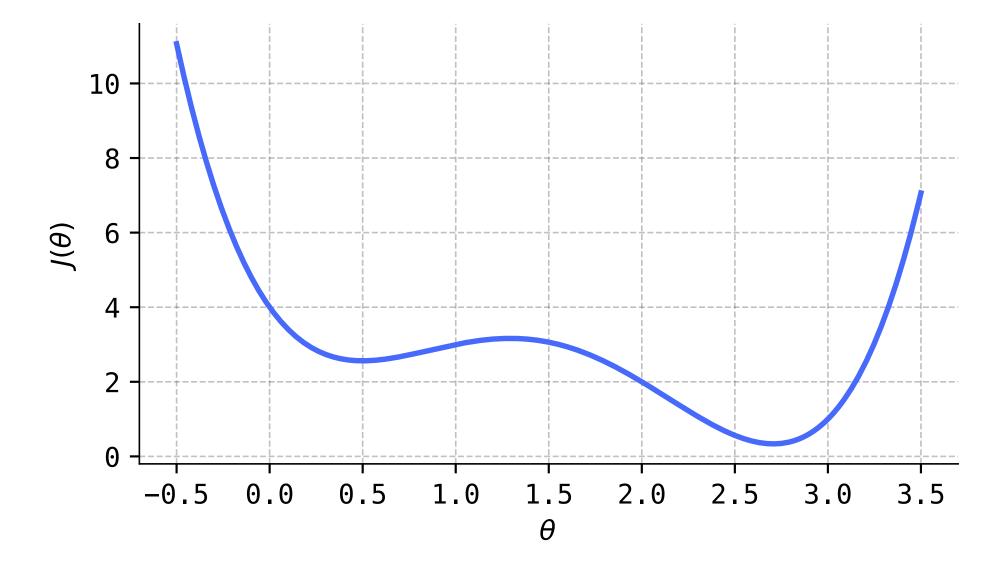# 11. Optimization

# 11.1. Gradient descent in one variable

### Definition 11.1

Let $J : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function. A vector $\boldsymbol{\theta}^\star$ is a *local minimizer* of $J(\boldsymbol{\theta})$ provided that

$$J(\boldsymbol{\theta}^\star) \leq J(\boldsymbol{\theta})$$

for all $\boldsymbol{\theta}$ in a neighborhood of $\boldsymbol{\theta}^\star$; if this inequality holds for *all* $\boldsymbol{\theta}$, then $\boldsymbol{\theta}^\star$ is called a *global minimizer* of $J(\boldsymbol{\theta})$. If we flip the inequality the other direction, then we obtain the definitions of *local* and *global maximizers*. Collectively, local and global minimizers and maximizers of $J(\boldsymbol{\theta})$ are called *extremizers*, and the values $J(\boldsymbol{\theta}^\star)$ of the function where $\boldsymbol{\theta}^\star$ is an extremizer are called *extrema* or *extreme values*.

## 🔔 Algorithm 11.1 (Single-variable gradient descent)

**Input:** A differentiable objective function $J : \mathbb{R} \to \mathbb{R}$, an initial guess $\theta_0 \in \mathbb{R}$ for a local minimizer $\theta^\star$, a learning rate $\alpha > 0$, and the number $N$ of gradient steps.
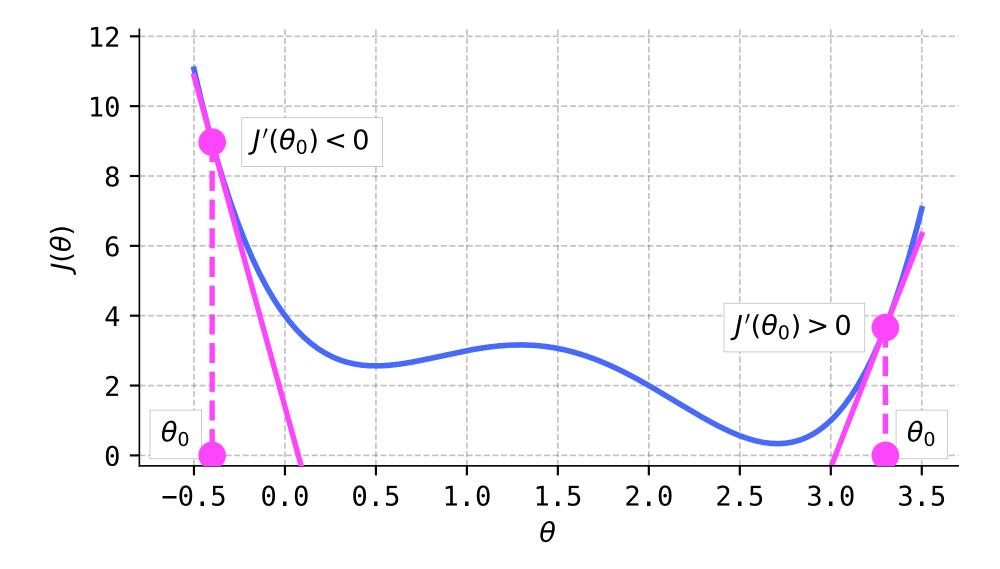
**Output:** An approximation to a local minimizer $\theta^\star$.

$\theta := \theta_0$

For $t$ from $0$ to $N - 1$, do:
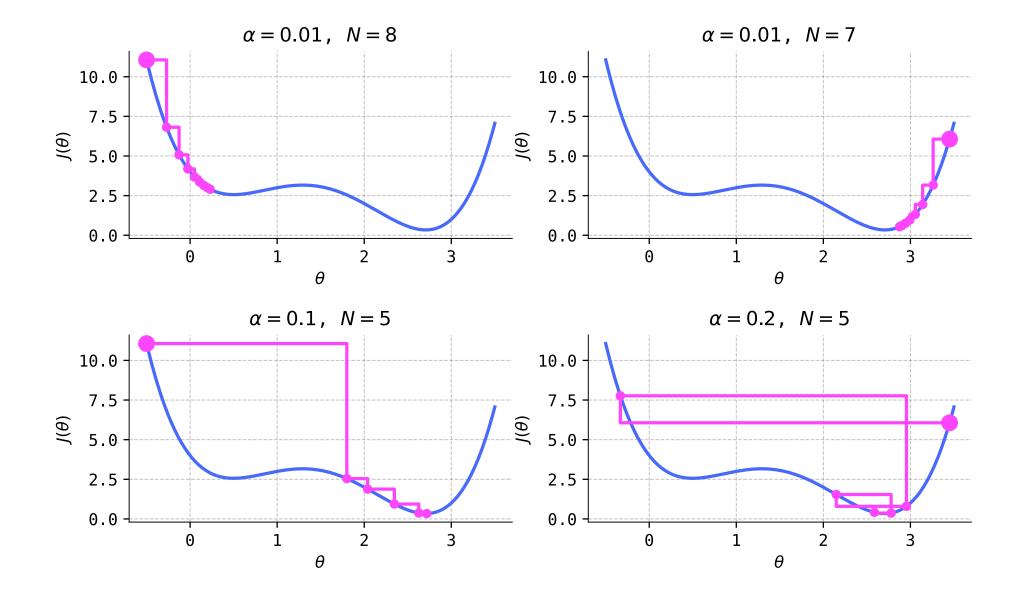
$\quad \theta := \theta - \alpha J'(\theta)$
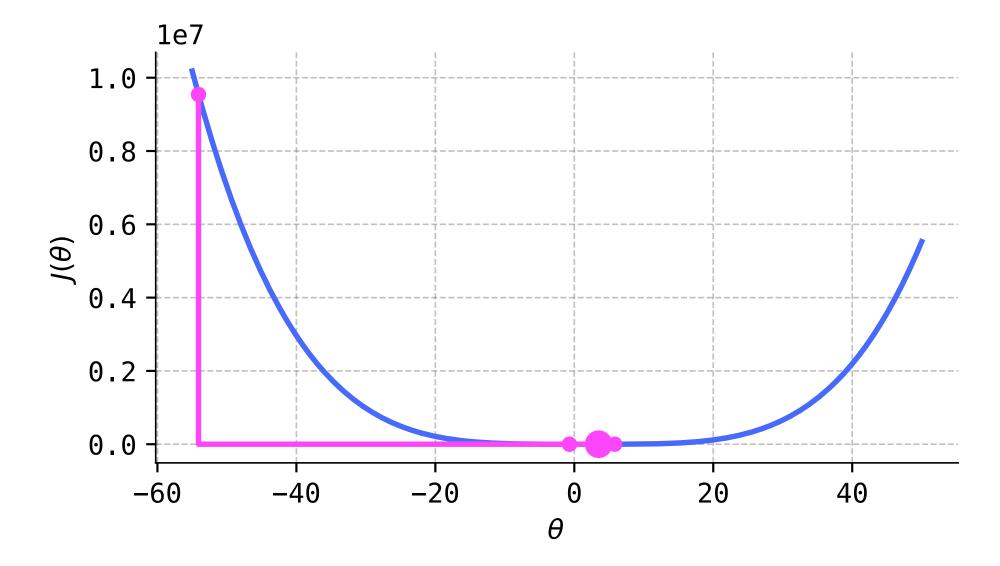
Return $\theta$

**🔔 Observation 11.1**

- The negative derivative $-J'(\theta)$ always "points downhill."
- When the gradient descent algorithm works, it locates a minimizer by following the negative derivative "downhill."

**🔔 Problem Prompt**

Do problem 1 on the worksheet.

**🔔 Problem Prompt**

Do problems 2 and 3 on the worksheet.

## 🔔 Algorithm 11.2 (Single-variable gradient descent with learning rate decay)

**Input:** A differentiable objective function $J : \mathbb{R} \to \mathbb{R}$, an initial guess $\theta_0 \in \mathbb{R}$ for a local minimizer $\theta^\star$, a learning rate $\alpha > 0$, a decay rate $\beta \in [0, 1)$, and the number $N$ of gradient steps.
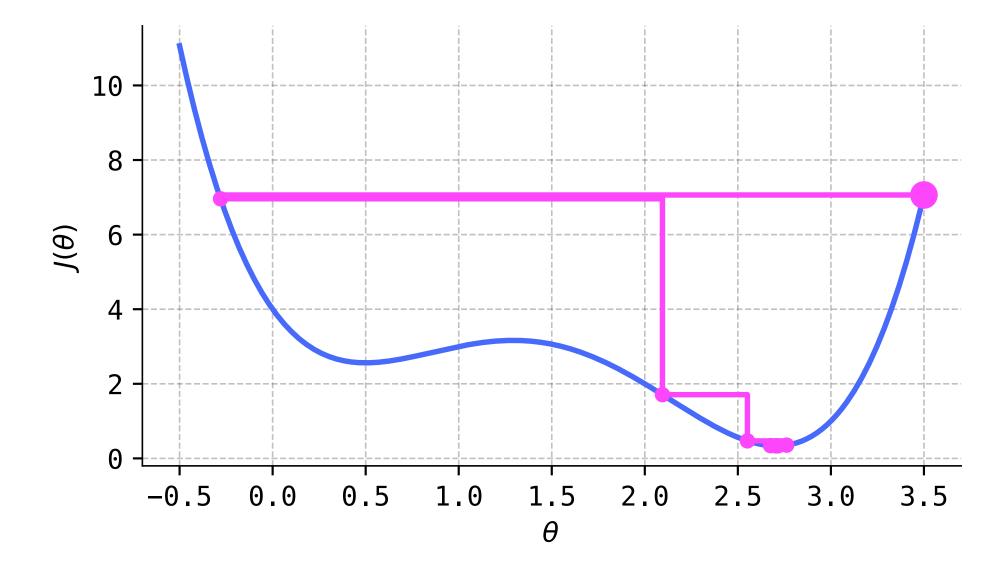
**Output:** An approximation to a local minimizer $\theta^\star$.

$\theta := \theta_0$
For $t$ from $0$ to $N - 1$, do:
$\qquad \theta := \theta - \alpha(1 - \beta)^{t+1} J'(\theta)$
Return $\theta$

**🔔 Problem Prompt**

Do problem 4 on the worksheet.

# 11.2. Differential geometry

convex ⇒ minimizer          concave ⇒ maximizer

$J'(0) = 0, \; J''(0) > 0$

$J'(0) = 0, \; J''(0) < 0$

## Definition 11.2

Let $J : \mathbb{R}^n \to \mathbb{R}$ be a function of class $C^2$, $\boldsymbol{\theta} \in \mathbb{R}^n$ a point, and $\mathbf{v} \in \mathbb{R}^n$ a vector. We define the *directional first derivative of $J$ at $\boldsymbol{\theta}$ in the direction $\mathbf{v}$* to be

$$J'_{\mathbf{v}}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \left. \frac{\mathrm{d}}{\mathrm{d}t} \right|_{t=0} J(t\mathbf{v} + \boldsymbol{\theta}),$$

while we define the *directional second derivative* to be

$$J''_{\mathbf{v}}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \left. \frac{\mathrm{d}^2}{\mathrm{d}t^2} \right|_{t=0} J(t\mathbf{v} + \boldsymbol{\theta}).$$

In this context, the vector $\mathbf{v}$ is called the *directional vector*.

**Problem Prompt**

Do problem 5 on the worksheet.

🔔 **Definition 11.3**

Let $J : \mathbb{R}^n \to \mathbb{R}$ be a function of class $C^2$ and $\boldsymbol{\theta} \in \mathbb{R}^n$ a point. We define the *gradient vector* to be

$$\nabla J(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \left[ \frac{\partial J}{\partial \theta_i}(\boldsymbol{\theta}) \right] = \begin{bmatrix} \dfrac{\partial J}{\partial \theta_1}(\boldsymbol{\theta}) \\ \vdots \\ \dfrac{\partial J}{\partial \theta_n}(\boldsymbol{\theta}) \end{bmatrix} \in \mathbb{R}^n,$$

while we define the the *Hessian matrix* to be

$$\nabla^2 J(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \left[ \frac{\partial^2 J}{\partial \theta_i \partial \theta_j}(\boldsymbol{\theta}) \right] = \begin{bmatrix} \dfrac{\partial^2 J}{\partial \theta_1^2}(\boldsymbol{\theta}) & \cdots & \dfrac{\partial^2 J}{\partial \theta_1 \partial \theta_n}(\boldsymbol{\theta}) \\ \vdots & \ddots & \vdots \\ \dfrac{\partial^2 J}{\partial \theta_n \partial \theta_1}(\boldsymbol{\theta}) & \cdots & \dfrac{\partial^2 J}{\partial \theta_n^2}(\boldsymbol{\theta}) \end{bmatrix} \in \mathbb{R}^{n \times n}.$$

> **🔔 Theorem 11.1 (Slopes, curvatures, and partial derivatives)**
>
> Let $J : \mathbb{R}^n \to \mathbb{R}$ be a function of class $C^2$, $\boldsymbol{\theta} \in \mathbb{R}^n$ a point, and $\mathbf{v} \in \mathbb{R}^n$ a directional vector.
>
> 1. We have
>
> $$J'_{\mathbf{v}}(\boldsymbol{\theta}) = \mathbf{v}^\mathsf{T} \nabla J(\boldsymbol{\theta}).$$
>
> 2. We have
>
> $$J''_{\mathbf{v}}(\boldsymbol{\theta}) = \mathbf{v}^\mathsf{T} \nabla^2 J(\boldsymbol{\theta}) \mathbf{v}.$$

**Problem Prompt**

Do problem 6 on the worksheet.