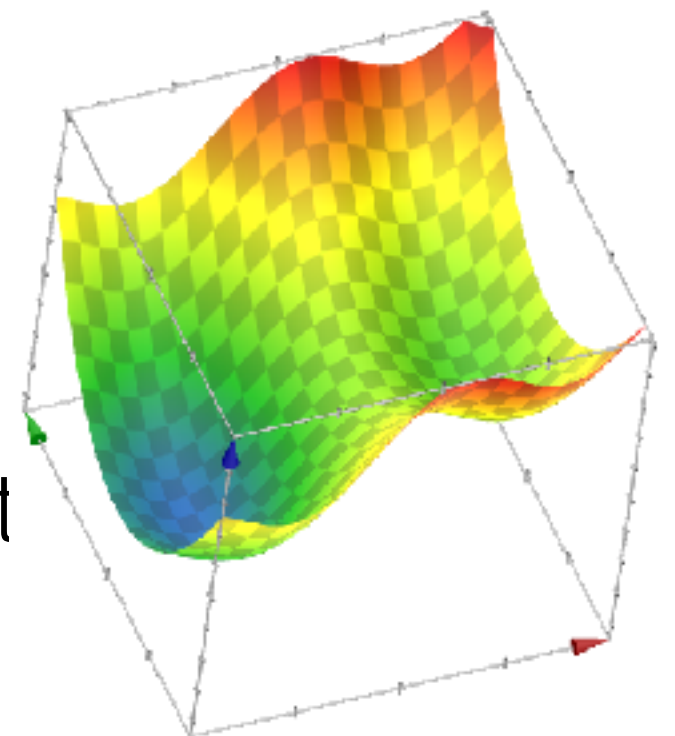
# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
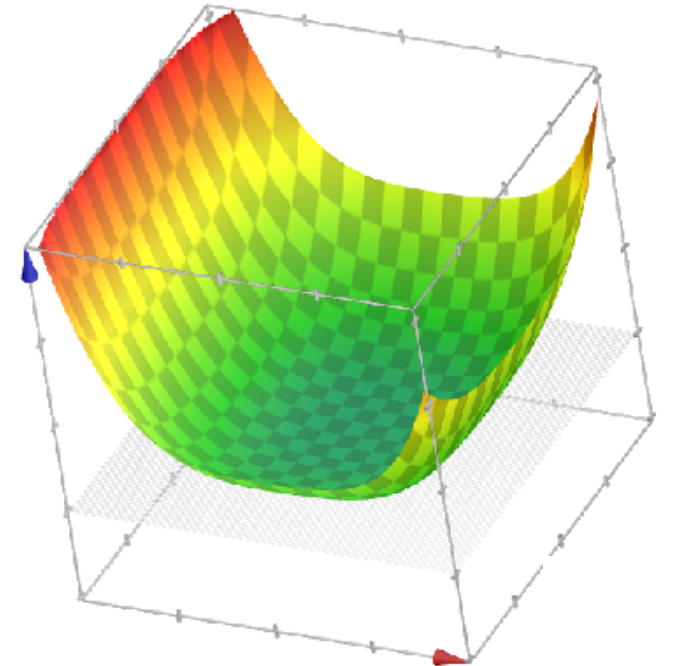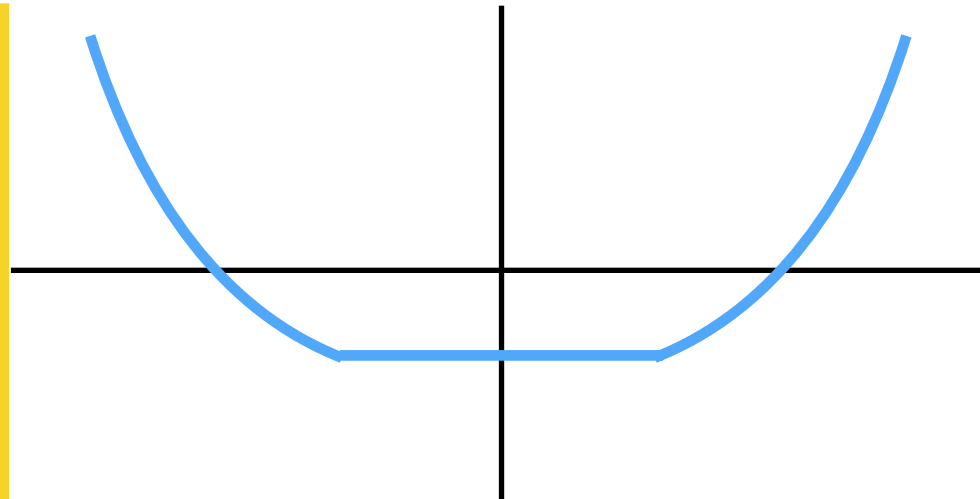


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
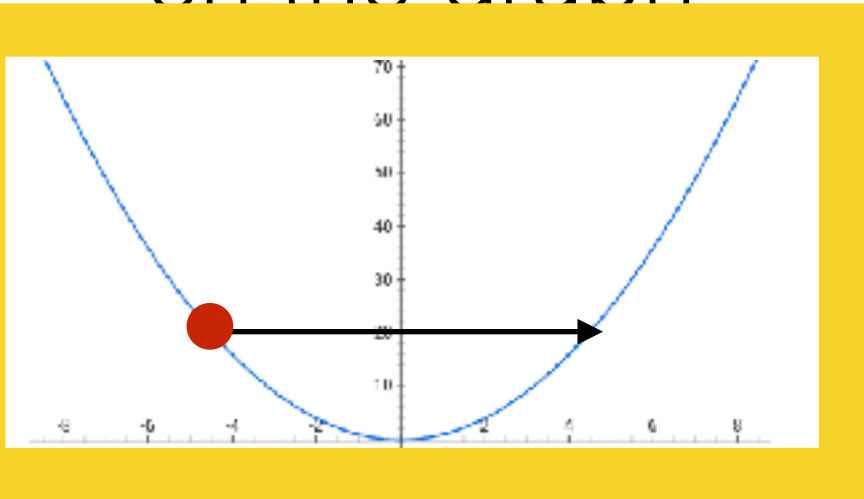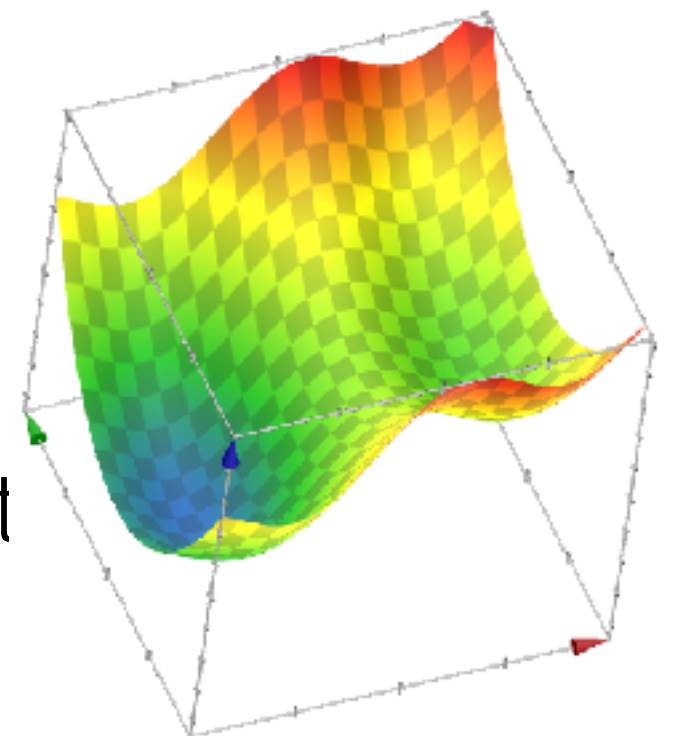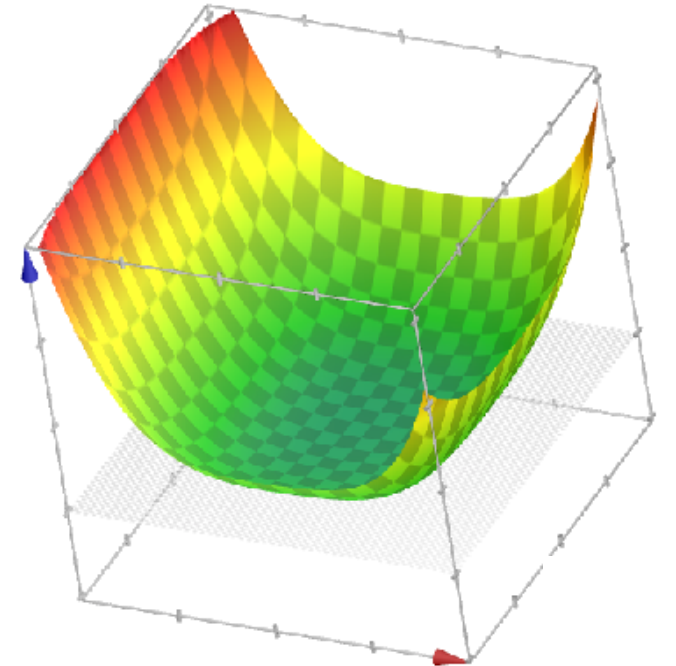


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
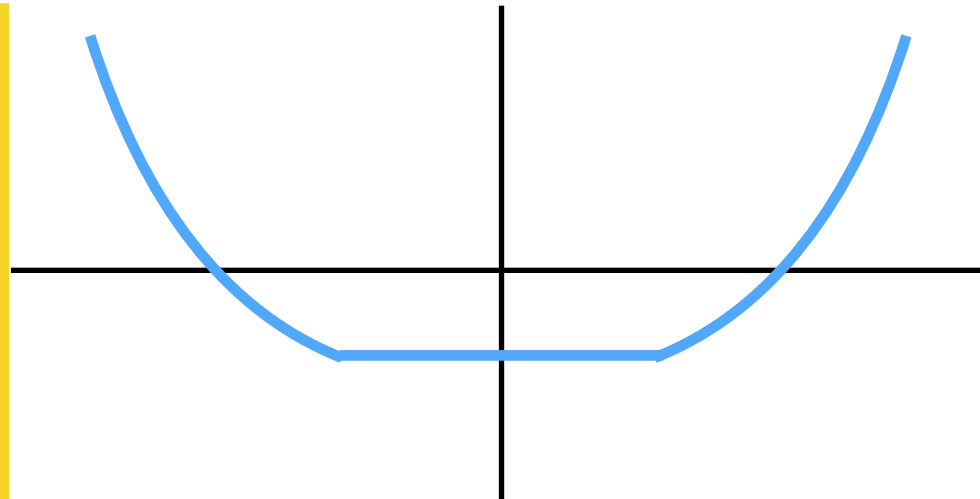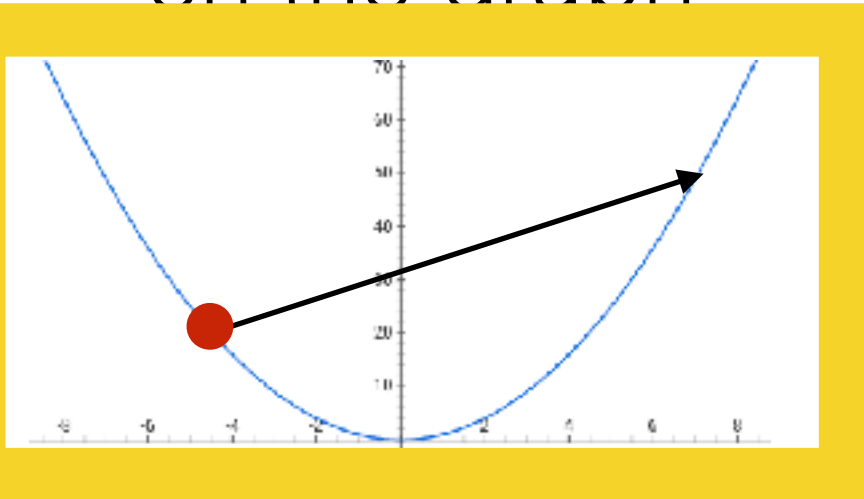


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
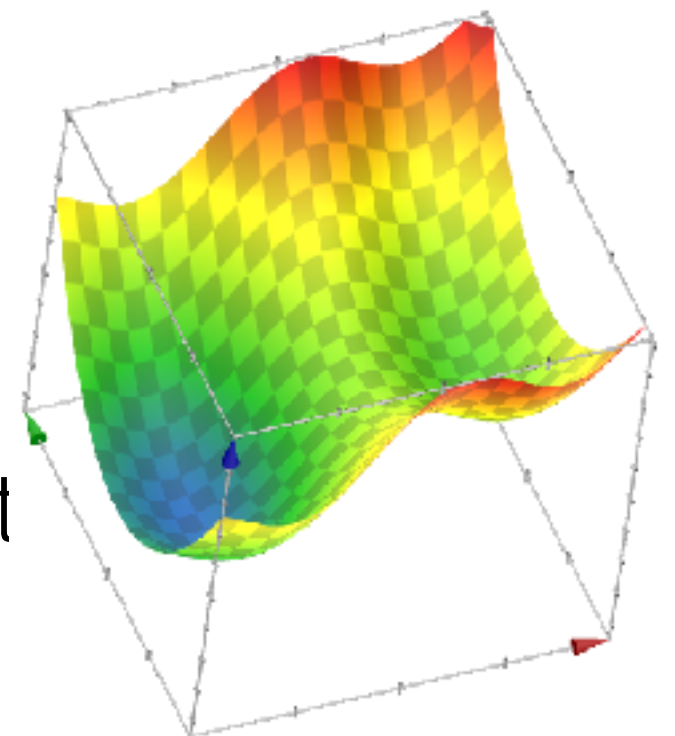  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
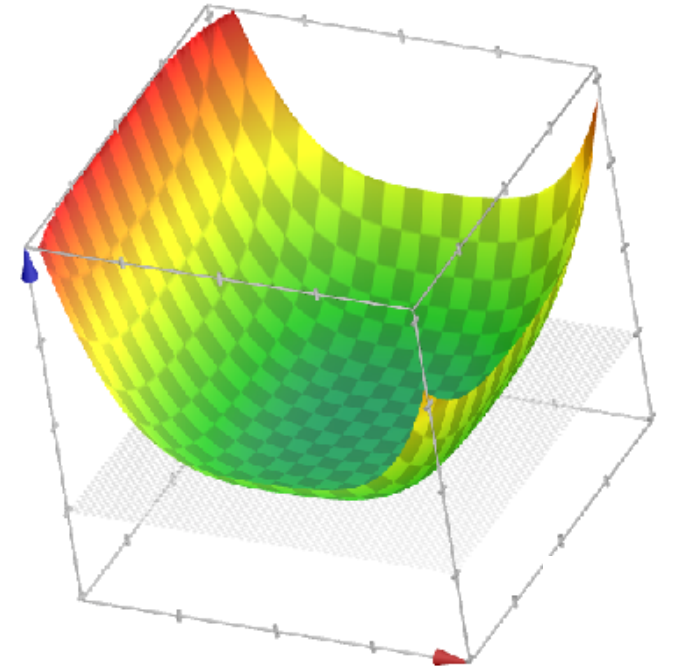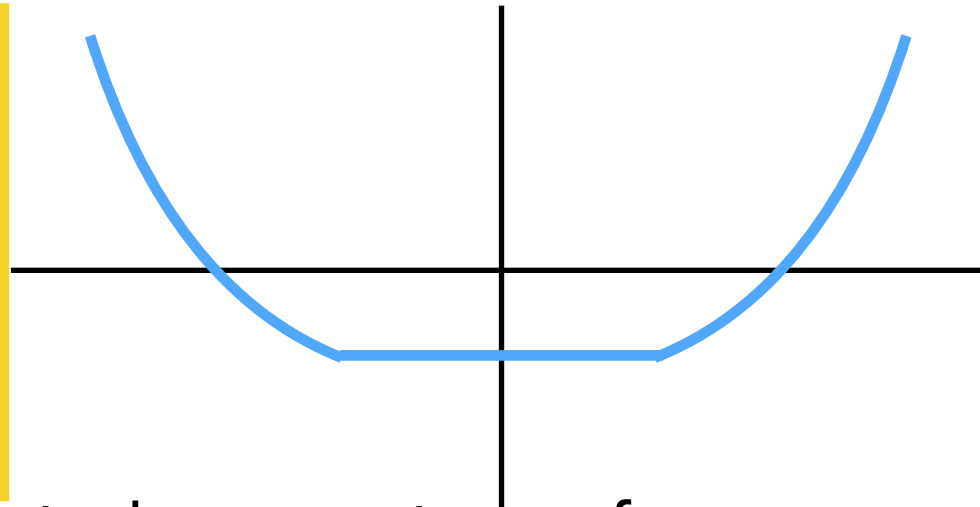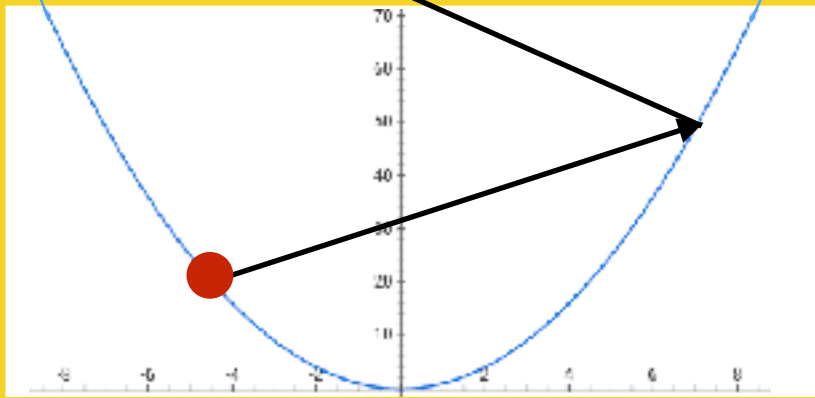
- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
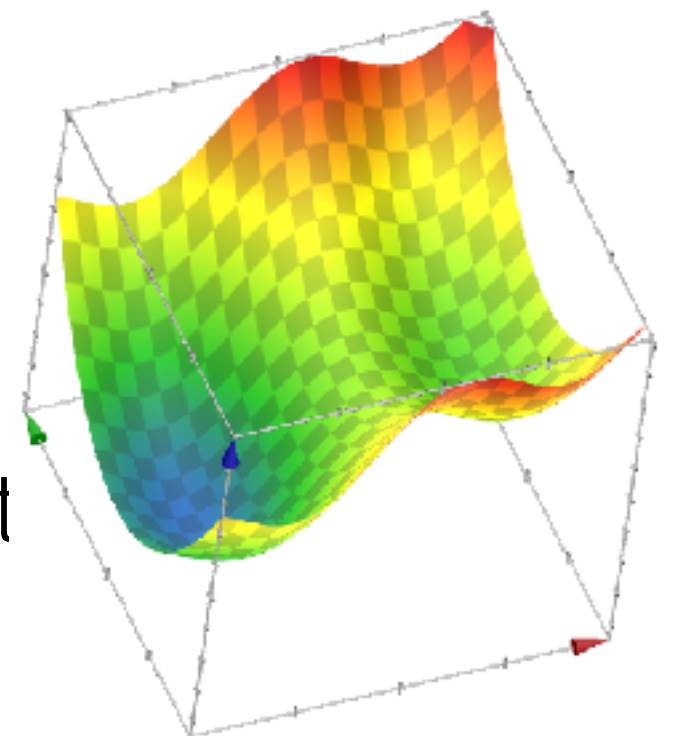  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
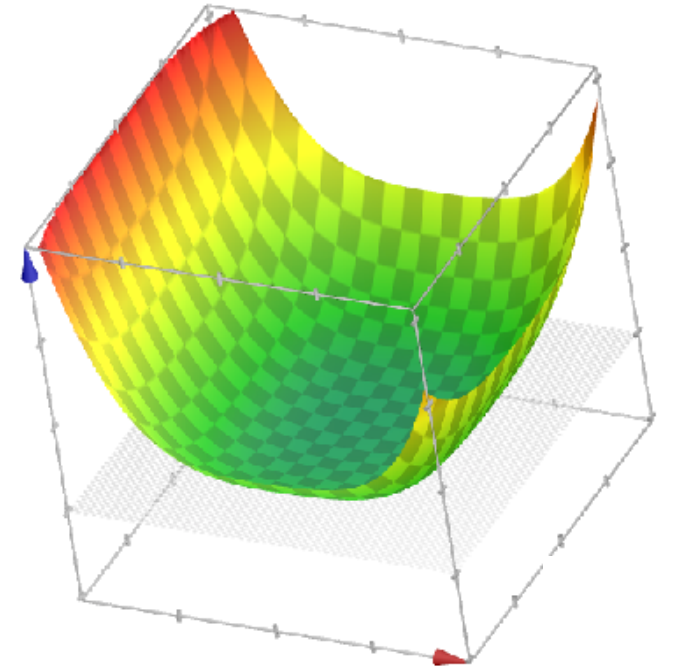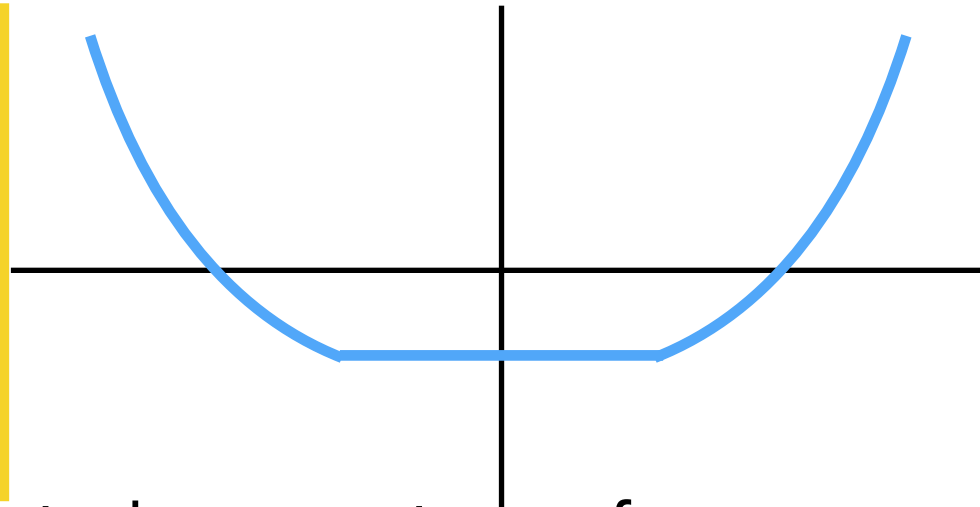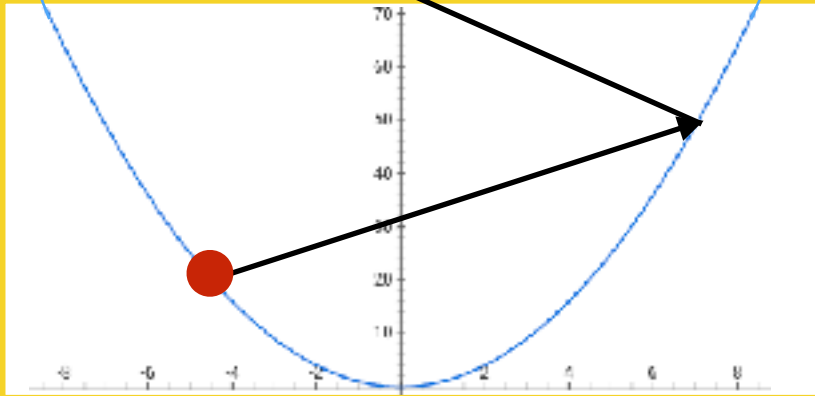


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
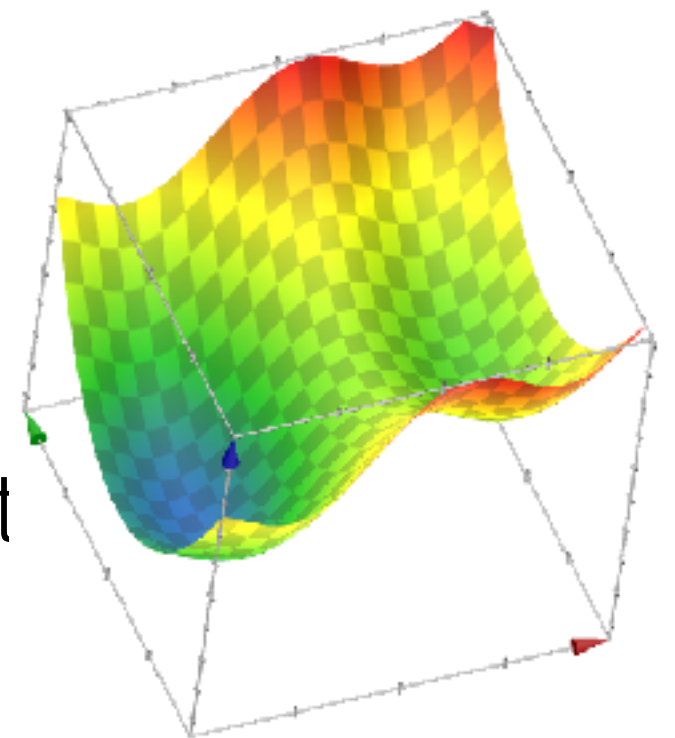  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
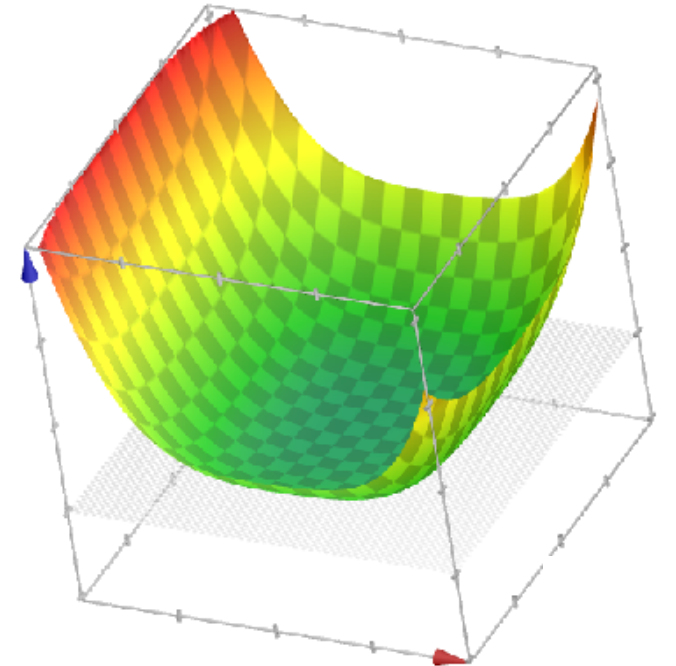






- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
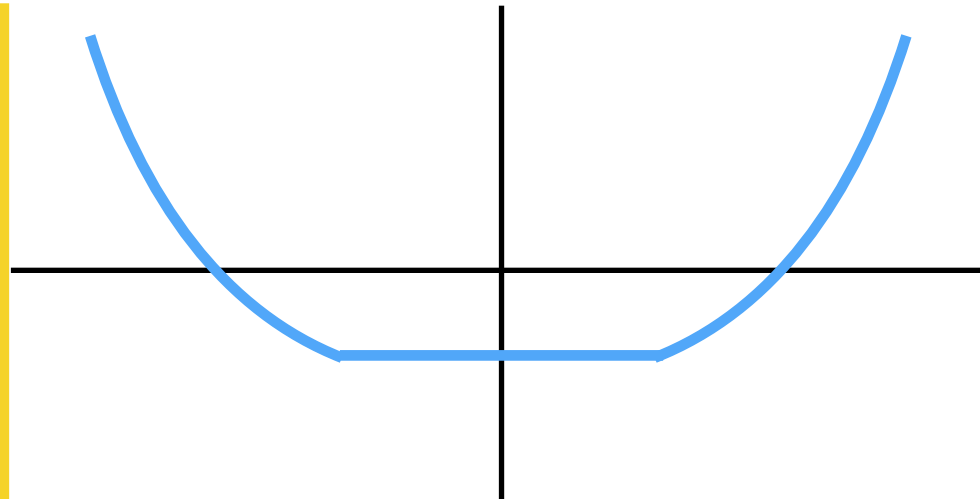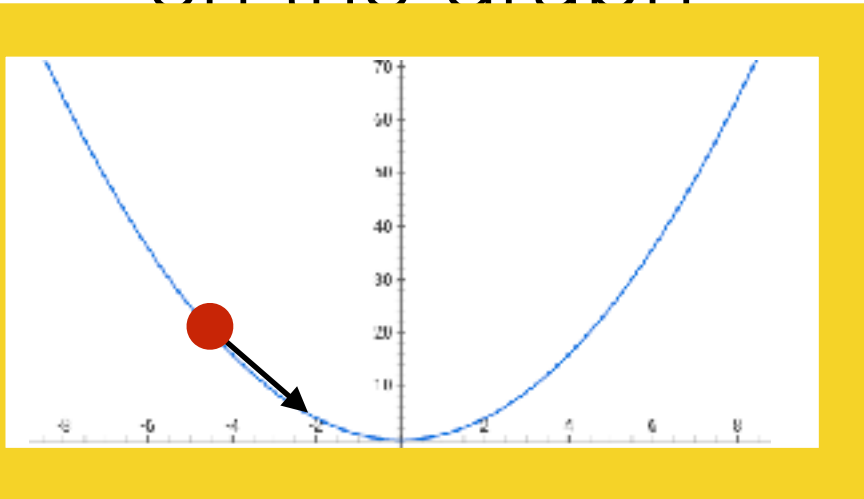




- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
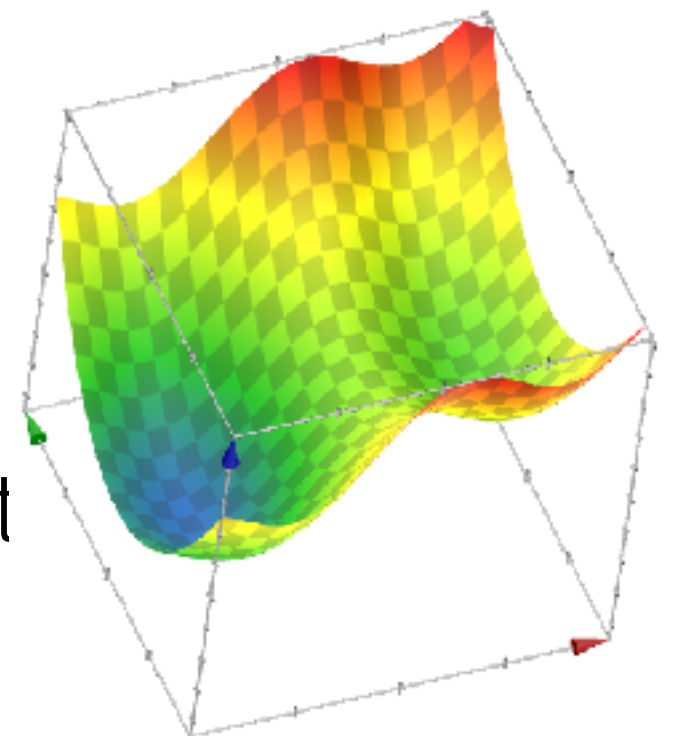  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Gradient descent properties

- A function $f$ on $\mathbb{R}^m$ is convex if any line segment connecting two points of the graph of $f$ lies above or on the graph
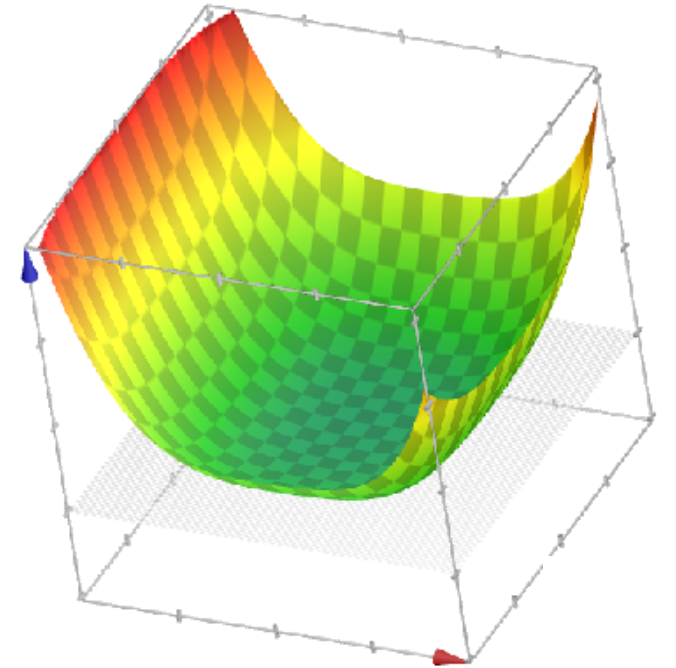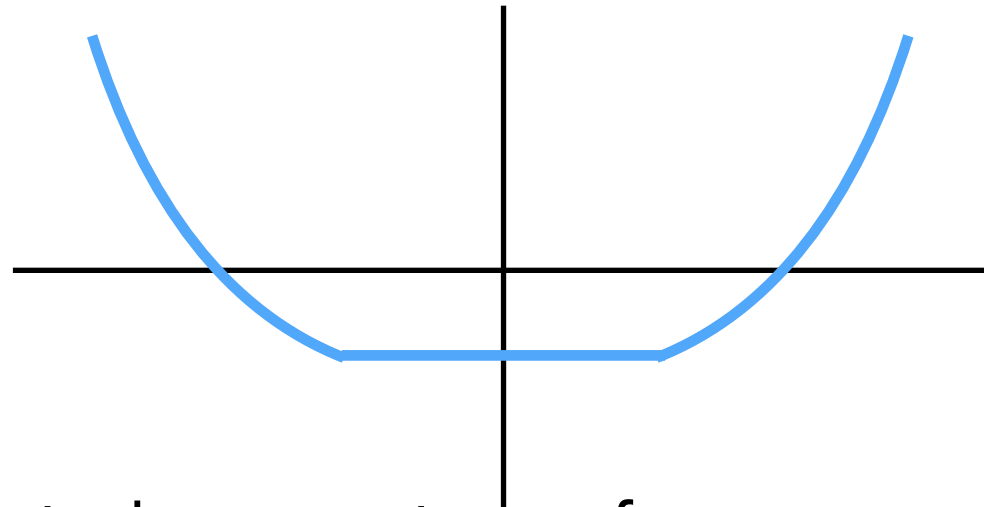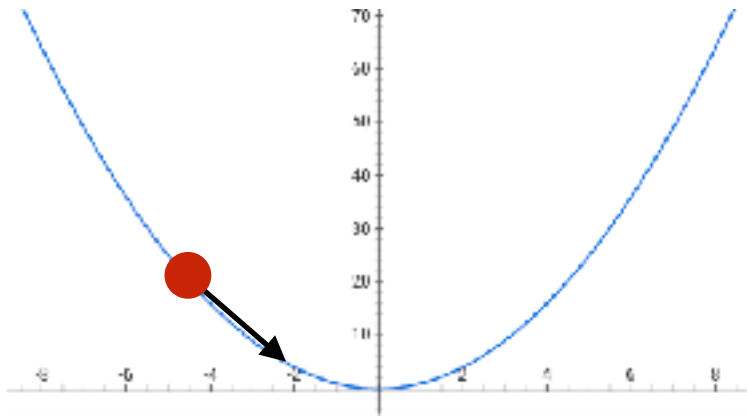


- **Theorem**: Gradient descent performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is sufficiently "smooth" and convex
    - $f$ has at least one global optimum
    - $\eta$ is sufficiently small
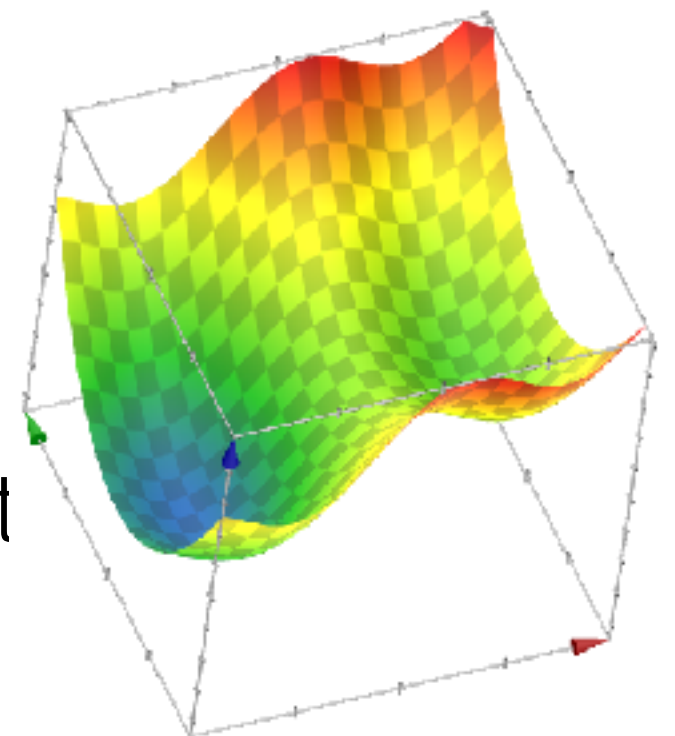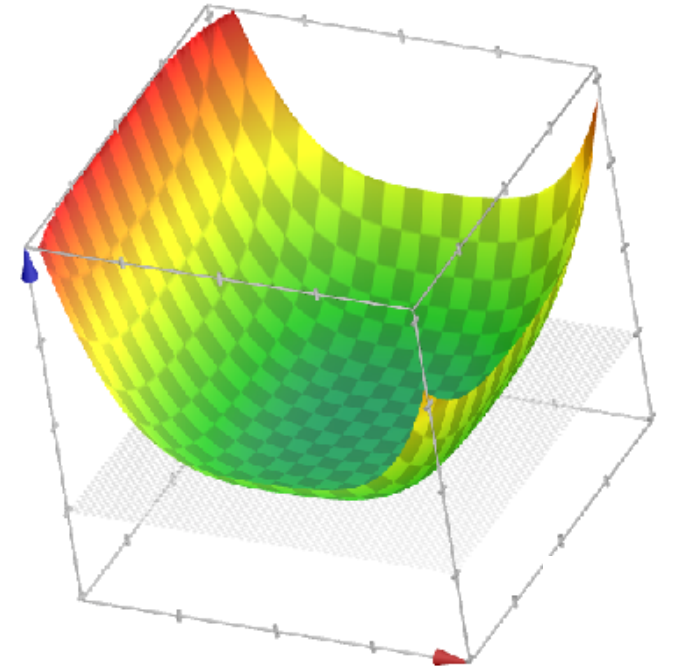  - **Conclusion**: If run long enough, gradient descent will return a value within $\tilde{\epsilon}$ of a global optimum $\Theta$

# Optimizing ridge regression

# Optimizing ridge regression

- Gradient descent

# Optimizing ridge regression

- Gradient descent

# Optimizing ridge regression

- Gradient descent

# Optimizing ridge regression

- Gradient descent

# Optimizing ridge regression

- Gradient descent

# Optimizing ridge regression

- Gradient descent

# Optimizing ridge regression

- Gradient descent vs. analytical/closed-form/direct solution

# Optimizing ridge regression

- Gradient descent vs. analytical/closed-form/direct solution

# Optimizing ridge regression

- Gradient descent vs. analytical/closed-form/direct solution

# Optimizing ridge regression

- Gradient descent vs. analytical/closed-form/direct solution



- Accuracy doesn't mean anything without running time

# Optimizing ridge regression

- Gradient descent vs. analytical/closed-form/direct solution



- Accuracy doesn't mean anything without running time
- Running time doesn't mean anything without accuracy

# Optimizing ridge regression

- Gradient descent vs. analytical/closed-form/direct solution



- Accuracy doesn't mean anything without running time

- Running time doesn't mean anything without accuracy

- Need to measure accuracy for the running time we have

# Optimizing ridge regression

- Gradient descent vs. analytical/closed-form/direct solution



- Accuracy doesn't mean anything without running time

- Running time doesn't mean anything without accuracy

- Need to measure accuracy for the running time we have
  - Recall: closed-
    form solution (if
    no offset)

# Optimizing ridge regression

- Gradient descent vs. analytical/closed-form/direct solution



- Accuracy doesn't mean anything without running time

- Running time doesn't mean anything without accuracy

- Need to measure accuracy for the running time we have
  - Recall: closed-form solution (if no offset)

$$\theta = (\tilde{X}^\top \tilde{X} + n\lambda I)^{-1} \tilde{X}^\top \tilde{Y}$$

# Optimizing ridge regression

- Gradient descent vs. analytical/closed-form/direct solution



- Accuracy doesn't mean anything without running time

- Running time doesn't mean anything without accuracy

- Need to measure accuracy for the running time we have
  - Recall: closed-form solution (if no offset)

$$\theta = (\tilde{X}^\top \tilde{X} + n\lambda I)^{-1} \tilde{X}^\top \tilde{Y}$$

# Optimizing ridge regression

- Gradient descent vs. analytical/closed-form/direct solution



- Accuracy doesn't mean anything without running time

- Running time doesn't mean anything without accuracy

- Need to measure accuracy for the running time we have
  - Recall: closed-form solution (if no offset)

$$\theta = (\tilde{X}^\top \tilde{X} + n\lambda I)^{-1} \tilde{X}^\top \tilde{Y}$$

$\underbrace{\qquad\qquad}_{\text{dxd}}$

# Optimizing ridge regression

- Gradient descent vs. analytical/closed-form/direct solution



- Accuracy doesn't mean anything without running time

- Running time doesn't mean anything without accuracy

- Need to measure accuracy for the running time we have

  - Recall: closed-form solution (if no offset)

$$\theta = (\underbrace{\tilde{X}^\top \tilde{X} + n\lambda I}_{\text{dxd}})^{-1} \tilde{X}^\top \tilde{Y}$$

Matrix inversion running time: $O(d^3)$

# Gradient descent for ridge regression

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

```
Gradient-Descent( Θ_init, η, f, ∇_Θ f )
  Initialize Θ^(0) = Θ_init
  Initialize t = 0
  repeat
   t = t + 1
   Θ^(t) = Θ^(t-1) - η∇_Θ f(Θ^(t-1))

  until stopping criterion
  Return Θ^(t)
```

$\text{Gradient-Descent}(\ \Theta_{\text{init}}, \eta, f, \nabla_\Theta f\ )$

$\quad\text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\quad\text{Initialize } t = 0$

$\quad\textbf{repeat}$

$\quad\quad t = t + 1$

$\quad\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta\nabla_\Theta f(\Theta^{(t-1)})$

$\quad\textbf{until}\text{ stopping criterion}$

$\quad\textbf{Return } \Theta^{(t)}$

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

Gradient-Descent( $\Theta_{\text{init}}, \eta, f, \nabla_\Theta f$ )

Initialize $\Theta^{(0)} = \Theta_{\text{init}}$

Initialize t = 0

**repeat**

  t = t + 1

  $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** stopping criterion

**Return** $\Theta^{(t)}$

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

Gradient-Descent( $\Theta_{\text{init}}, \eta, f, \nabla_\Theta f$ )

Initialize $\theta^{(0)} = \theta_{\text{init}}$

Initialize t = 0

**repeat**

  t = t + 1

  $\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** stopping criterion

**Return** $\Theta^{(t)}$

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

Gradient-Descent( $\Theta_{\text{init}}, \eta, f, \nabla_\Theta f$ )

Initialize $\theta^{(0)} = \theta_{\text{init}}$

Initialize t = 0

**repeat**

t = t + 1

$$\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$$

**until** stopping criterion

**Return** $\Theta^{(t)}$

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

$\texttt{Gradient-Descent(}\ \Theta_{\text{init}}, \eta, f, \nabla_\Theta f\ \texttt{)}$

$\quad \texttt{Initialize}\ \theta^{(0)} = \theta_{\text{init}}$

$\quad \texttt{Initialize t = 0}$

$\quad \textbf{repeat}$

$\quad\quad \texttt{t = t + 1}$

$\quad\quad \Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

$\quad \textbf{until}\ \texttt{stopping criterion}$

$\quad \textbf{Return}\ \Theta^{(t)}$

8

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

$\texttt{Gradient-Descent(}\ \Theta_{\text{init}}, \eta, f, \nabla_\Theta f\ \texttt{)}$

$\quad \texttt{Initialize}\ \theta^{(0)} = \theta_{\text{init}}$

$\quad \texttt{Initialize t = 0}$

**repeat**

$\quad \texttt{t = t + 1}$

$\quad \boxed{\Theta^{(t)} = \Theta^{(t-1)}} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** stopping criterion

**Return** $\Theta^{(t)}$

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

$\texttt{Gradient-Descent(}\ \Theta_{\text{init}}, \eta, f, \nabla_\Theta f\ \texttt{)}$

$\quad \texttt{Initialize}\ \theta^{(0)} = \theta_{\text{init}}$

$\quad \texttt{Initialize t = 0}$

**repeat**

$\quad \texttt{t = t + 1}$

$\quad \boxed{\theta^{(t)} = \theta^{(t-1)}} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

**until** stopping criterion

**Return** $\Theta^{(t)}$

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

```
Gradient-Descent( Θ_init, η, f, ∇_Θ f )
  Initialize θ^(0) = θ_init
  Initialize t = 0
  repeat
   t = t + 1
```

$$\theta^{(t)} = \theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$$

```
  until stopping criterion
  Return Θ^(t)
```

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

$\texttt{Gradient-Descent(}\ \Theta_{\text{init}}, \eta, f, \nabla_\Theta f\ \texttt{)}$

  Initialize $\theta^{(0)} = \theta_{\text{init}}$

  Initialize t = 0

  **repeat**

   t = t + 1

   $\theta^{(t)} = \theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$

  **until** stopping criterion

  **Return** $\Theta^{(t)}$

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

```
Gradient-Descent(
```
$\Theta_{\text{init}}, \eta, f, \nabla_\Theta f$
```
)
```

```
  Initialize
```
$\theta^{(0)} = \theta_{\text{init}}$

```
  Initialize t = 0
```

**repeat**

```
  t = t + 1
```

$$\theta^{(t)} = \theta^{(t-1)} - \eta \left\{ \frac{1}{n} \sum_{i=1}^{n} 2\left[\theta^{(t-1)\top} x^{(i)} - y^{(i)}\right] x^{(i)} + 2\lambda\theta^{(t-1)} \right\}$$

**until** stopping criterion
**Return** $\Theta^{(t)}$

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

```
Gradient-Descent( Θ_init, η, f, ∇_Θ f )
```

$\quad$ Initialize $\theta^{(0)} = \theta_{\text{init}}$

$\quad$ Initialize t = 0

$\quad$ **repeat**

$\quad\quad$ t = t + 1

$$\theta^{(t)} = \theta^{(t-1)} - \eta \left\{ \frac{1}{n} \sum_{i=1}^{n} 2 \left[ \theta^{(t-1)\top} x^{(i)} - y^{(i)} \right] x^{(i)} + 2\lambda \theta^{(t-1)} \right\}$$

Exercise: Check!

$\quad$ **until** stopping criterion

$\quad$ **Return** $\Theta^{(t)}$

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

```
Gradient-Descent( Θ_init, η, f, ∇_Θ f )
```
$\text{Initialize } \theta^{(0)} = \theta_{\text{init}}$

$\text{Initialize t = 0}$

**repeat**

$\text{t = t + 1}$

$$\theta^{(t)} = \theta^{(t-1)} - \eta \left\{ \frac{1}{n} \sum_{i=1}^{n} 2 \left[ \theta^{(t-1)\top} x^{(i)} - y^{(i)} \right] x^{(i)} + 2\lambda \theta^{(t-1)} \right\}$$

Exercise: Check!

**until** stopping criterion

**Return** $\Theta^{(t)}$

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

```
Gradient-Descent( Θinit, η, f, ∇Θf )
```
$\quad$ `Initialize` $\theta^{(0)} = \theta_{\text{init}}$

$\quad$ `Initialize t = 0`

$\quad$ **repeat**

$\qquad$ `t = t + 1`

$$\theta^{(t)} = \theta^{(t-1)} - \eta \left\{ \frac{1}{n} \sum_{i=1}^{n} 2\left[\theta^{(t-1)\top} x^{(i)} - y^{(i)}\right] x^{(i)} + 2\lambda \theta^{(t-1)} \right\}$$

Exercise: Check!

$\quad$ **until** `stopping criterion`

$\quad$ **Return** $\Theta^{(t)}$

8

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

$\texttt{Gradient-Descent(}\ \Theta_{\text{init}}, \eta, f, \nabla_\Theta f\ \texttt{)}$

  $\texttt{Initialize}\ \theta^{(0)} = \theta_{\text{init}}$

  $\texttt{Initialize t = 0}$

  **repeat**

  $\texttt{t = t + 1}$

  $\theta^{(t)} = \theta^{(t-1)} - \eta\left\{ \frac{1}{n} \sum_{i=1}^{n} 2\left[\theta^{(t-1)\top} x^{(i)} - y^{(i)}\right] x^{(i)} + 2\lambda\theta^{(t-1)} \right\}$

  **until** $\texttt{stopping criterion}$

  **Return** $\theta^{(t)}$

Exercise: Check!

8

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

Gradient-Descent( $\Theta_{\text{init}}, \eta, f, \nabla_\Theta f$ )

Initialize $\theta^{(0)} = \theta_{\text{init}}$

Initialize $t = 0$

**repeat**

Exercise: Check!

$t = t + 1$

$$\theta^{(t)} = \theta^{(t-1)} - \eta \left\{ \frac{1}{n} \sum_{i=1}^{n} 2 \left[ \theta^{(t-1)\top} x^{(i)} - y^{(i)} \right] x^{(i)} + 2\lambda \theta^{(t-1)} \right\}$$

**until** stopping criterion

**Return** $\theta^{(t)}$

8

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

`RidgeRegression-Gradient-Descent(`$\theta_{\text{init}}, \eta$`)`

$\quad$`Initialize` $\theta^{(0)} = \theta_{\text{init}}$

$\quad$`Initialize t = 0`

$\quad$**repeat**

Exercise: Check!

$\quad\quad$`t = t + 1`

$$\theta^{(t)} = \theta^{(t-1)} - \eta\left\{\frac{1}{n}\sum_{i=1}^{n} 2\left[\theta^{(t-1)\top}x^{(i)} - y^{(i)}\right]x^{(i)} + 2\lambda\theta^{(t-1)}\right\}$$

$\quad$**until** `stopping criterion`

$\quad$**Return** $\theta^{(t)}$

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

```
RidgeRegression-Gradient-Descent(θ_init, η)
```
Initialize $\theta^{(0)} = \theta_{\text{init}}$

Initialize t = 0

**repeat**

  t = t + 1

$$\theta^{(t)} = \theta^{(t-1)} - \eta\left\{\frac{1}{n}\sum_{i=1}^{n}2\big[\theta^{(t-1)\top}x^{(i)} - y^{(i)}\big]x^{(i)} + 2\lambda\theta^{(t-1)}\right\}$$

Exercise: Check!

**until** stopping criterion

**Return** $\theta^{(t)}$

- No more matrix inversion! (see lab)

8

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

```
RidgeRegression-Gradient-Descent($\theta_{\text{init}}, \eta$)
  Initialize $\theta^{(0)} = \theta_{\text{init}}$
  Initialize t = 0
```
**repeat**
```
  t = t + 1
```
$$\theta^{(t)} = \theta^{(t-1)} - \eta\left\{\frac{1}{n}\sum_{i=1}^{n} 2\left[\theta^{(t-1)\top}x^{(i)} - y^{(i)}\right]x^{(i)} + 2\lambda\theta^{(t-1)}\right\}$$

**until** stopping criterion
**Return** $\theta^{(t)}$

**Exercise: Check!**

- No more matrix inversion! (see lab)
- But have to look at all $n$ data points every step

8

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

RidgeRegression-Gradient-Descent($\theta_{\text{init}}, \eta$)

  Initialize $\theta^{(0)} = \theta_{\text{init}}$

  Initialize t = 0

  **repeat**

  t = t + 1

  $$\theta^{(t)} = \theta^{(t-1)} - \eta\left\{\frac{1}{n}\sum_{i=1}^{n}2\big[\theta^{(t-1)\top}x^{(i)} - y^{(i)}\big]x^{(i)} + 2\lambda\theta^{(t-1)}\right\}$$

  **until** stopping criterion

  **Return** $\theta^{(t)}$

Exercise: Check!

- No more matrix inversion! (see lab)
- But have to look at all *n* data points every step
- How to better handle large *n*?

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

```
RidgeRegression-Gradient-Descent($\theta_{\mathrm{init}}, \eta$)
    Initialize $\theta^{(0)} = \theta_{\mathrm{init}}$
    Initialize t = 0
```
**repeat**

Exercise: Check!

$$\mathtt{t \ = \ t \ + \ 1}$$
$$\theta^{(t)} = \theta^{(t-1)} - \eta \left\{ \frac{1}{n} \sum_{i=1}^{n} 2\left[ \theta^{(t-1)\top} x^{(i)} - y^{(i)} \right] x^{(i)} + 2\lambda \theta^{(t-1)} \right\}$$

**until** stopping criterion

**Return** $\theta^{(t)}$

- No more matrix inversion! (see lab)
- But have to look at all $n$ data points every step
- How to better handle large $n$?



8

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

```
RidgeRegression-Gradient-Descent(θ_init, η)
   Initialize θ⁽⁰⁾ = θ_init
   Initialize t = 0
```

**repeat**

$$\texttt{t = t + 1}$$

$$\theta^{(t)} = \theta^{(t-1)} - \eta \left\{ \frac{1}{n} \sum_{i=1}^{n} 2 \left[ \theta^{(t-1)\top} x^{(i)} - y^{(i)} \right] x^{(i)} + 2\lambda \theta^{(t-1)} \right\}$$

**until** stopping criterion

**Return** $\theta^{(t)}$

Exercise: Check!

- No more matrix inversion! (see lab)
- But have to look at all $n$ data points every step
- How to better handle large $n$?



*Pollution level* — $y$

*Satellite reading* — $x_1$

8

# Gradient descent for ridge regression

- Gradient descent with $f$ = ridge regression objective
  - For the moment, assume no offset (can extend)

```
RidgeRegression-Gradient-Descent($\theta_{\mathrm{init}}, \eta$)
  Initialize $\theta^{(0)} = \theta_{\mathrm{init}}$
  Initialize t = 0
  repeat
    t = t + 1
```

$$\theta^{(t)} = \theta^{(t-1)} - \eta \left\{ \frac{1}{n} \sum_{i=1}^{n} 2 \left[ \theta^{(t-1)\top} x^{(i)} - y^{(i)} \right] x^{(i)} + 2\lambda \theta^{(t-1)} \right\}$$

```
  until stopping criterion
  Return $\theta^{(t)}$
```

**Exercise: Check!**

- No more matrix inversion! (see lab)
- But have to look at all $n$ data points every step
- How to better handle large $n$?

8

Plot with axes: $y$ (Pollution level) vertical, $x_1$ (Satellite reading) horizontal, showing data points and a fitted line.

# Stochastic gradient descent

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\Theta)$$

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\Theta)$$

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\Theta)$$

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\Theta)$$

- Stay tuned for more examples

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):
$$J_{\mathrm{linreg}}(\Theta) = J_{\mathrm{linreg}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:
$$f(\Theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

Stochastic-Gradient-Descent$(\Theta_{\mathrm{init}}, \eta, T)$

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

Stochastic-Gradient-Descent($\Theta_{\text{init}}, \eta, T$)

   Initialize $\Theta^{(0)} = \Theta_{\text{init}}$

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

Stochastic-Gradient-Descent$(\Theta_{\text{init}}, \eta, T)$

Initialize $\Theta^{(0)} = \Theta_{\text{init}}$

**for** t = 1 **to** T

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

```
Stochastic-Gradient-Descent(Θ_init, η, T)
  Initialize Θ^(0) = Θ_init
  for t = 1 to T
    randomly select i from {1,…,n}
```
(with equal probability)

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n}\sum_{i=1}^{n}(\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n}\sum_{i=1}^{n} f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

```
Stochastic-Gradient-Descent(Θ_init, η, T)
  Initialize Θ^(0) = Θ_init
  for t = 1 to T
    randomly select i from {1,…,n}  (with equal
                                     probability)
    Θ^(t) = Θ^(t−1) − η(t)∇_Θ f_i(Θ^(t−1))
```

9

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n}\sum_{i=1}^{n}(\theta^{\top}x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n}\sum_{i=1}^{n}f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

```
Stochastic-Gradient-Descent(Θ_init, η, T)
  Initialize Θ^(0) = Θ_init
  for t = 1 to T
    randomly select i from {1,…,n}
```
(with equal probability)

$$\Theta^{(t)} = \Theta^{(t-1)} - \eta(t)\nabla_{\Theta}f_i(\Theta^{(t-1)})$$

Compare to gradient descent update:
$$\Theta^{(t)} = \Theta^{(t-1)} - \eta\nabla_{\Theta}f(\Theta^{(t-1)})$$

9

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

```
Stochastic-Gradient-Descent(Θ_init, η, T)
  Initialize Θ^(0) = Θ_init
  for t = 1 to T
    randomly select i from {1,…,n}
```
(with equal probability)

$$\Theta^{(t)} = \Theta^{(t-1)} - \eta(t) \nabla_\Theta f_i(\Theta^{(t-1)})$$

Compare to gradient descent update:

$$\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$$

9

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n}\sum_{i=1}^{n}(\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n}\sum_{i=1}^{n} f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

```
Stochastic-Gradient-Descent(Θ_init, η, T)
  Initialize Θ^(0) = Θ_init
  for t = 1 to T
    randomly select i from {1,…,n}  (with equal
                                     probability)
```
$$\Theta^{(t)} = \Theta^{(t-1)} - \boxed{\eta(t)}\nabla_\Theta f_i(\Theta^{(t-1)})$$

Compare to gradient descent update:
$$\Theta^{(t)} = \Theta^{(t-1)} - \boxed{\eta}\nabla_\Theta f(\Theta^{(t-1)})$$

9

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n}\sum_{i=1}^{n}(\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n}\sum_{i=1}^{n} f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

```
Stochastic-Gradient-Descent(Θ_init, η, T)
  Initialize Θ^(0) = Θ_init
  for t = 1 to T
    randomly select i from {1,…,n}  (with equal
                                     probability)
```
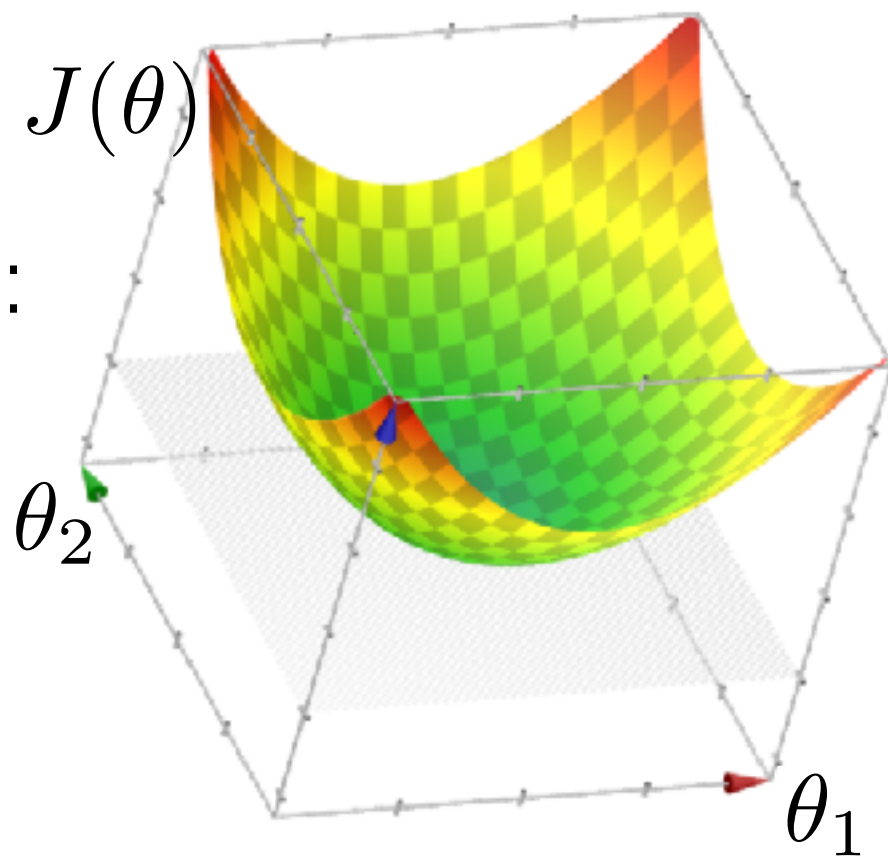
$$\Theta^{(t)} = \Theta^{(t-1)} - \eta(t)\nabla_\Theta f_i(\Theta^{(t-1)})$$

Compare to gradient descent update:

$$\Theta^{(t)} = \Theta^{(t-1)} - \eta\nabla_\Theta f(\Theta^{(t-1)})$$

9

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n}\sum_{i=1}^{n}(\theta^{\top}x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n}\sum_{i=1}^{n}f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

```
Stochastic-Gradient-Descent(Θinit, η, T)
  Initialize Θ(0) = Θinit
  for t = 1 to T
    randomly select i from {1,…,n}
```

$$\Theta^{(t)} = \Theta^{(t-1)} - \eta(t)\nabla_{\Theta}f_i(\Theta^{(t-1)})$$

(with equal probability)

Compare to gradient descent update:

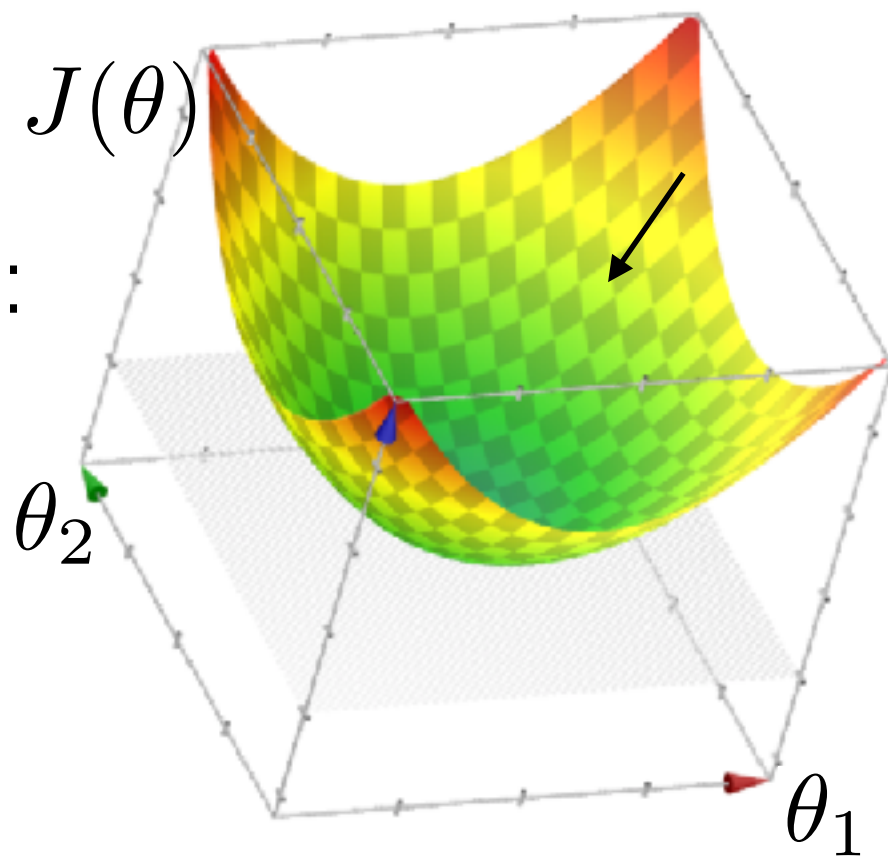$$\Theta^{(t)} = \Theta^{(t-1)} - \eta\nabla_{\Theta}f(\Theta^{(t-1)})$$

9

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n}\sum_{i=1}^{n}(\theta^{\top}x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n}\sum_{i=1}^{n}f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

```
Stochastic-Gradient-Descent(Θ_init, η, T)
  Initialize Θ^(0) = Θ_init
  for t = 1 to T
    randomly select i from {1,…,n}
    Θ^(t) = Θ^(t-1) − η(t)∇_Θ f_i(Θ^(t-1))
  Return Θ^(t)
```

$\text{Initialize } \Theta^{(0)} = \Theta_{\text{init}}$

$\text{randomly select i from } \{1,...,n\}$ (with equal probability)

$\Theta^{(t)} = \Theta^{(t-1)} - \eta(t)\nabla_\Theta f_i(\Theta^{(t-1)})$

**Return** $\Theta^{(t)}$

Compare to gradient descent update:
$$\Theta^{(t)} = \Theta^{(t-1)} - \eta\nabla_\Theta f(\Theta^{(t-1)})$$

# Stochastic gradient descent

- Linear regression objective (with $\lambda = 0$):

$$J_{\text{linreg}}(\Theta) = J_{\text{linreg}}(\theta, \theta_0) = \frac{1}{n} \sum_{i=1}^{n} (\theta^\top x^{(i)} + \theta_0 - y^{(i)})^2$$

- A common machine learning objective:

$$f(\Theta) = \frac{1}{n} \sum_{i=1}^{n} f_i(\Theta)$$

- Stay tuned for more examples
- Compare to training error defn.

```
Stochastic-Gradient-Descent(Θ_init, η, T)
  Initialize Θ^(0) = Θ_init
  for t = 1 to T
    randomly select i from {1,…,n}  (with equal
                                     probability)
    Θ^(t) = Θ^(t-1) − η(t)∇_Θ f_i(Θ^(t-1))
  Return Θ^(t)
```

Compare to gradient descent update:
$$\Theta^{(t)} = \Theta^{(t-1)} - \eta \nabla_\Theta f(\Theta^{(t-1)})$$

- Commonly used with "minibatches"
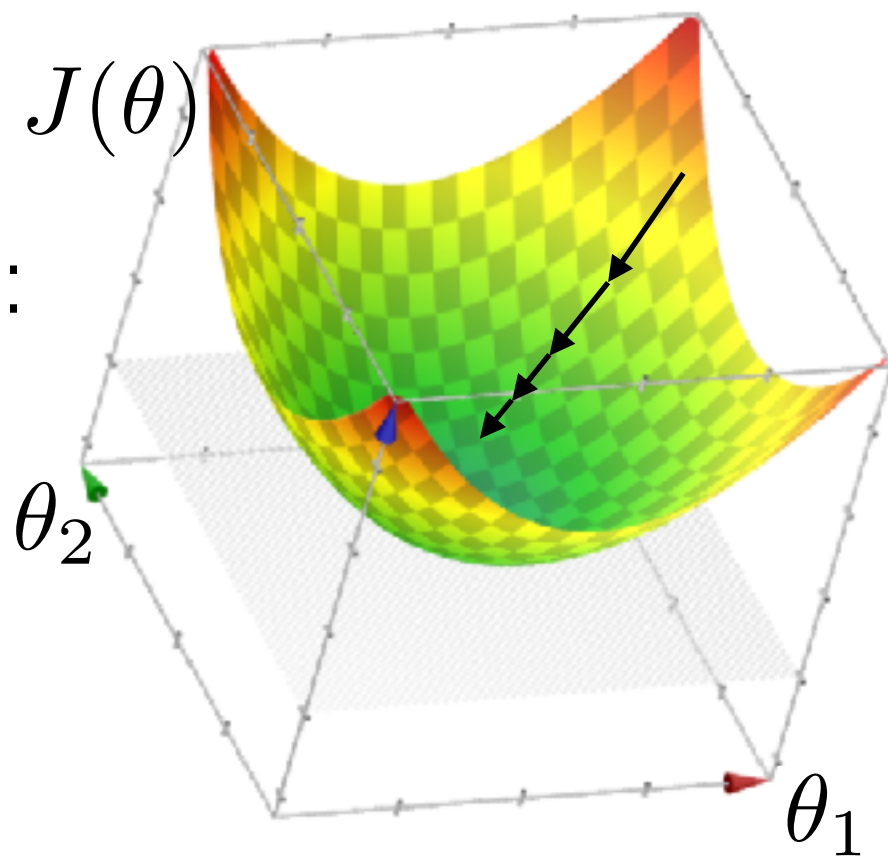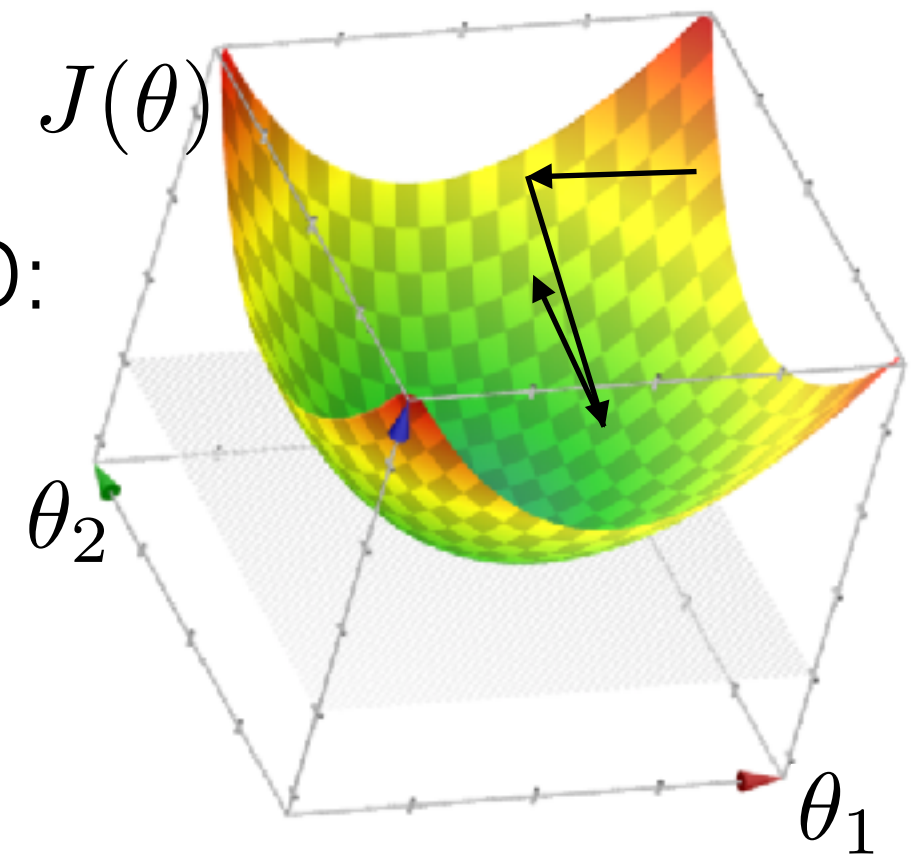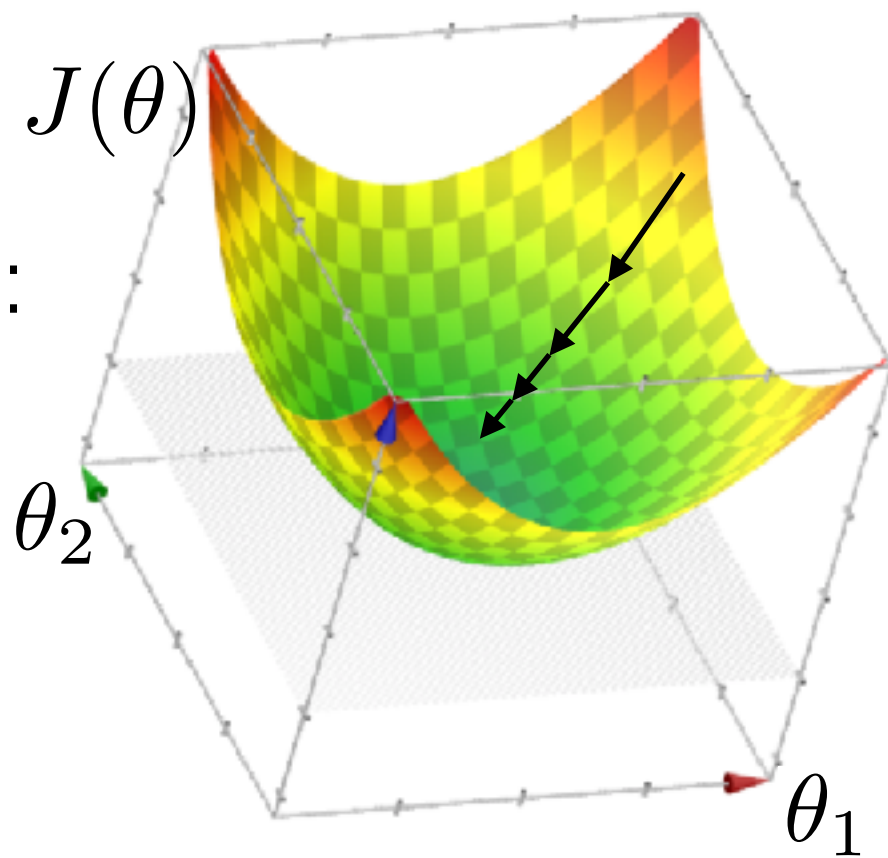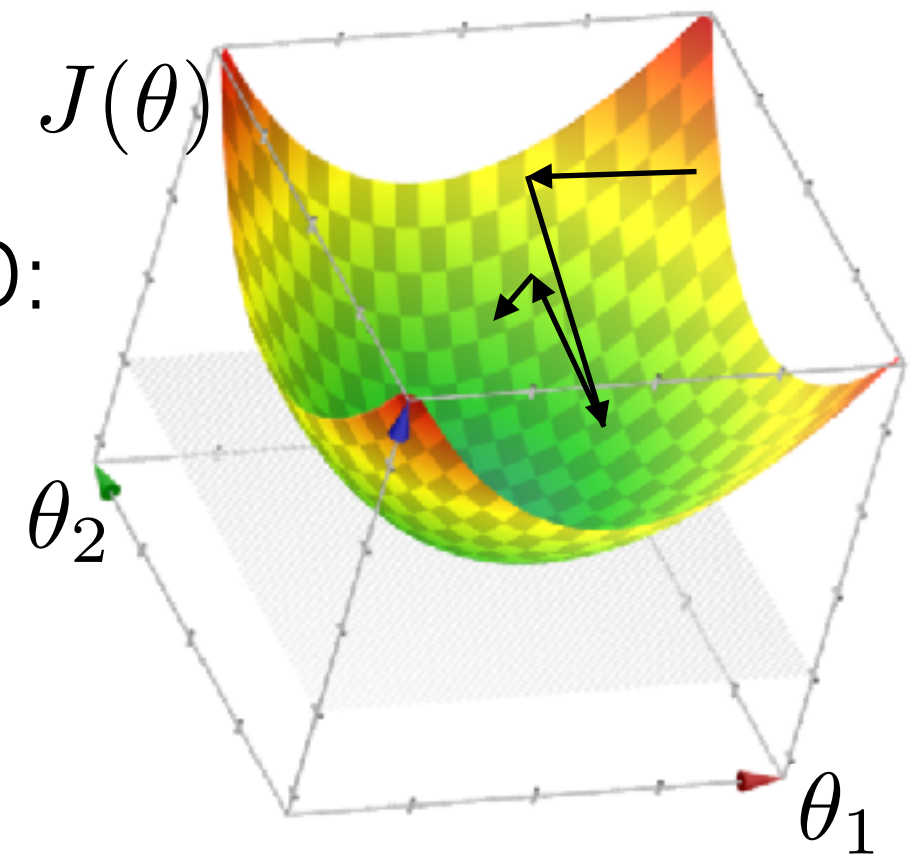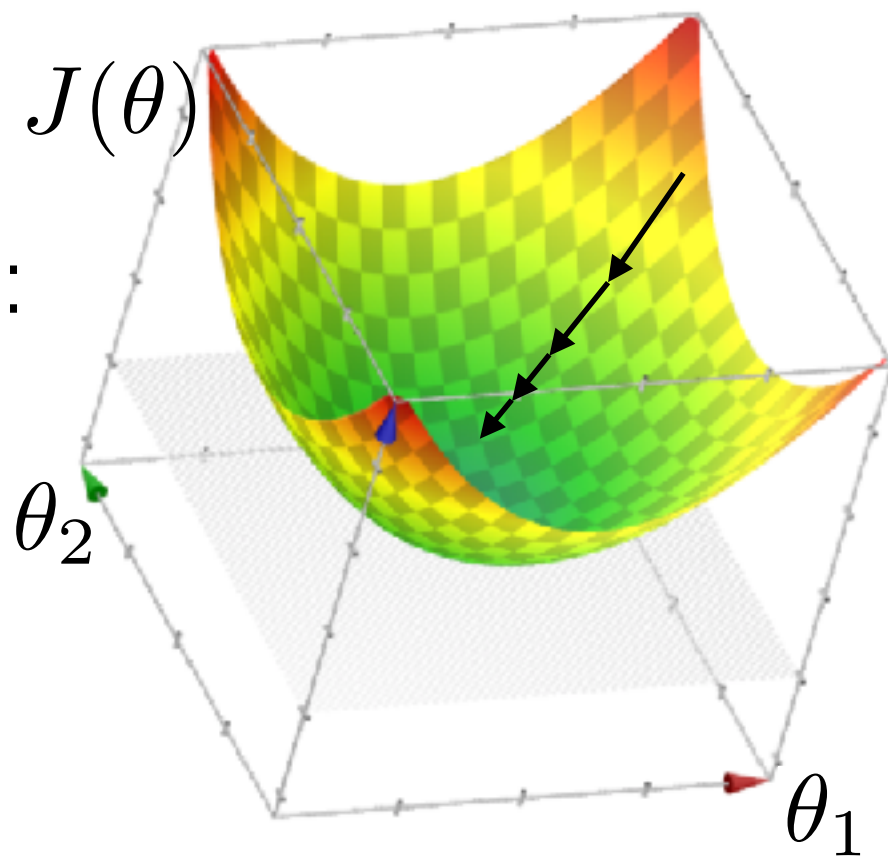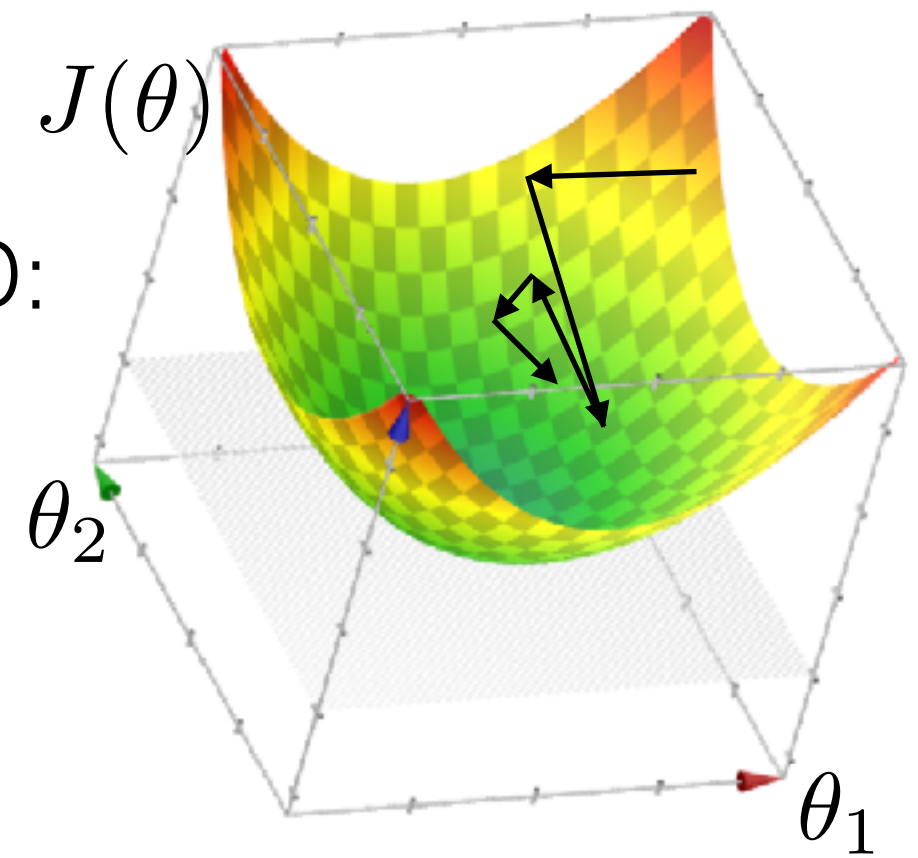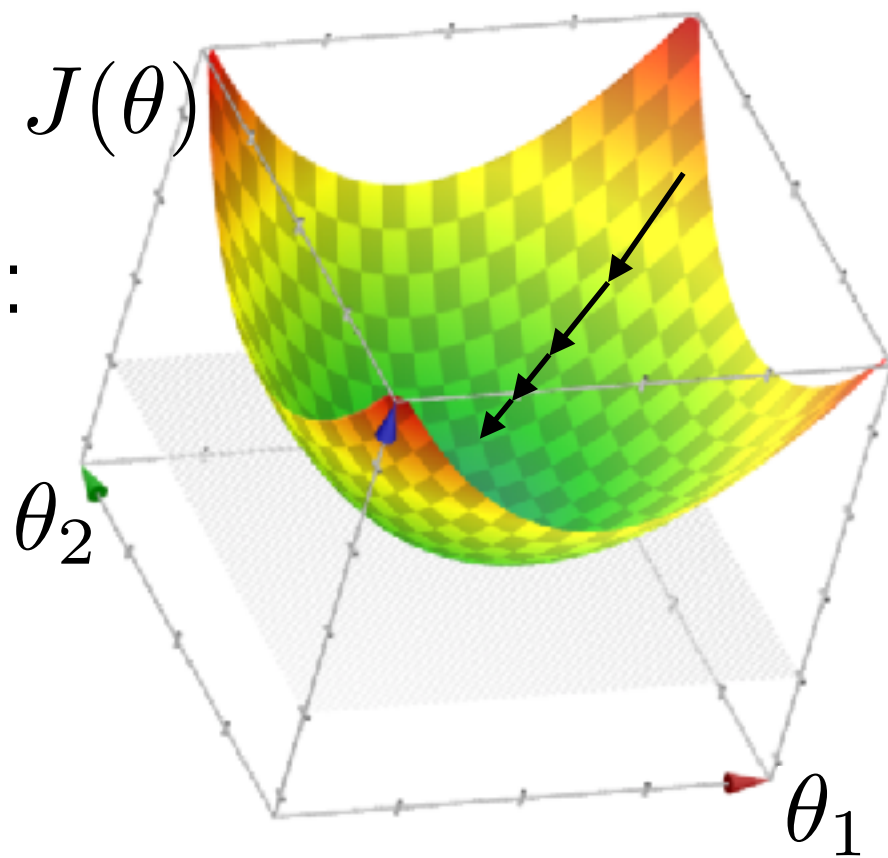
# Stochastic gradient descent (SGD) properties



- GD:

# Stochastic gradient descent (SGD) properties



- GD:

# Stochastic gradient descent (SGD) properties

- GD:



$J(\theta)$

$\theta_2$

$\theta_1$

# Stochastic gradient descent (SGD) properties



- GD:

# Stochastic gradient descent (SGD) properties

$J(\theta)$

- GD:



$\theta_2$

$\theta_1$

# Stochastic gradient descent (SGD) properties



- GD:

# Stochastic gradient descent (SGD) properties

- GD:



- SGD:

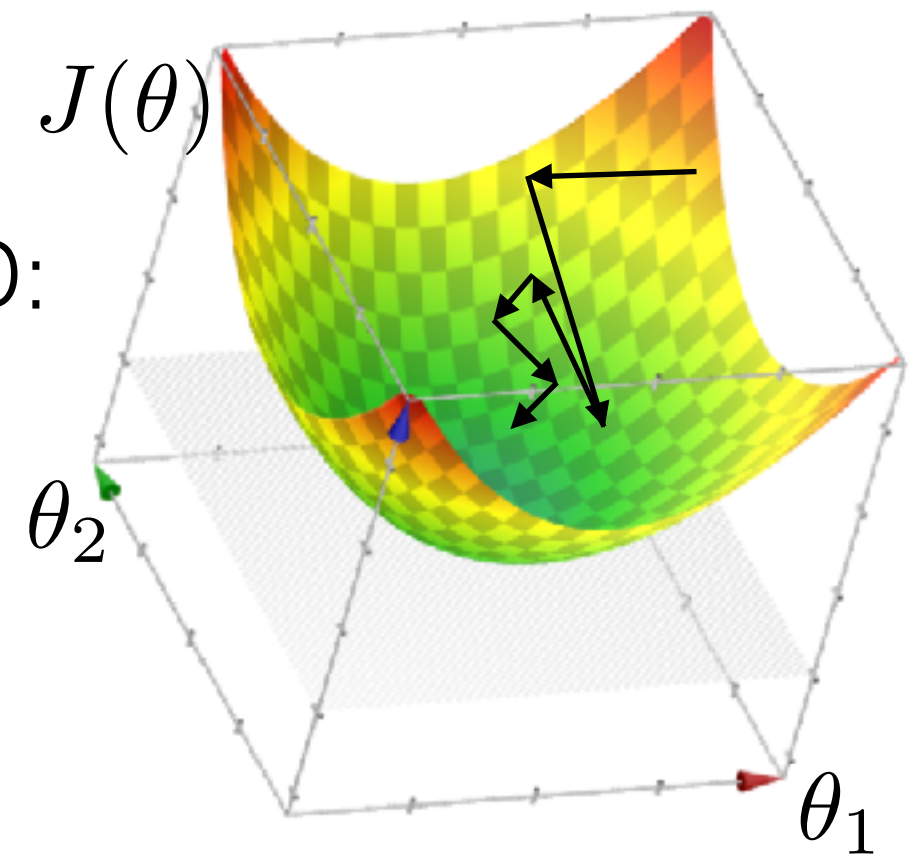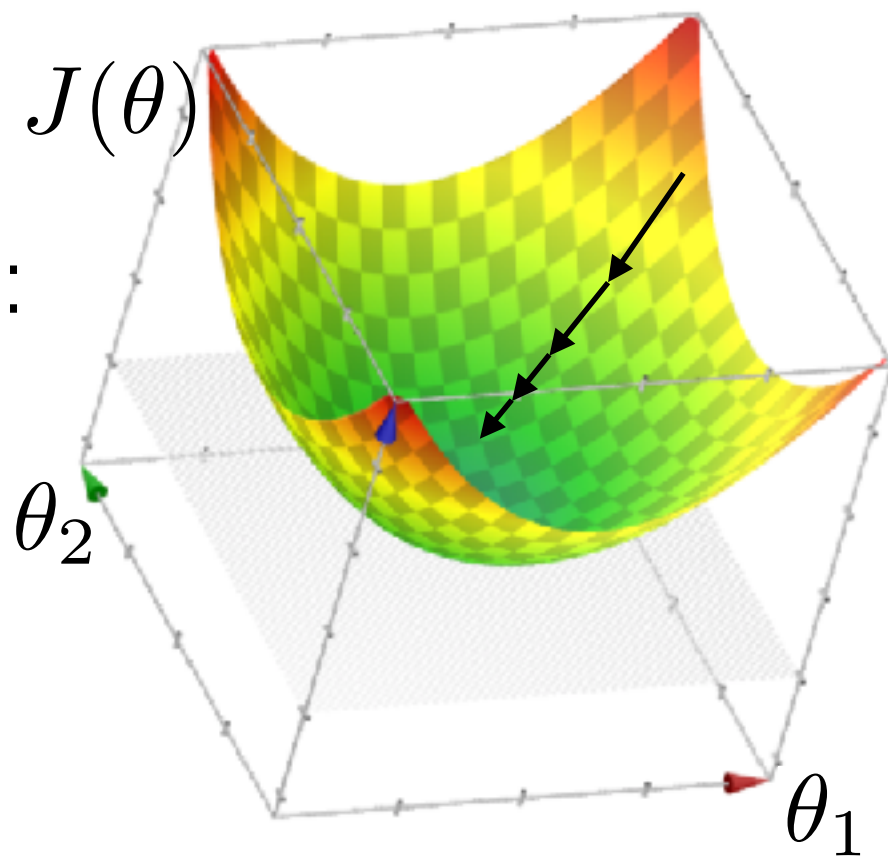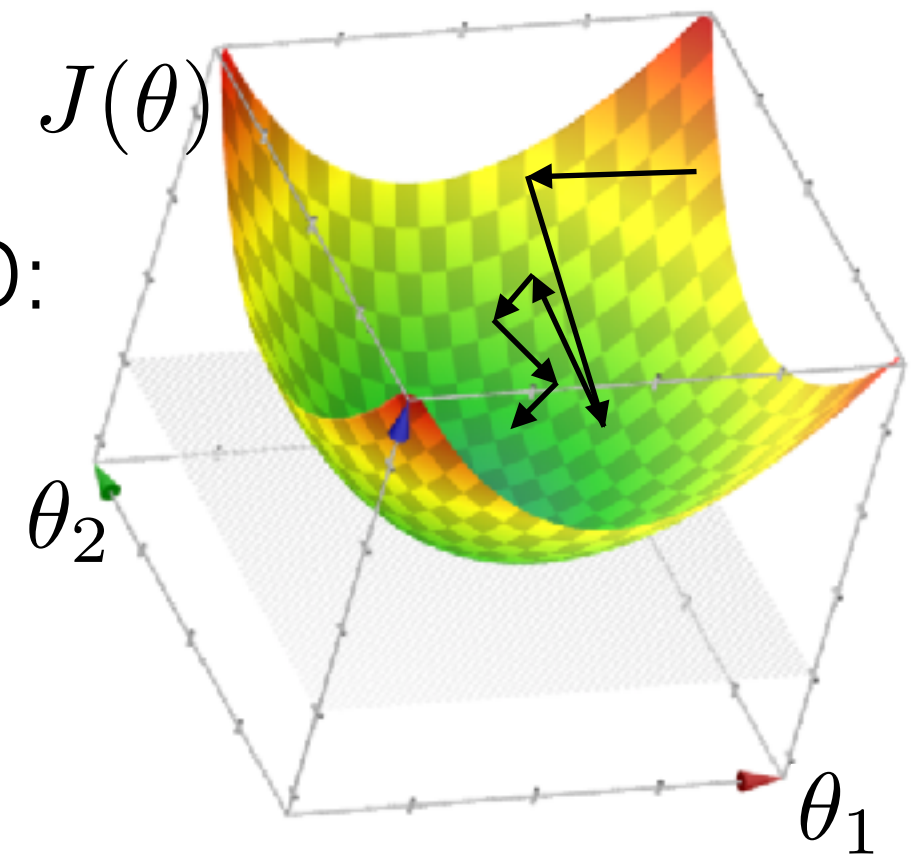# Stochastic gradient descent (SGD) properties

- GD:



- SGD:

# Stochastic gradient descent (SGD) properties

- GD:



- SGD:

# Stochastic gradient descent (SGD) properties

- GD:



- SGD:

# Stochastic gradient descent (SGD) properties



- GD:

- SGD:
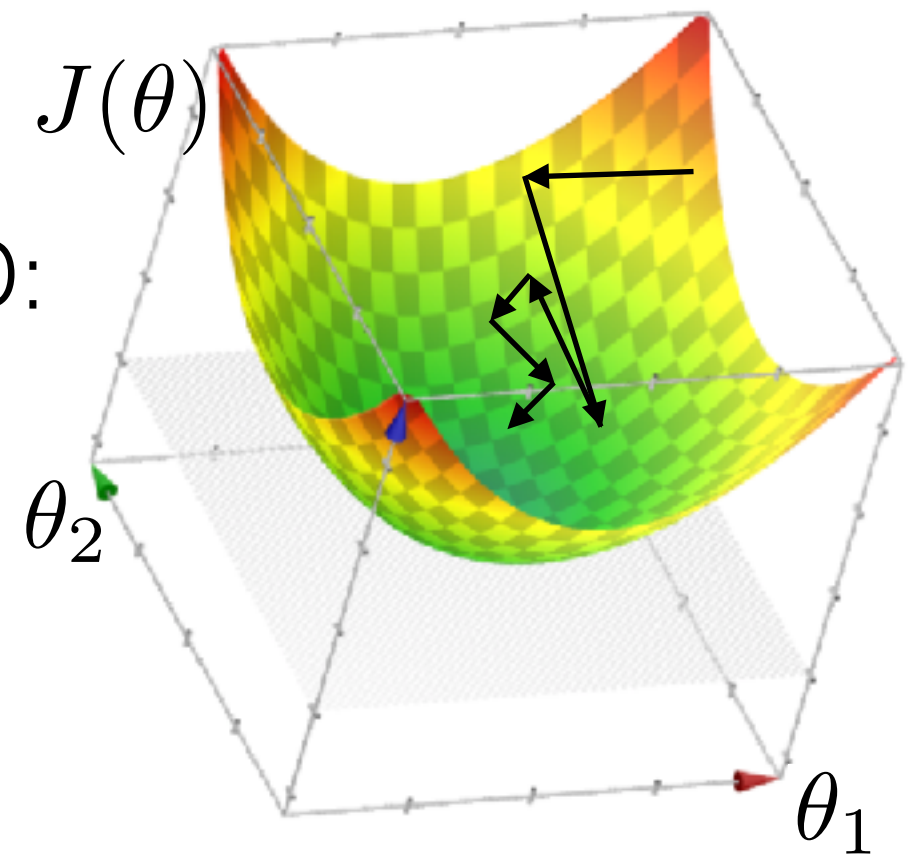
# Stochastic gradient descent (SGD) properties

- GD:



- SGD:

# Stochastic gradient descent (SGD) properties

- GD:



- SGD:

# Stochastic gradient descent (SGD) properties

- GD:



- SGD:

# Stochastic gradient descent (SGD) properties



- GD:
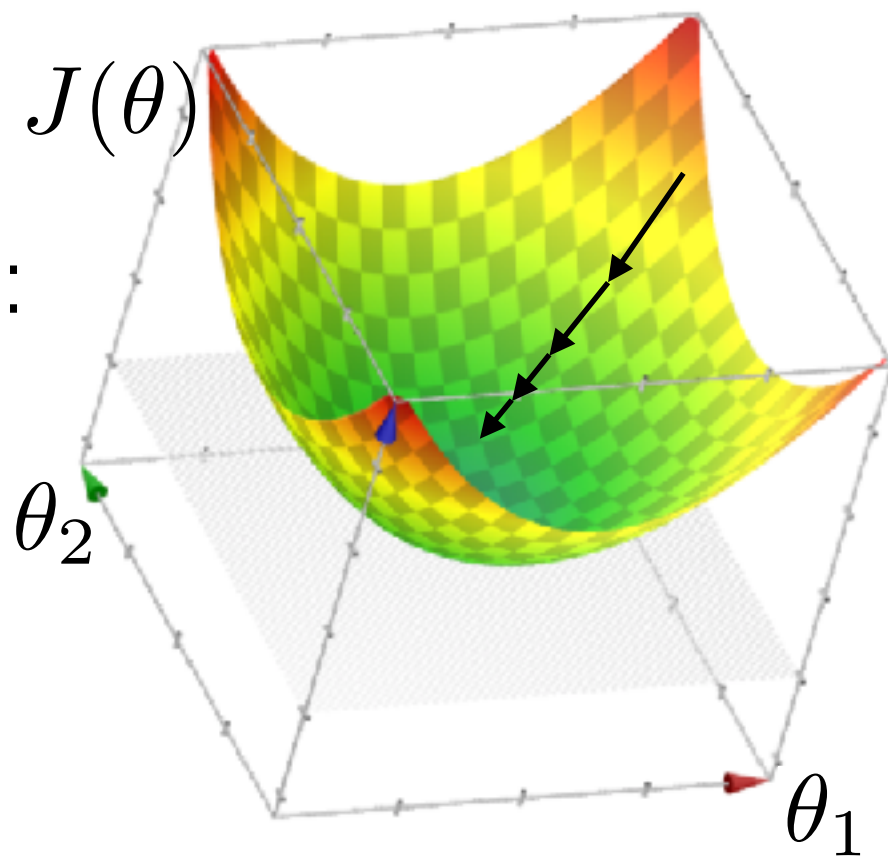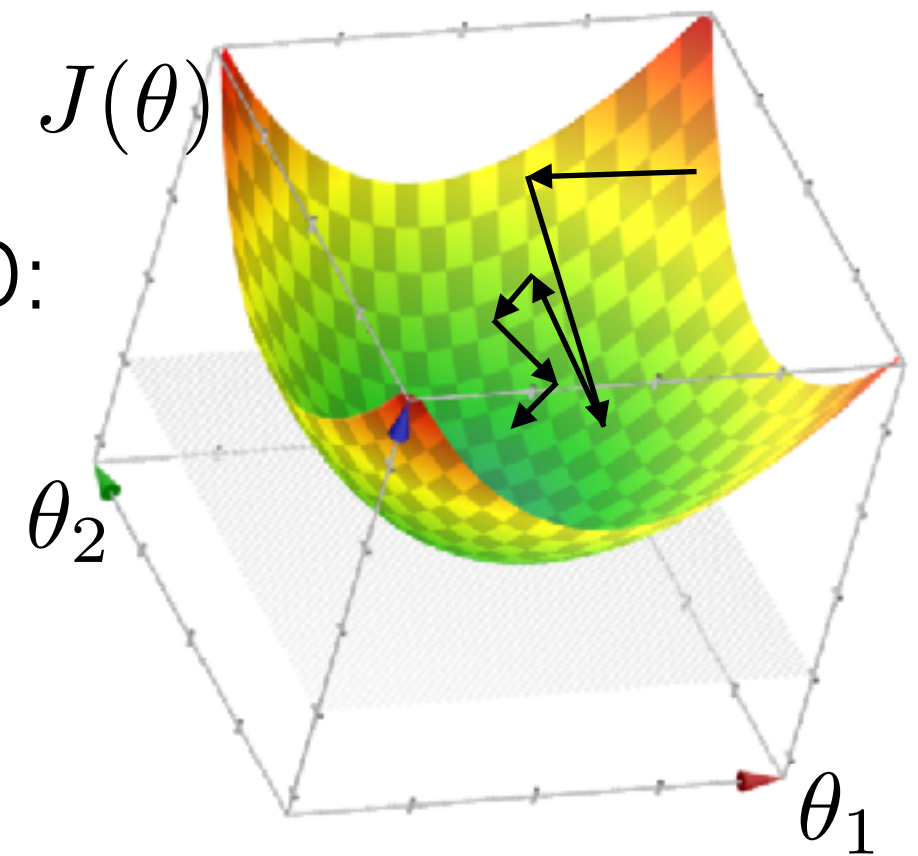
- SGD:

- **Theorem**: SGD performance

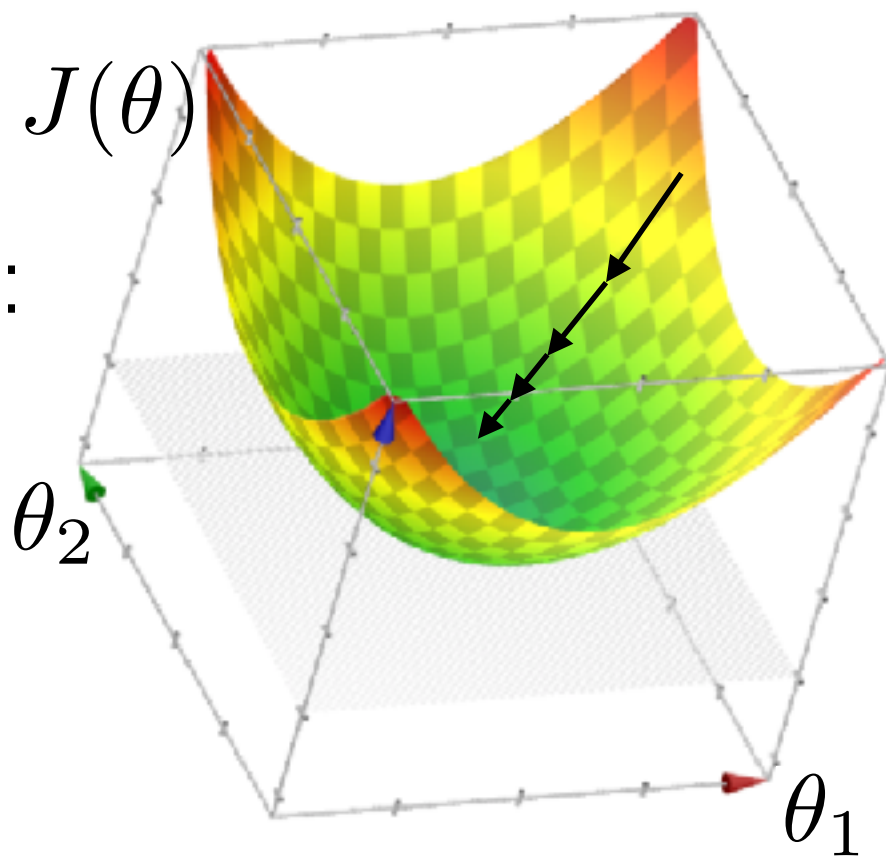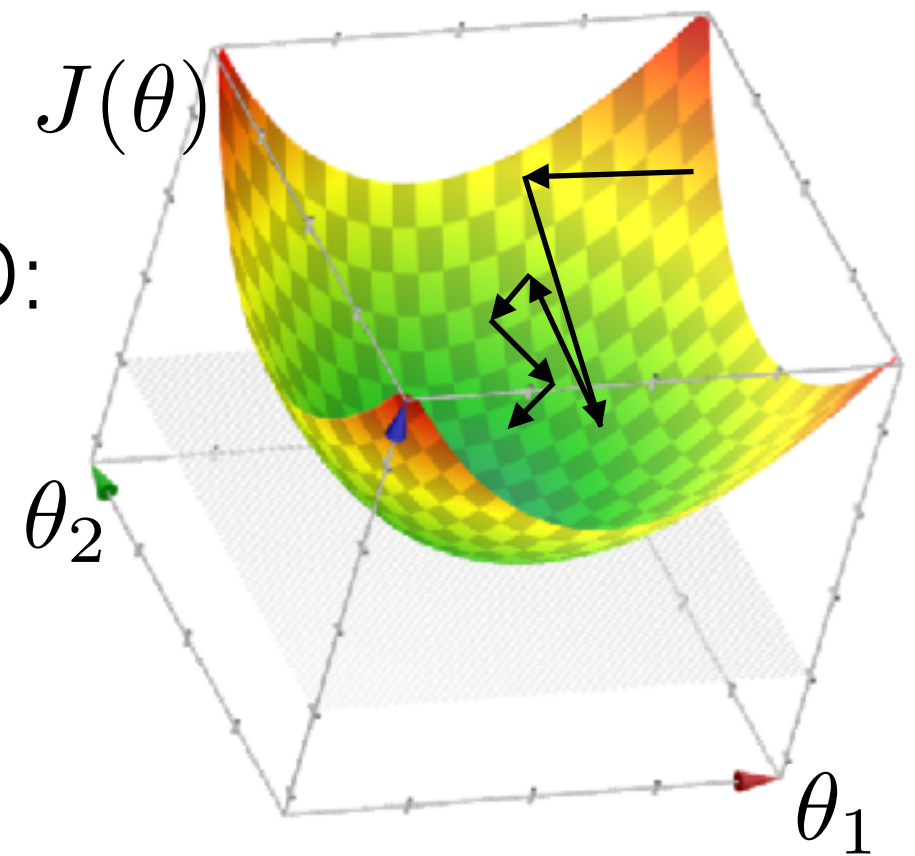# Stochastic gradient descent (SGD) properties



- GD:

- SGD:

- **Theorem**: SGD performance
  - **Assumptions**:

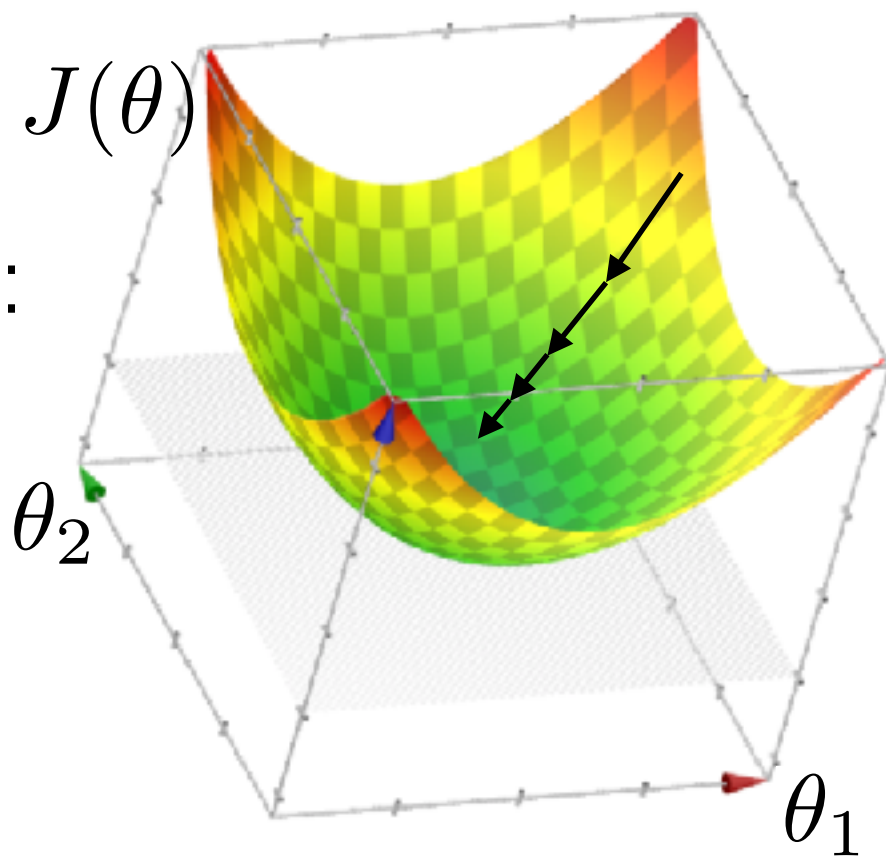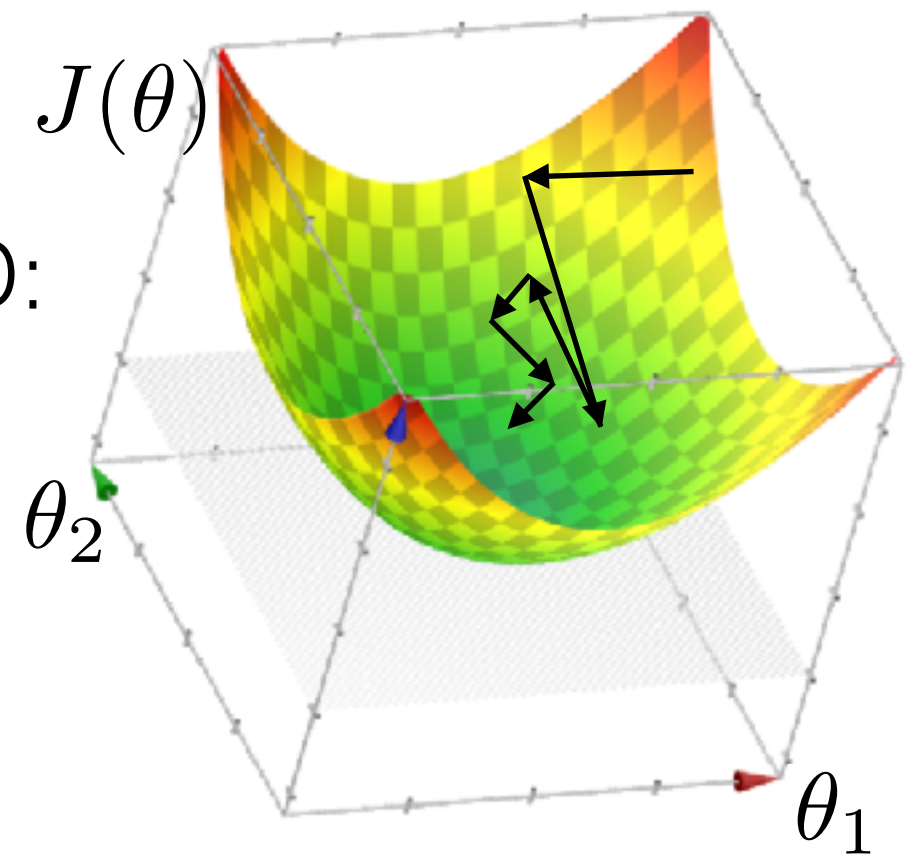# Stochastic gradient descent (SGD) properties



- GD:

- SGD:

- **Theorem**: SGD performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)

# Stochastic gradient descent (SGD) properties

$J(\theta)$

- GD:

$\theta_2$

$\theta_1$

$J(\theta)$

- SGD:
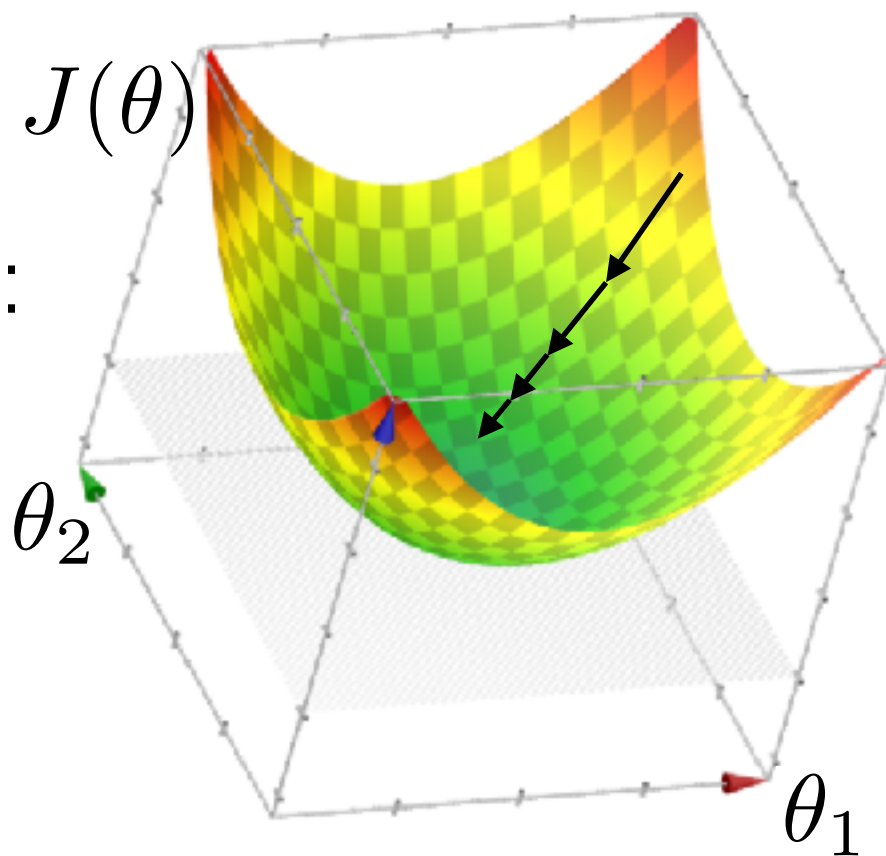
$\theta_2$

$\theta_1$

- **Theorem**: SGD performance
  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is "nice" & convex, has a unique global minimizer

# Stochastic gradient descent (SGD) properties



- GD:

- SGD:

- **Theorem**: SGD performance
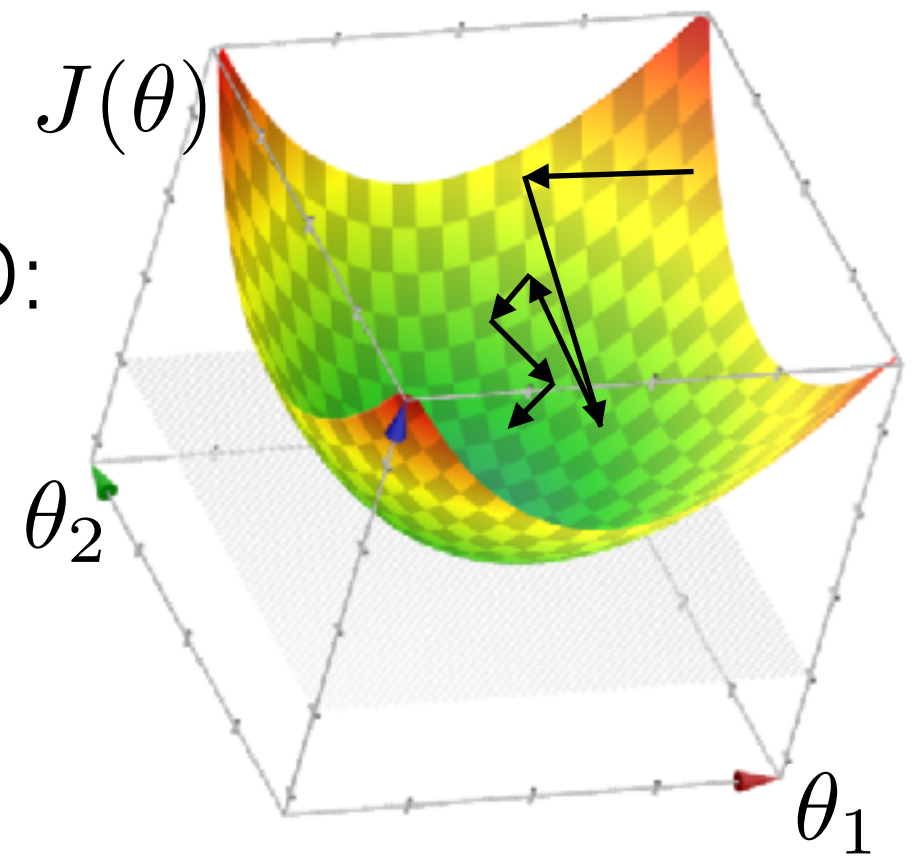  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is "nice" & convex, has a unique global minimizer
    - $\sum\limits_{t=1}^{\infty} \eta(t) = \infty, \sum\limits_{t=1}^{\infty} (\eta(t))^2 < \infty$

# Stochastic gradient descent (SGD) properties

- GD:



$J(\theta)$, $\theta_2$, $\theta_1$

- SGD:



$J(\theta)$, $\theta_2$, $\theta_1$
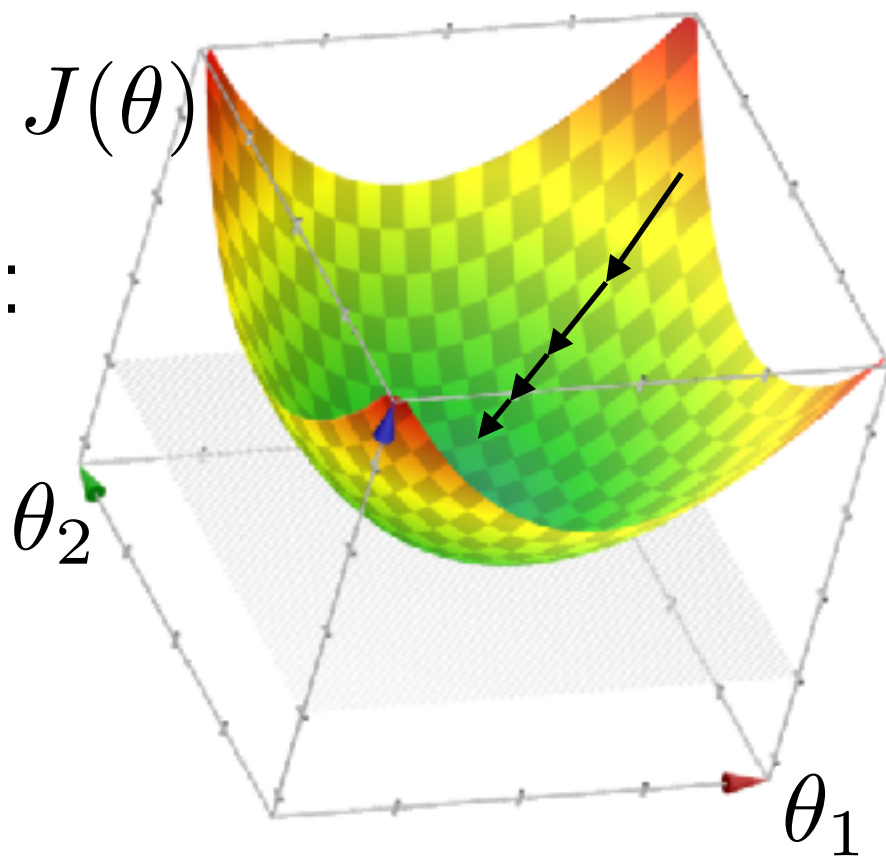
- **Theorem**: SGD performance

  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)

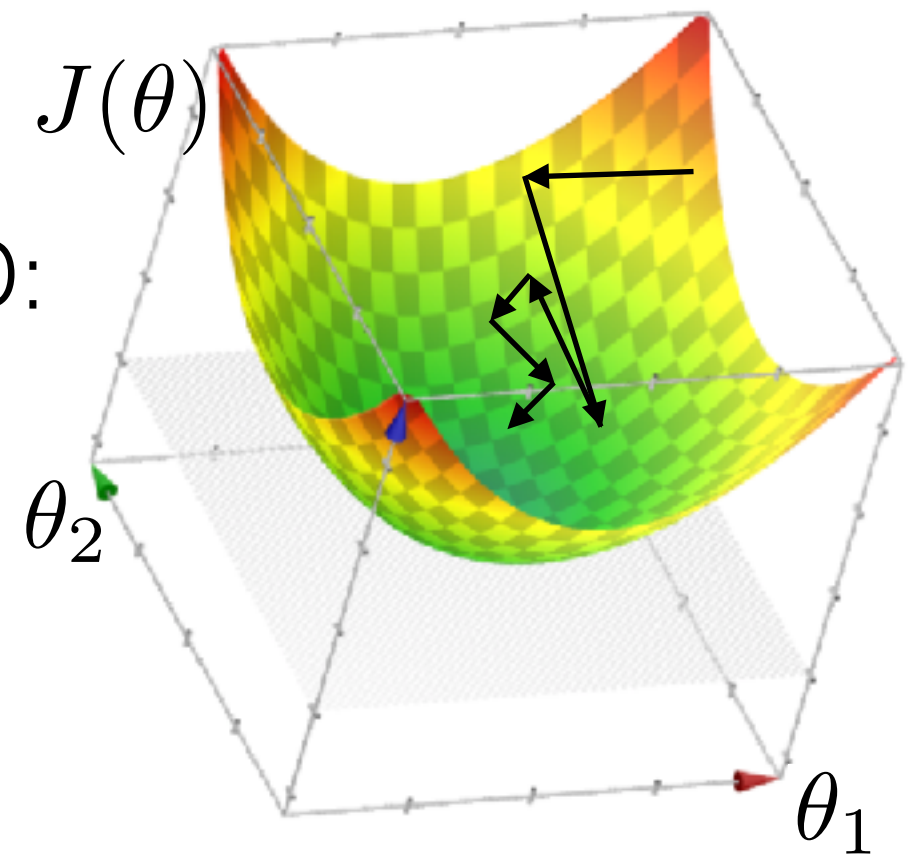    - $f$ is "nice" & convex, has a unique global minimizer

    - $\sum_{t=1}^{\infty} \eta(t) = \infty, \sum_{t=1}^{\infty} (\eta(t))^2 < \infty$

    - e.g. $\eta(t) = \alpha(\tau_0 + t)^{-\kappa} (\kappa \in (0.5, 1])$

# Stochastic gradient descent (SGD) properties



- GD:

- SGD:

- **Theorem**: SGD performance
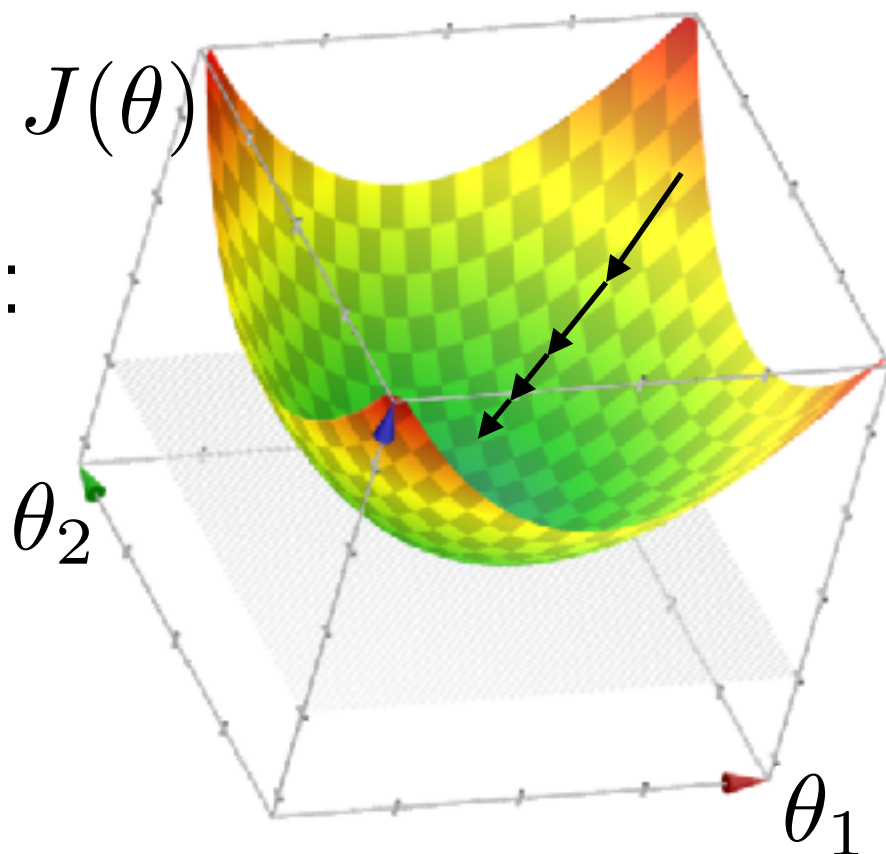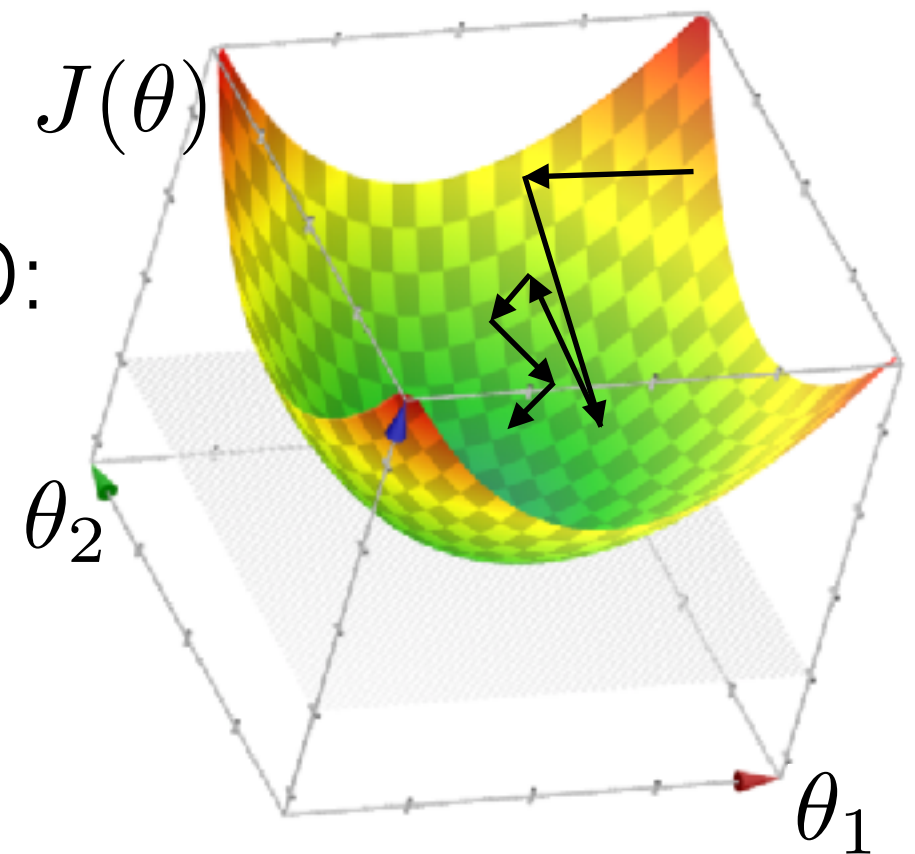  - **Assumptions**: (Choose any $\tilde{\epsilon} > 0$)
    - $f$ is "nice" & convex, has a unique global minimizer
    - $\sum_{t=1}^{\infty} \eta(t) = \infty, \sum_{t=1}^{\infty} (\eta(t))^2 < \infty$
      - e.g. $\eta(t) = \alpha(\tau_0 + t)^{-\kappa} (\kappa \in (0.5, 1])$
  - **Conclusion**: If run long enough, stochastic gradient descent will return a value within $\tilde{\epsilon}$ of the global minimizer