

# **Math for Machine Learning**

**Week 2.2: Eigendecomposition and PSD Matrices**

**By: Samuel Deng**

# **Logistics & Announcements**

# Lesson Overview

**Linear dynamical systems example.** Motivation for eigendecomposition as a way to make repeated matrix multiplication easier.

**Eigendecomposition.** Definition of eigenvectors, eigenvalues.

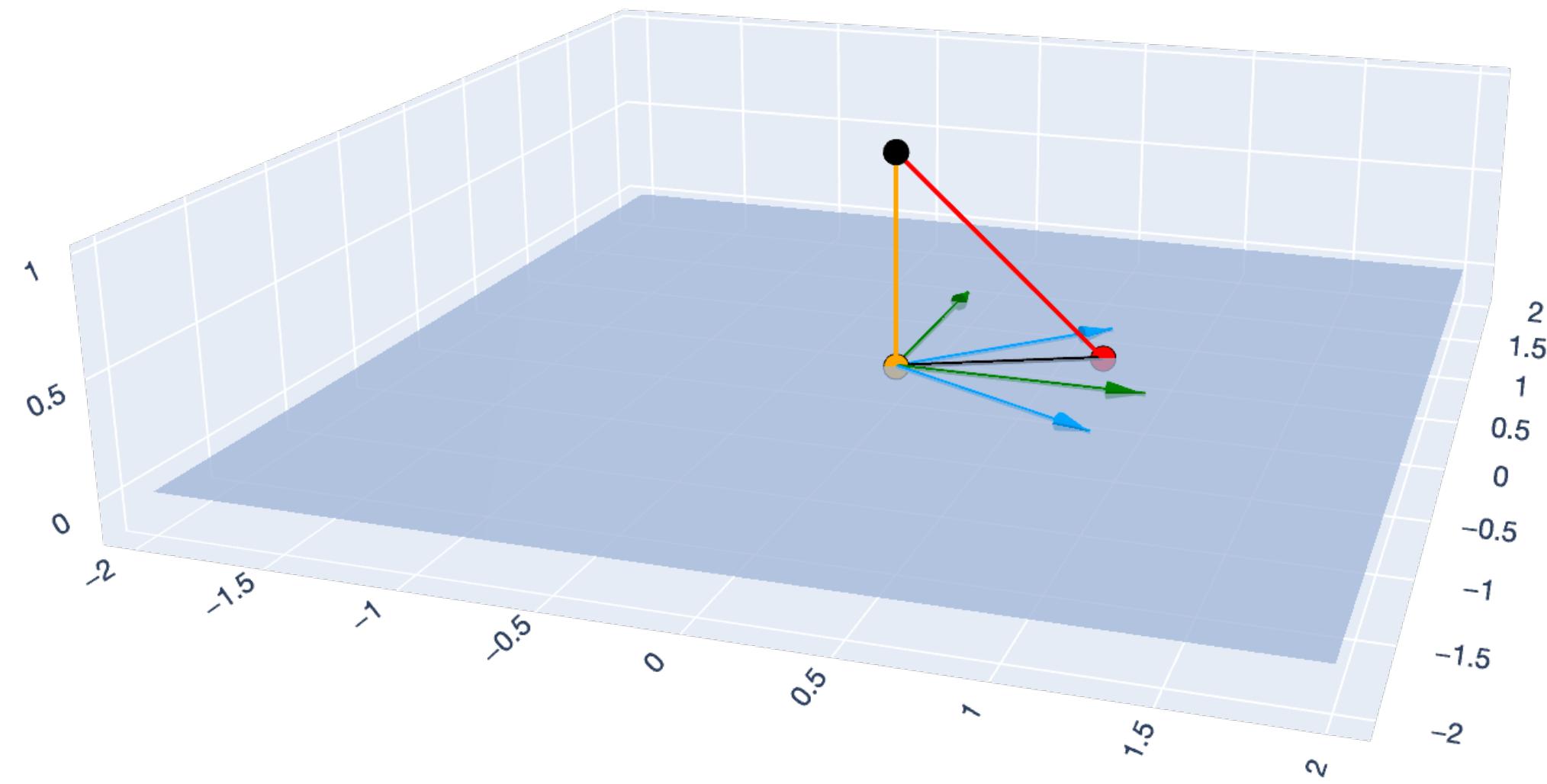
**Eigendecomposition and SVD.** The eigendecomposition drops out of the SVD.

**Spectral Theorem.** Symmetric matrices are always diagonalizable.

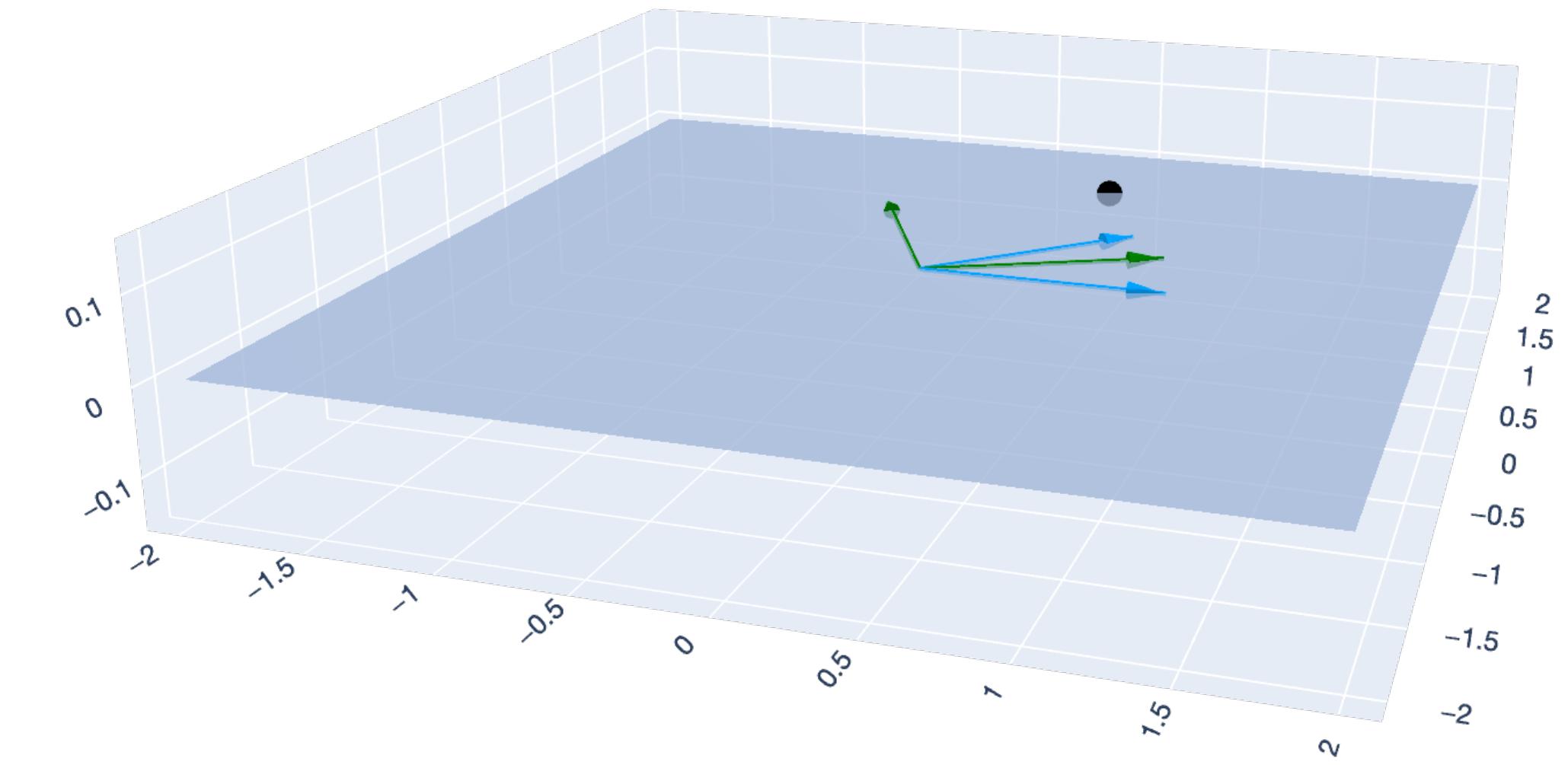
**Positive semidefinite matrices/positive definite matrices.** Definition and some visual examples through the corresponding quadratic forms.

# Lesson Overview

## Big Picture: Least Squares



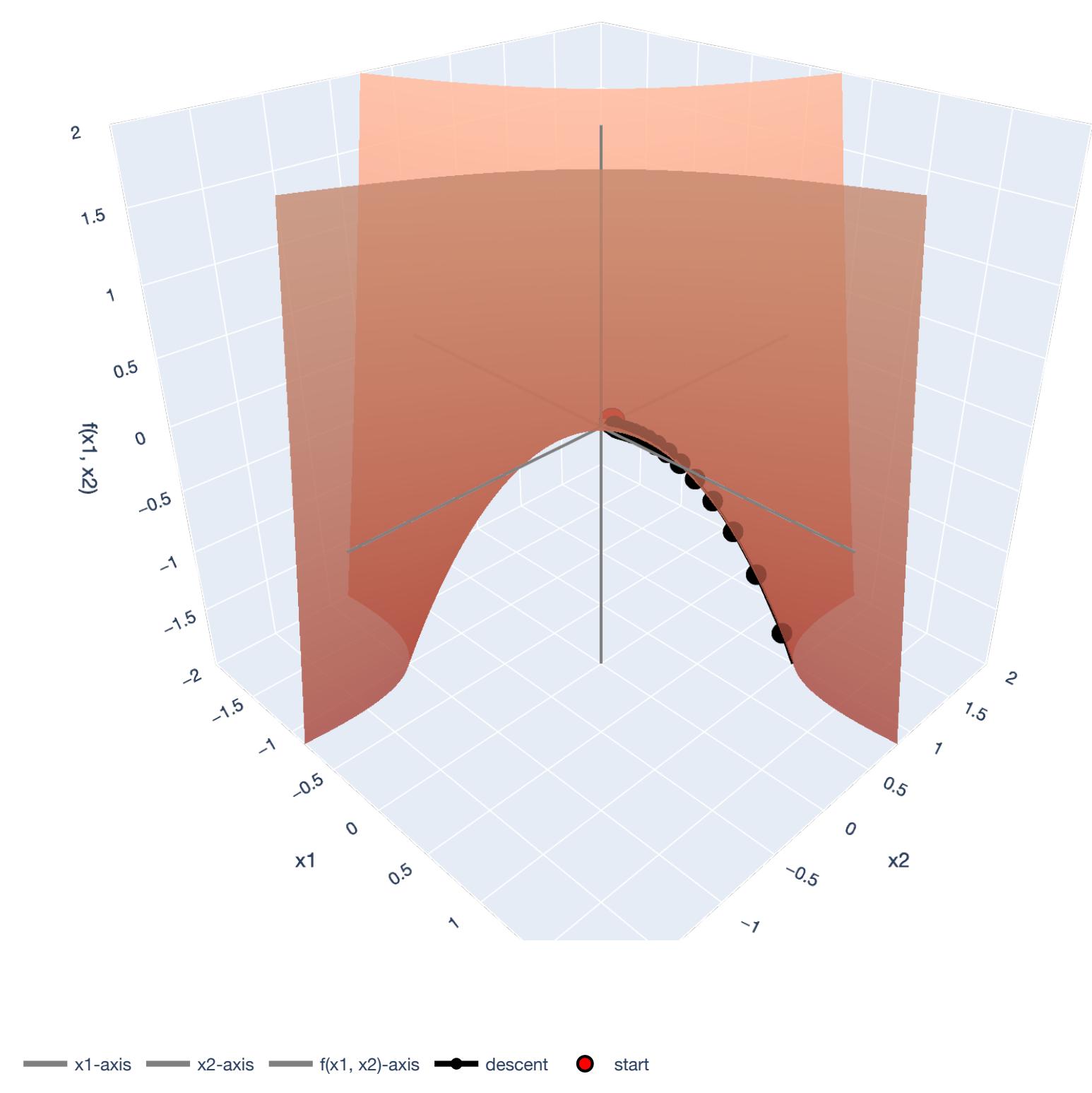
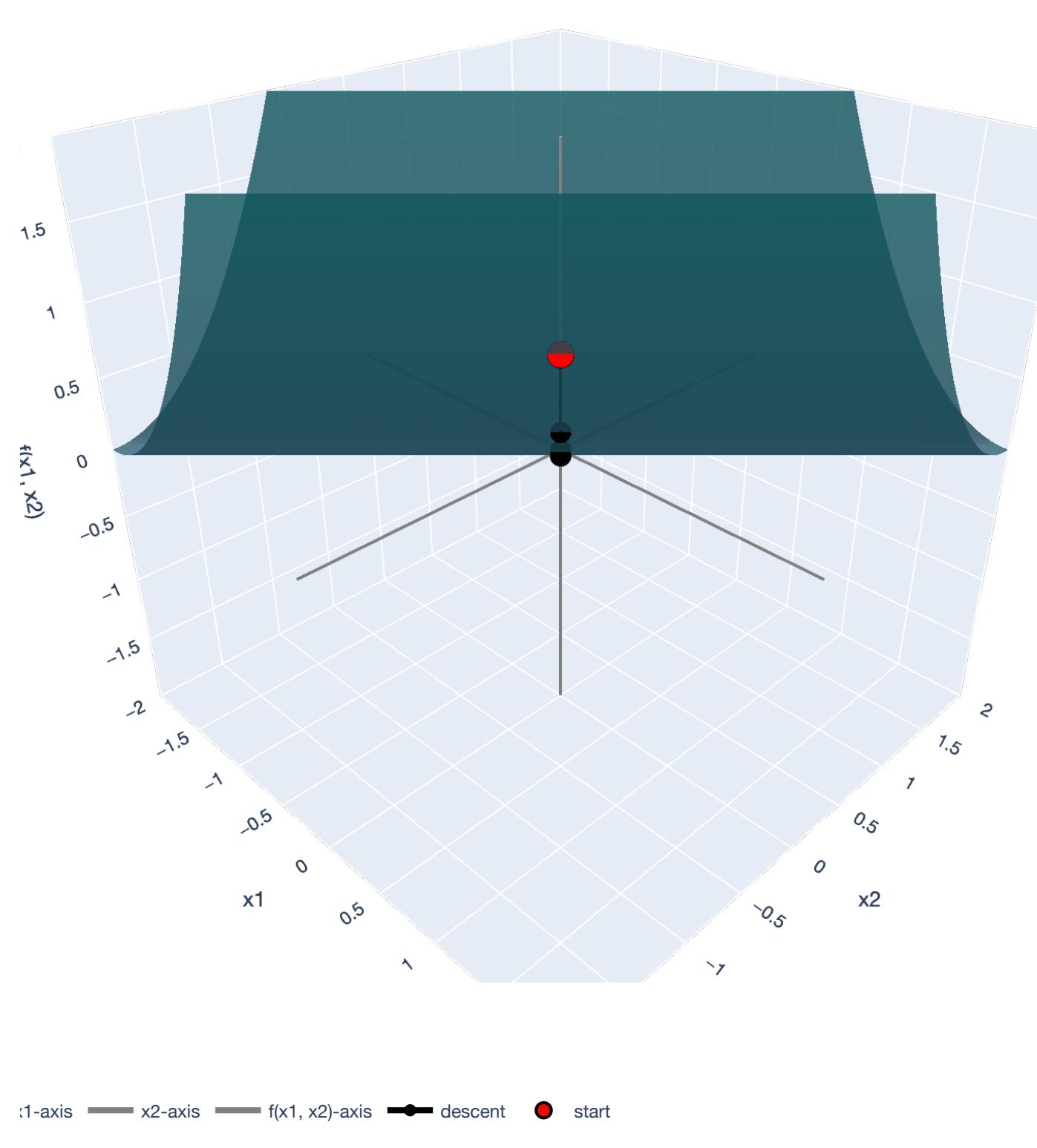
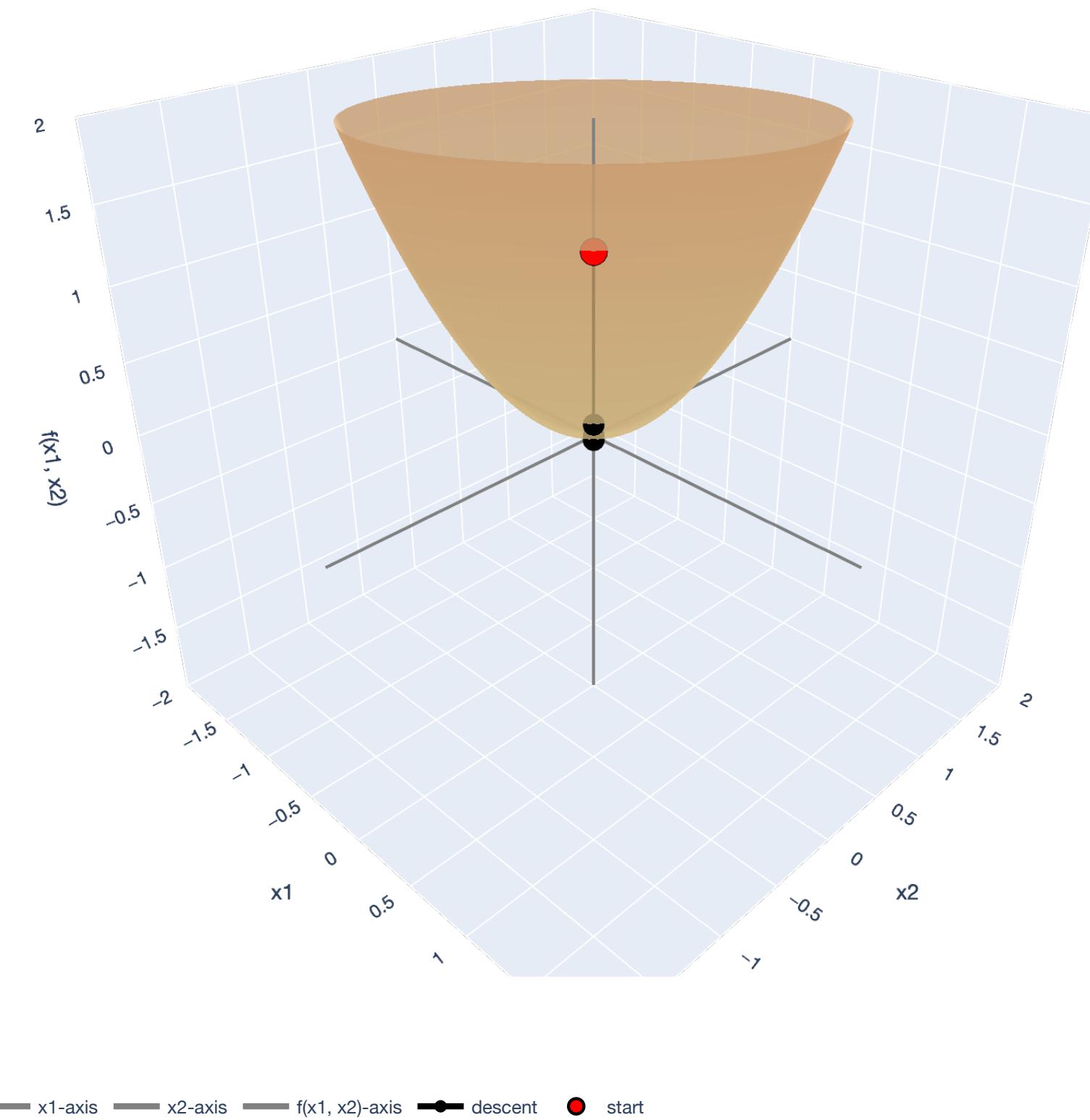
—  $x_1$  —  $x_2$  —  $u_1$  —  $u_2$  —  $y - \hat{y}$  —  $\hat{y} - y$  ●  $y$  ○  $\hat{y}$  ●  $\sim y$



—  $x_1$  —  $x_2$  —  $u_1$  —  $u_2$  ●  $y$

# Lesson Overview

## Big Picture: Gradient Descent



# Least Squares

## A Quick Review

# Regression Setup

**Observed:** Matrix of *training samples*  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and vector of *training labels*  $\mathbf{y} \in \mathbb{R}^d$ .

$$\mathbf{X} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{x}_1 & \dots & \mathbf{x}_d \\ \downarrow & & \downarrow \end{bmatrix} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix}.$$

**Unknown:** *Weight vector*  $\mathbf{w} \in \mathbb{R}^d$  with weights  $w_1, \dots, w_d$ .

**Goal:** For each  $i \in [n]$ , we predict:  $\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i = w_1 x_{i1} + \dots + w_d x_{id} \in \mathbb{R}$ .

Choose a weight vector that “fits the training data”:  $\mathbf{w} \in \mathbb{R}^d$  such that  $y_i \approx \hat{y}_i$  for  $i \in [n]$ , or:

$$\mathbf{X}\mathbf{w} = \hat{\mathbf{y}} \approx \mathbf{y}.$$

# Regression Setup

**Goal:** For each  $i \in [n]$ , we predict:  $\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i = w_1 x_{i1} + \dots + w_d x_{id} \in \mathbb{R}$ .

Choose a weight vector that “fits the training data”:  $\hat{\mathbf{w}} \in \mathbb{R}^d$  such that  $y_i \approx \hat{y}_i$  for  $i \in [n]$ , or:

$$\mathbf{X}\hat{\mathbf{w}} = \hat{\mathbf{y}} \approx \mathbf{y}.$$

To find  $\hat{\mathbf{w}}$ , we follow the *principle of least squares*.

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

# SVD and Pseudoinverse

## Review

Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be a matrix, and let  $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$  be its full SVD.

If  $n \geq d$ , the matrix  $(\Sigma^\top\Sigma)^{-1}\Sigma^\top \in \mathbb{R}^{d \times n}$  is the [\(Moore-Penrose\) pseudoinverse](#) of the matrix  $\Sigma$ , denoted  $\Sigma^+ := (\Sigma^\top\Sigma)^{-1}\Sigma^\top$ .

If  $d > n$ , the matrix  $\Sigma^+ := \Sigma^\top(\Sigma\Sigma^\top)^{-1}$  is the pseudoinverse.

More generally, the matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  with full SVD  $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$  has the [\(Moore-Penrose\) pseudoinverse](#):  $\mathbf{X}^+ := \mathbf{V}\Sigma^+\mathbf{U}^\top$ .

# Least Squares: SVD Perspective

## Unified Picture

We want to solve  $\mathbf{X}\mathbf{w} = \mathbf{y}$ .

If  $n = d$  and  $\text{rank}(\mathbf{X}) = d\dots$

We can solve exactly.

Choose

$$\hat{\mathbf{w}} = \mathbf{X}^{-1}\mathbf{y},$$

which is an exact solution.

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

Choose

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{X}^+ \mathbf{y},$$

the best approximate solution:

$$\|\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}\|^2 \leq \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

If  $n > d$  and  $\text{rank}(\mathbf{X}) = d\dots$

We approximate by least squares:

Choose

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{X}^+ \mathbf{y},$$

the minimum norm solution:

If  $n < d$  and  $\text{rank}(\mathbf{X}) = n\dots$

We can solve exactly, but there are infinitely many solutions.

Choose

$$\hat{\mathbf{w}} = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y} = \mathbf{X}^+ \mathbf{y},$$

the minimum norm solution:

$$\|\hat{\mathbf{w}}\|^2 \leq \|\mathbf{w}\|^2.$$

# Least Squares: SVD Perspective

## Unified Picture

We want to solve  $\mathbf{X}\mathbf{w} = \mathbf{y}$ . Use  $\hat{\mathbf{w}} = \mathbf{X}^+\mathbf{y}$ !

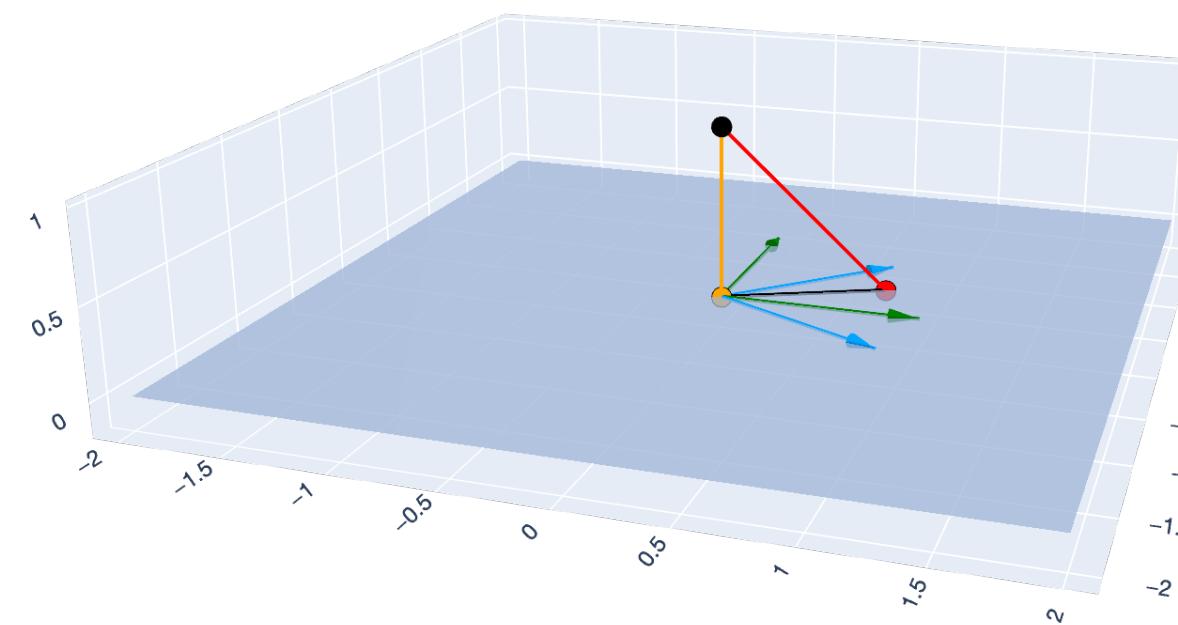
If  $n > d$  and  $\text{rank}(\mathbf{X}) = d$ ...

We approximate by least squares:

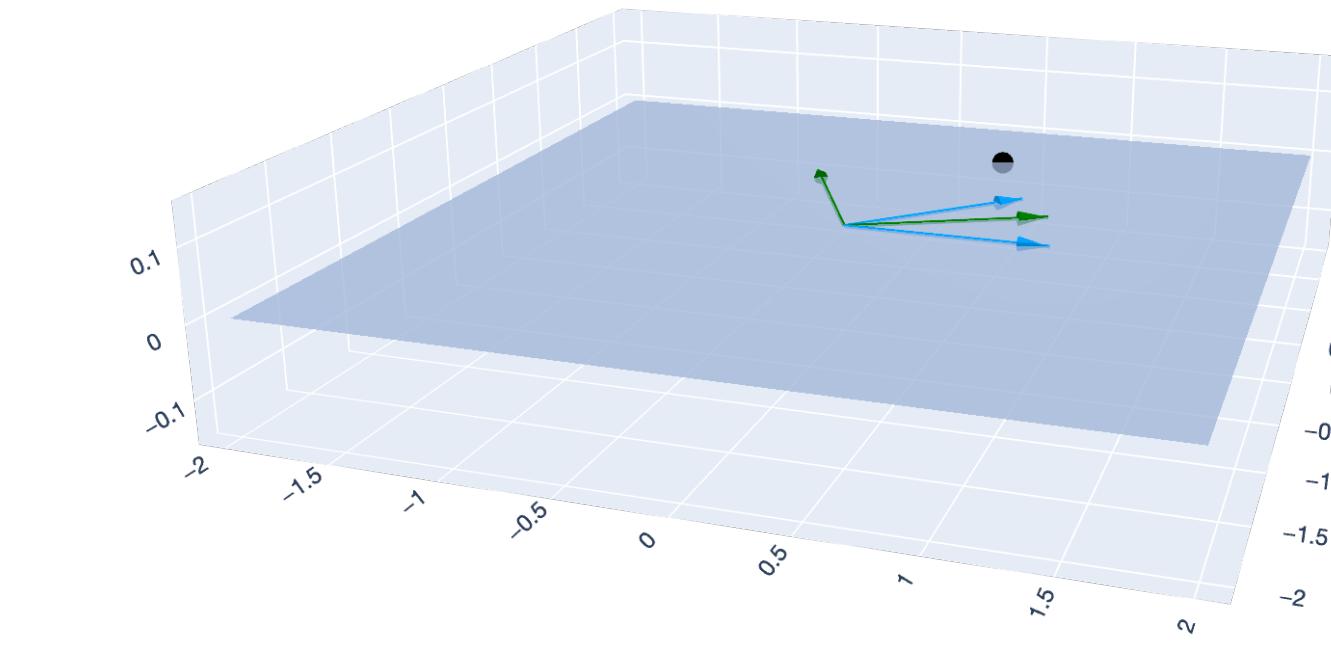
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$

If  $n < d$  and  $\text{rank}(\mathbf{X}) = n$ ...

We can solve exactly, but there are infinitely many solutions.



Legend:  $\mathbf{x}_1$   $\mathbf{x}_2$   $\mathbf{u}_1$   $\mathbf{u}_2$   $\mathbf{y} - \hat{\mathbf{y}}$   $\mathbf{u}_1 - \hat{\mathbf{y}}$   $\mathbf{u}_2 - \hat{\mathbf{y}}$   $\mathbf{y}$   $\hat{\mathbf{y}}$   $\mathbf{y} - \hat{\mathbf{y}}$



Legend:  $\mathbf{x}_1$   $\mathbf{x}_2$   $\mathbf{u}_1$   $\mathbf{u}_2$   $\mathbf{y}$

# Singular Value Decomposition (SVD)

## Matrix Decompositions

$$\underbrace{\mathbf{X}}_{n \times d} = \underbrace{\mathbf{U}}_{n \times n} \underbrace{\Sigma}_{n \times d} \underbrace{\mathbf{V}^\top}_{d \times d}.$$

$\mathbf{U}$  is orthogonal, i.e.  $\mathbf{U}^\top \mathbf{U} = \mathbf{U} \mathbf{U}^\top = \mathbf{I}$ .

$\mathbf{V}$  is orthogonal, i.e.  $\mathbf{V}^\top \mathbf{V} = \mathbf{V} \mathbf{V}^\top = \mathbf{I}$ .

$\Sigma \in \mathbb{R}^{n \times d}$  is a diagonal matrix with **singular values**  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d \geq 0$  on the diagonal.  $\text{rank}(\mathbf{X})$  is equal to the number of  $\sigma_i > 0$ .

*What other matrix  
decompositions are out there?*

# Eigendecomposition

## Motivation: Linear Dynamical System

# Population Change

## Example of a linear dynamical system

Consider the following example.

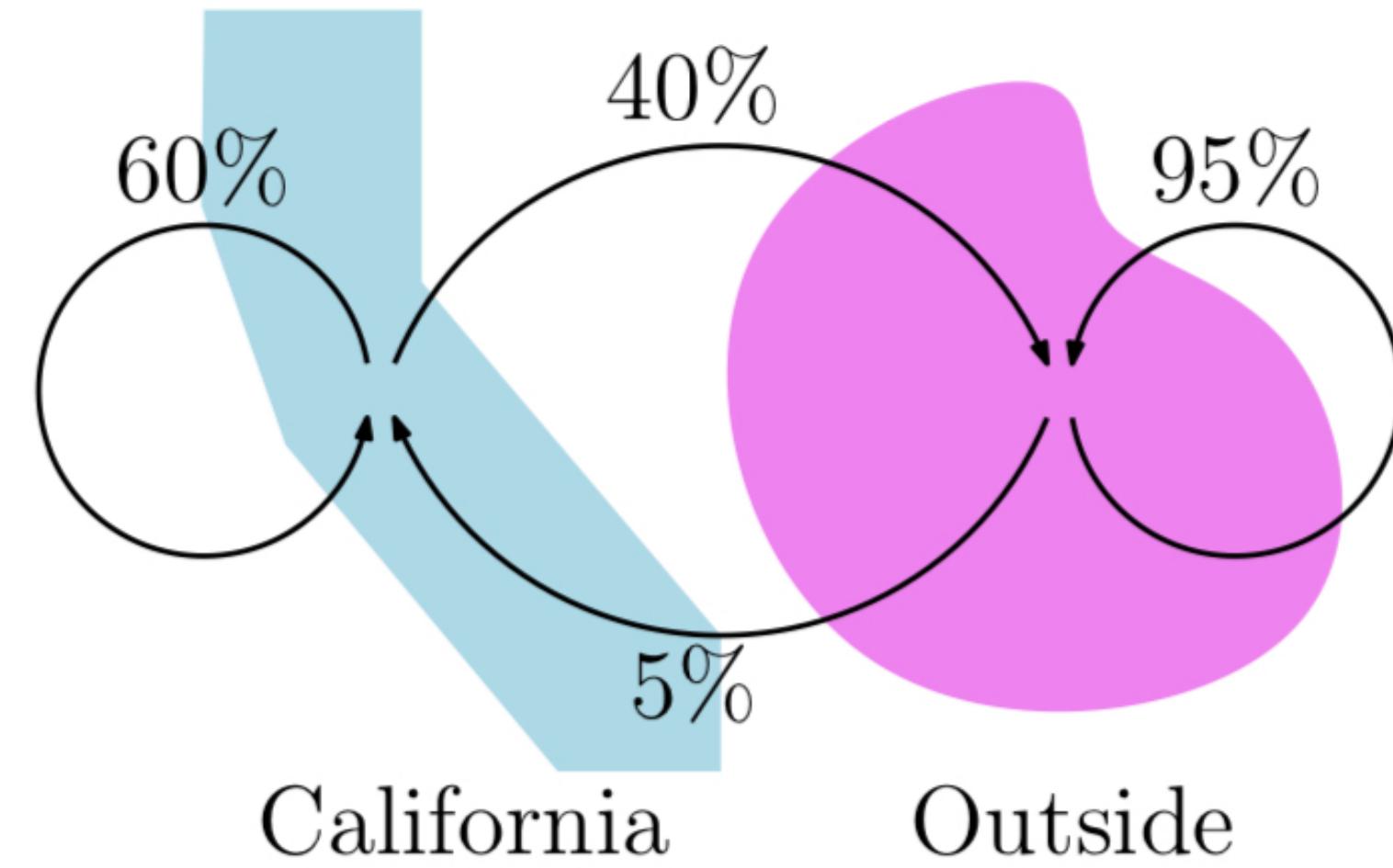
Suppose that

- of those who start a year in California, 60% stay in California and 40% move out of California by the end of the year.
- of those who start a year outside California, 95% stay out and 5% move to California by the end of the year.

If we know how many people are in California  $x_{in}$  and how many people are outside of California  $x_{out}$ , then we can find the number of people inside and outside of California at the end of the year:

$$\# \text{ inside} = 0.6x_{in} + 0.05x_{out}$$

$$\# \text{ outside} = 0.4x_{in} + 0.95x_{out}$$



Example and graphic from Daniel Hsu's course:  
*Computational Linear Algebra* (Fall 2022)

# Population Change

## Modeling with a transition matrix

Consider the following example.

Suppose that

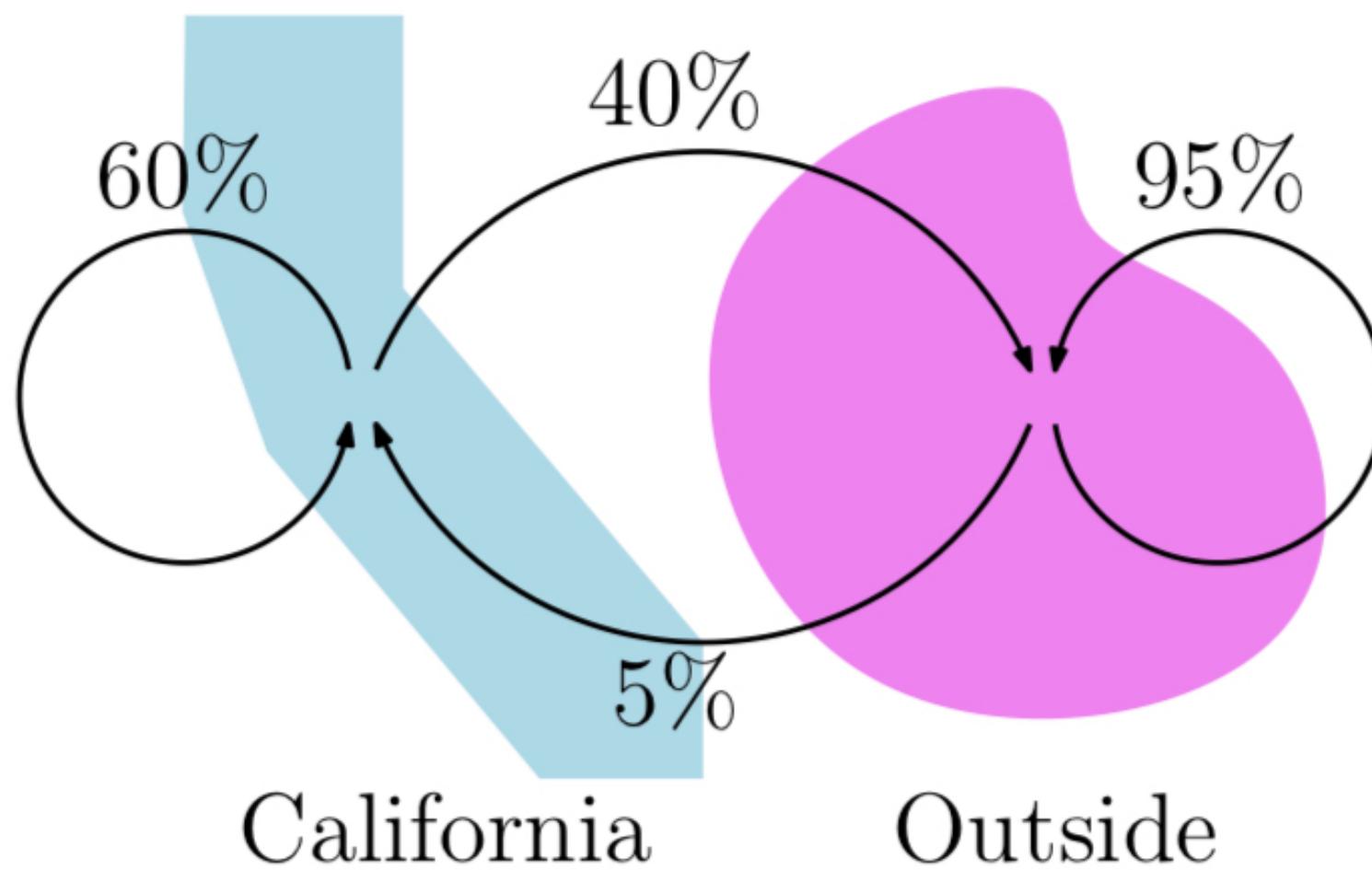
- of those who start a year in California, 60% stay in California and 40% move out of California by the end of the year.
- of those who start a year outside California, 95% stay out and 5% move to California by the end of the year.

We can model this with a *transition matrix*

$$\mathbf{A} = \begin{bmatrix} \text{in} \rightarrow \text{in} & \text{out} \rightarrow \text{in} \\ \text{in} \rightarrow \text{out} & \text{out} \rightarrow \text{out} \end{bmatrix} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix}$$

and a system of linear equations:

$$\mathbf{Ax} = \begin{bmatrix} \text{in} \rightarrow \text{in} & \text{out} \rightarrow \text{in} \\ \text{in} \rightarrow \text{out} & \text{out} \rightarrow \text{out} \end{bmatrix} \begin{bmatrix} x_{\text{in}} \\ x_{\text{out}} \end{bmatrix} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix} x_{\text{in}} \\ x_{\text{out}} \end{bmatrix}$$



Example and graphic from Daniel Hsu's course:  
*Computational Linear Algebra* (Fall 2022)

# Population Change

## Modeling with a transition matrix

Consider the transition matrix

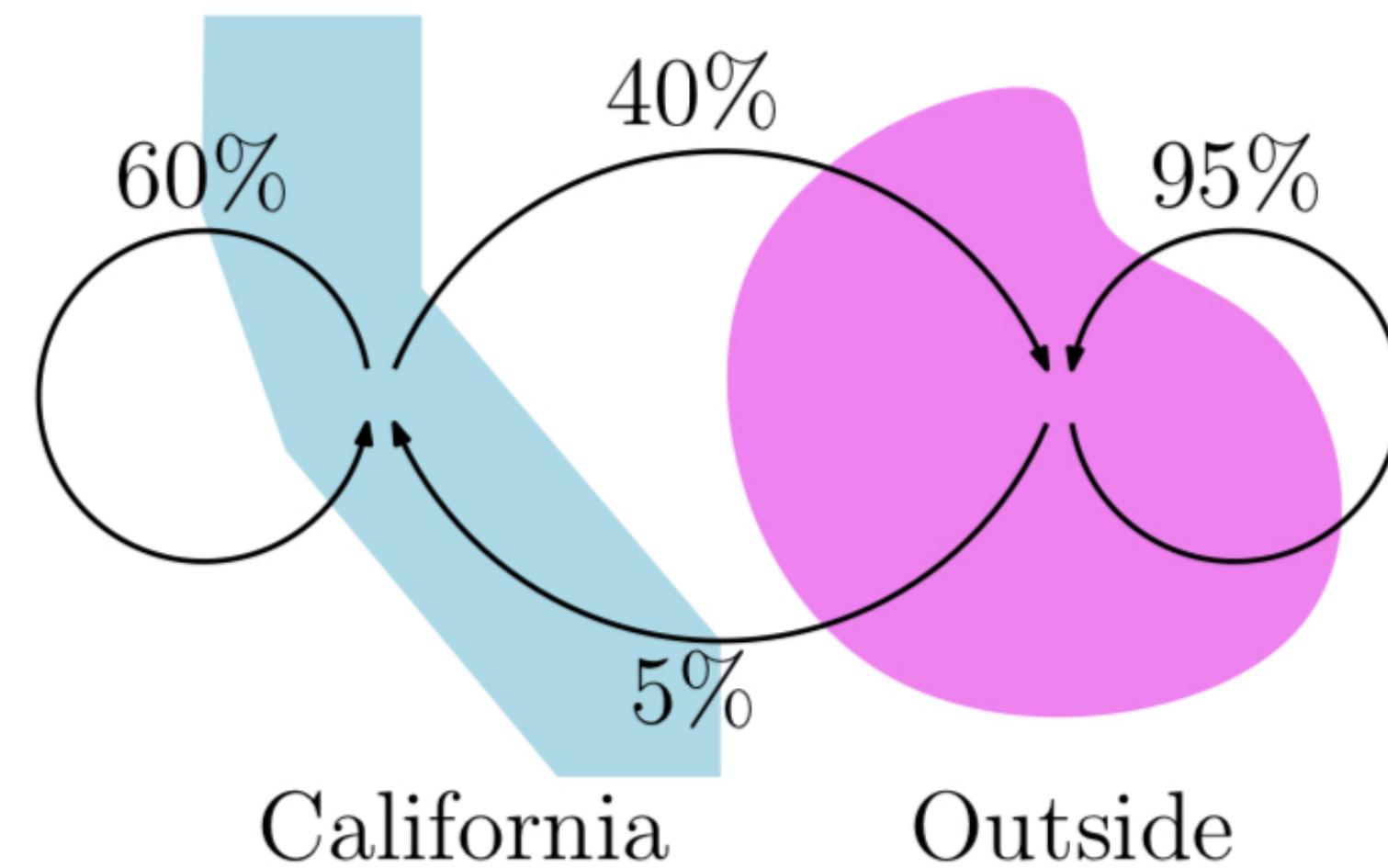
$$A = \begin{bmatrix} \text{in} \rightarrow \text{in} & \text{out} \rightarrow \text{in} \\ \text{in} \rightarrow \text{out} & \text{out} \rightarrow \text{out} \end{bmatrix} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix}$$

with a corresponding system of linear equations:

$$Ax = \begin{bmatrix} \text{in} \rightarrow \text{in} & \text{out} \rightarrow \text{in} \\ \text{in} \rightarrow \text{out} & \text{out} \rightarrow \text{out} \end{bmatrix} \begin{bmatrix} x_{\text{in}} \\ x_{\text{out}} \end{bmatrix} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix} x_{\text{in}} \\ x_{\text{out}} \end{bmatrix}.$$

The vector  $Ax \in \mathbb{R}^2$  gives the number of people inside and outside of California after a year has passed, from the initial populations in  $x \in \mathbb{R}^2$ .

*How to find the number of people inside/outside of California after  $t$  years have passed?*



Example and graphic from Daniel Hsu's course:  
*Computational Linear Algebra* (Fall 2022)

# Population Change

## Modeling with a transition matrix

Consider the transition matrix

$$\mathbf{A} = \begin{bmatrix} \text{in} \rightarrow \text{in} & \text{out} \rightarrow \text{in} \\ \text{in} \rightarrow \text{out} & \text{out} \rightarrow \text{out} \end{bmatrix} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix}$$

with a corresponding system of linear equations:

$$\mathbf{Ax} = \begin{bmatrix} \text{in} \rightarrow \text{in} & \text{out} \rightarrow \text{in} \\ \text{in} \rightarrow \text{out} & \text{out} \rightarrow \text{out} \end{bmatrix} \begin{bmatrix} x_{\text{in}} \\ x_{\text{out}} \end{bmatrix} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix} x_{\text{in}} \\ x_{\text{out}} \end{bmatrix}.$$

The vector  $\mathbf{Ax}^{(0)} \in \mathbb{R}^2$  gives the number of people inside and outside of California after a year has passed, from the initial populations in  $\mathbf{x}^{(0)} \in \mathbb{R}^2$ .

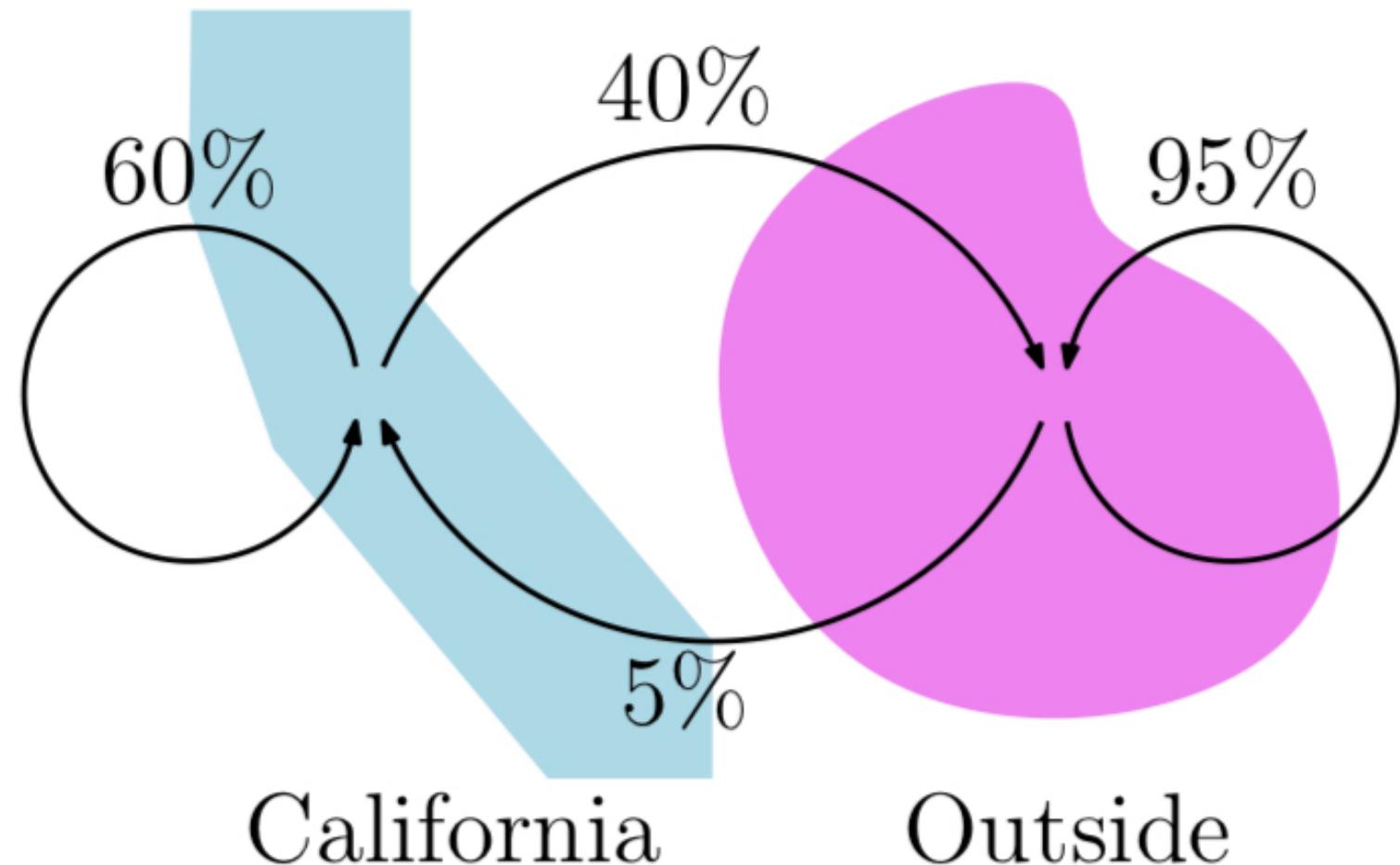
*How to find the number of people inside/outside of California after  $t$  years have passed?*

$$\mathbf{x}^{(1)} = \mathbf{Ax}^{(0)}$$

$$\mathbf{x}^{(2)} = \mathbf{Ax}^{(1)} = \mathbf{A}\mathbf{Ax}^{(0)} = \mathbf{A}^2\mathbf{x}^{(0)}$$

⋮

$$\mathbf{x}^{(t)} = \underbrace{\mathbf{A}\mathbf{A}\dots\mathbf{A}}_{t \text{ products}} \mathbf{x}^{(0)} = \mathbf{A}^t \mathbf{x}^{(0)}$$



Example and graphic from Daniel Hsu's course:  
*Computational Linear Algebra* (Fall 2022)

# Population Change

## Modeling with a transition matrix

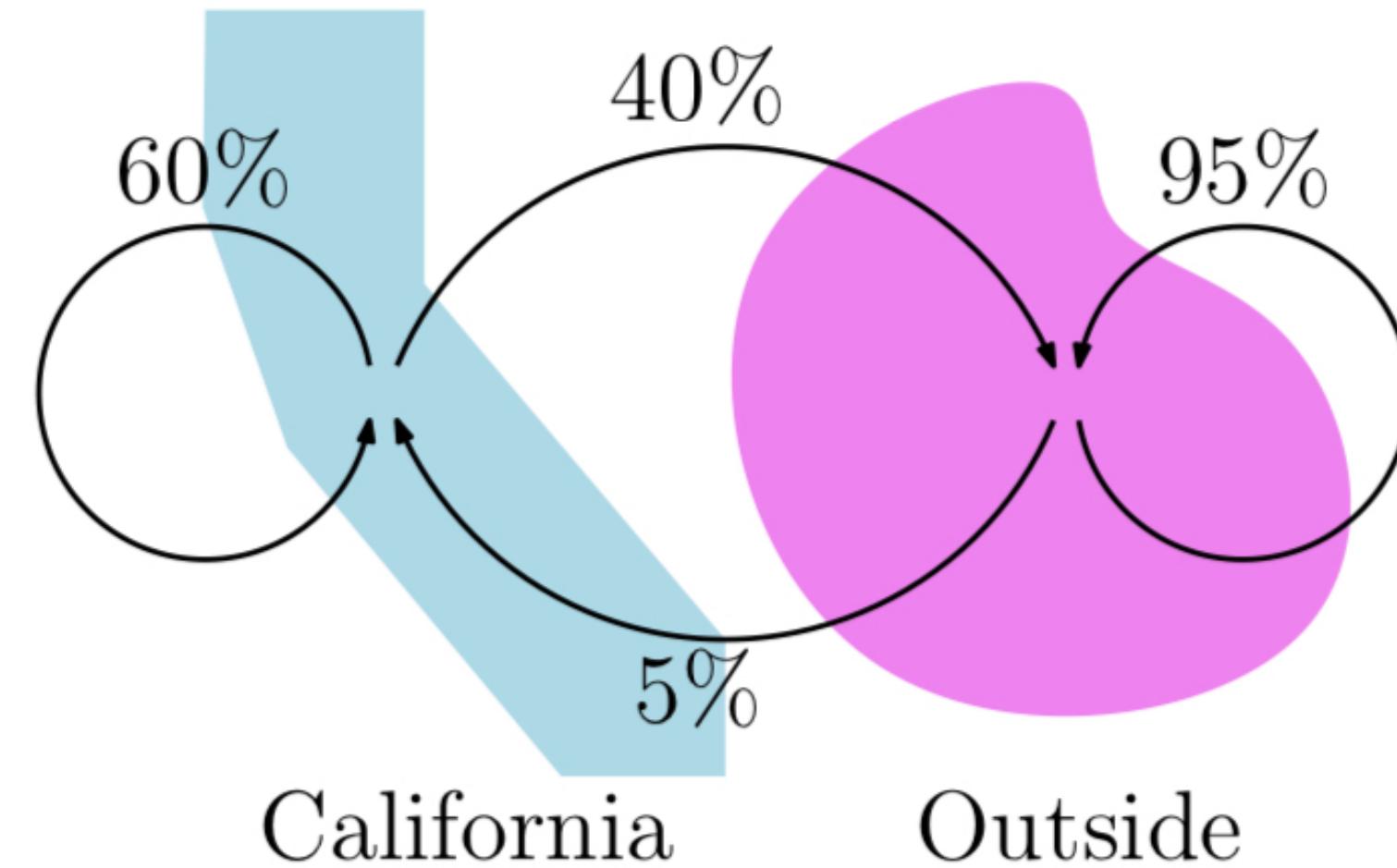
$$\mathbf{Ax} = \begin{bmatrix} \text{in} \rightarrow \text{in} & \text{out} \rightarrow \text{in} \\ \text{in} \rightarrow \text{out} & \text{out} \rightarrow \text{out} \end{bmatrix} \begin{bmatrix} x_{in} \\ x_{out} \end{bmatrix} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix} x_{in} \\ x_{out} \end{bmatrix}$$

Concretely, suppose there are 300 million outside of California and 40 million inside of California at the start of a year. Then,

$$\mathbf{x}^{(0)} = \begin{bmatrix} 40 \\ 300 \end{bmatrix}$$

*What are the populations inside and outside of CA after  $t$  years?*

$$\mathbf{x}^{(t)} = \mathbf{A}^t \mathbf{x}^{(0)} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix}^t \begin{bmatrix} 40 \\ 300 \end{bmatrix}$$



Example and graphic from Daniel Hsu's course:  
*Computational Linear Algebra* (Fall 2022)

# Population Change

Annoying computation 😤

*What are the populations inside and outside of CA after  $t$  years?*

$$\mathbf{x}^{(t)} = \mathbf{A}^t \mathbf{x}^{(0)} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix}^t \begin{bmatrix} 40 \\ 300 \end{bmatrix}$$

Try calculating this...

$$\begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \cdots \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix} 40 \\ 300 \end{bmatrix}$$

# Population Change

Easy computation 😊

Assume I gave you a couple of vectors,  $\mathbf{u} = (1,8)$  and  $\mathbf{v} = (-1,1)$ . These two vectors have the properties:

$$\mathbf{A}\mathbf{u} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix} 1 \\ 8 \end{bmatrix} = \begin{bmatrix} 1 \\ 8 \end{bmatrix}$$

$$\mathbf{A}\mathbf{v} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \frac{11}{20} \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

# Population Change

Easy computation 😊

Assume I gave you a couple of vectors,  $\mathbf{u} = (1, 8)$  and  $\mathbf{v} = (-1, 1)$ . These two vectors have the properties:

$$\mathbf{A}\mathbf{u} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix} 1 \\ 8 \end{bmatrix} = \begin{bmatrix} 1 \\ 8 \end{bmatrix}$$

$$\mathbf{A}\mathbf{v} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \frac{11}{20} \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

Now, the repeated multiplication looks like:

$$\mathbf{A}^t\mathbf{u} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix}^t \begin{bmatrix} 1 \\ 8 \end{bmatrix} = (1)^t \begin{bmatrix} 1 \\ 8 \end{bmatrix} = \begin{bmatrix} 1 \\ 8 \end{bmatrix}$$

$$\mathbf{A}^t\mathbf{v} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix}^t \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \left(\frac{11}{20}\right)^t \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

# Population Change

## Using $\mathbf{u}$ and $\mathbf{v}$ for initial population

Assume I gave you a couple of vectors,  $\mathbf{u} = (1, 8)$  and  $\mathbf{v} = (-1, 1)$ . These two vectors have the properties:

$$\mathbf{A}\mathbf{u} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix} 1 \\ 8 \end{bmatrix} = \begin{bmatrix} 1 \\ 8 \end{bmatrix}$$

$$\mathbf{A}\mathbf{v} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \frac{11}{20} \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

Now, the repeated multiplication looks like:

$$\mathbf{A}^t \mathbf{u} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix}^t \begin{bmatrix} 1 \\ 8 \end{bmatrix} = (1)^t \begin{bmatrix} 1 \\ 8 \end{bmatrix} = \begin{bmatrix} 1 \\ 8 \end{bmatrix} \implies \mathbf{A}^t \mathbf{u} = \mathbf{u}$$

$$\mathbf{A}^t \mathbf{v} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix}^t \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \left(\frac{11}{20}\right)^t \begin{bmatrix} -1 \\ 1 \end{bmatrix} \implies \mathbf{A}^t \mathbf{v} = \left(\frac{11}{20}\right)^t \mathbf{v}$$

# Population Change

## Using $\mathbf{u}$ and $\mathbf{v}$ for initial population

For  $\mathbf{u} = (1, 8)$  and  $\mathbf{v} = (-1, 1)$ ,

$$\mathbf{A}^t \mathbf{u} = \mathbf{u}$$

$$\mathbf{A}^t \mathbf{v} = \left( \frac{11}{20} \right)^t \mathbf{v}$$

Notice that  $\mathbf{u}, \mathbf{v}$  are a basis for  $\mathbb{R}^2$ . Then, if we rewrite  $\mathbf{x}^{(0)}$  as a linear combination of  $\mathbf{u}$  and  $\mathbf{v}$ , i.e.

$$\mathbf{x}^{(0)} = a\mathbf{u} + b\mathbf{v},$$

we can obtain  $\mathbf{x}^{(t)}$  with the following computation:

$$\mathbf{x}^{(t)} = \mathbf{A}^t \mathbf{x}^{(0)} = \mathbf{A}^t(a\mathbf{u} + b\mathbf{v}) = a\mathbf{A}^t \mathbf{u} + b\mathbf{A}^t \mathbf{v} = a\mathbf{u} + b(11/20)^t \mathbf{v}.$$

# Population Change

## Using $\mathbf{u}$ and $\mathbf{v}$ for initial population

For  $\mathbf{u} = (1, 8)$  and  $\mathbf{v} = (-1, 1)$ ,

$$\mathbf{A}^t \mathbf{u} = \mathbf{u}$$

$$\mathbf{A}^t \mathbf{v} = \left( \frac{11}{20} \right)^t \mathbf{v}$$

Notice that  $\mathbf{u}, \mathbf{v}$  are a basis for  $\mathbb{R}^2$ . Then, if we rewrite  $\mathbf{x}^{(0)}$  as a linear combination of  $\mathbf{u}$  and  $\mathbf{v}$ , i.e.

$$\mathbf{x}^{(0)} = a\mathbf{u} + b\mathbf{v},$$

we can obtain  $\mathbf{x}^{(t)}$  with the following computation:

$$\mathbf{x}^{(t)} = \mathbf{A}^t \mathbf{x}^{(0)} = \mathbf{A}^t(a\mathbf{u} + b\mathbf{v}) = a\mathbf{A}^t \mathbf{u} + b\mathbf{A}^t \mathbf{v} = a\mathbf{u} + b(11/20)^t \mathbf{v}.$$

In matrix form:

$$\mathbf{x}^{(t)} = \begin{bmatrix} \mathbf{u} & \mathbf{v} \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & (11/20)^t \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$

# Population Change

## Using $\mathbf{u}$ and $\mathbf{v}$ for initial population

For  $\mathbf{u} = (1, 8)$  and  $\mathbf{v} = (-1, 1)$ ,

$$\mathbf{x}^{(t)} = \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{u} & \mathbf{v} \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & (11/20)^t \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$

where

$$\mathbf{x}^{(0)} = a\mathbf{u} + b\mathbf{v}.$$

Writing  $\mathbf{x}^{(0)}$  in matrix form as well, we have:

$$\mathbf{x}^{(0)} = \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{u} & \mathbf{v} \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$

# Population Change

## Using $\mathbf{u}$ and $\mathbf{v}$ for initial population

For  $\mathbf{u} = (1, 8)$  and  $\mathbf{v} = (-1, 1)$ ,

$$\mathbf{x}^{(t)} = \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{u} & \mathbf{v} \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & (11/20)^t \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$

where

$$\mathbf{x}^{(0)} = a\mathbf{u} + b\mathbf{v}.$$

Writing  $\mathbf{x}^{(0)}$  in matrix form as well, we have:

$$\mathbf{x}^{(0)} = \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{u} & \mathbf{v} \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \mathbf{V} \begin{bmatrix} a \\ b \end{bmatrix}.$$

Because  $\mathbf{u}$  and  $\mathbf{v}$  are linearly independent,  $\mathbf{V} \in \mathbb{R}^{2 \times 2}$  has  $\text{rank}(\mathbf{V}) = 2$ , so we can invert:

$$\begin{bmatrix} a \\ b \end{bmatrix} = \mathbf{V}^{-1} \mathbf{x}^{(0)}.$$

# Population Change

## Using $\mathbf{u}$ and $\mathbf{v}$ for initial population

For  $\mathbf{u} = (1, 8)$  and  $\mathbf{v} = (-1, 1)$ ,

$$\mathbf{x}^{(t)} = \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{u} & \mathbf{v} \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & (11/20)^t \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$

where

$$\mathbf{x}^{(0)} = a\mathbf{u} + b\mathbf{v}.$$

Writing  $\mathbf{x}^{(0)}$  in matrix form as well, we have:

$$\mathbf{x}^{(0)} = \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{u} & \mathbf{v} \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \mathbf{V} \begin{bmatrix} a \\ b \end{bmatrix}.$$

Because  $\mathbf{u}$  and  $\mathbf{v}$  are linearly independent,  $\mathbf{V} \in \mathbb{R}^{2 \times 2}$  has  $\text{rank}(\mathbf{V}) = 2$ , so we can invert:

$$\begin{bmatrix} a \\ b \end{bmatrix} = \mathbf{V}^{-1} \mathbf{x}^{(0)}.$$

Therefore,

$$\mathbf{x}^{(t)} = \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{u} & \mathbf{v} \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & (11/20)^t \end{bmatrix} \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{u} & \mathbf{v} \\ \downarrow & \downarrow \end{bmatrix}^{-1} \mathbf{x}^{(0)} = \mathbf{V} \begin{bmatrix} 1 & 0 \\ 0 & (11/20)^t \end{bmatrix} \mathbf{V}^{-1} \mathbf{x}^{(0)}$$

# Population Change

## Using $\mathbf{u}$ and $\mathbf{v}$ for initial population

For  $\mathbf{u} = (1,8)$  and  $\mathbf{v} = (-1,1)$ ,

$$\mathbf{x}^{(t)} = \mathbf{V} \begin{bmatrix} 1 & 0 \\ 0 & (11/20)^t \end{bmatrix} \mathbf{V}^{-1} \mathbf{x}^{(0)}$$

where

$$\mathbf{V} = \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{u} & \mathbf{v} \\ \downarrow & \downarrow \end{bmatrix}.$$

# Population Change

## Comparison of hard and easy computation

Hard computation:

$$\mathbf{x}^{(t)} = \mathbf{A}^t \mathbf{x}^{(0)}$$

For initial populations  $\mathbf{x}^{(0)} = (40, 300)$ ,  
the population after  $t$  years is:

$$\mathbf{x}^{(t)} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix}^t \begin{bmatrix} 40 \\ 300 \end{bmatrix}.$$



Easy computation:

$$\mathbf{x}^{(t)} = \mathbf{V} \begin{bmatrix} 1 & 0 \\ 0 & (11/20)^t \end{bmatrix} \mathbf{V}^{-1} \mathbf{x}^{(0)}$$

For initial populations  $\mathbf{x}^{(0)} = (40, 300)$ , the  
population after  $t$  years is:

$$\mathbf{x}^{(t)} = \begin{bmatrix} 1 & -1 \\ 8 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & (11/20)^t \end{bmatrix} \begin{bmatrix} 1/9 & 1/9 \\ -8/9 & 1/9 \end{bmatrix} \begin{bmatrix} 40 \\ 300 \end{bmatrix}.$$



# Diagonal Matrices

## Why we like diagonal matrices

Multiplying diagonal matrices with themselves many times is easy:

$$\begin{bmatrix} 1 & 0 \\ 0 & (11/20)^t \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & (11/20) \end{bmatrix}^t.$$

# Diagonal Matrices

## Why we like diagonal matrices

Multiplying diagonal matrices with themselves many times is easy:

$$\begin{bmatrix} 1 & 0 \\ 0 & (11/20)^t \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & (11/20) \end{bmatrix}^t.$$

But this matrix depended on a basis of vectors that we got out of nowhere:

$$\mathbf{u} = (1, 8) \text{ and } \mathbf{v} = (-1, 1).$$

*In what cases (and how) can we obtain such nice bases?*

# Eigendecomposition

## Intuition and Definition

# Eigenvectors and eigenvalues

## Intuition

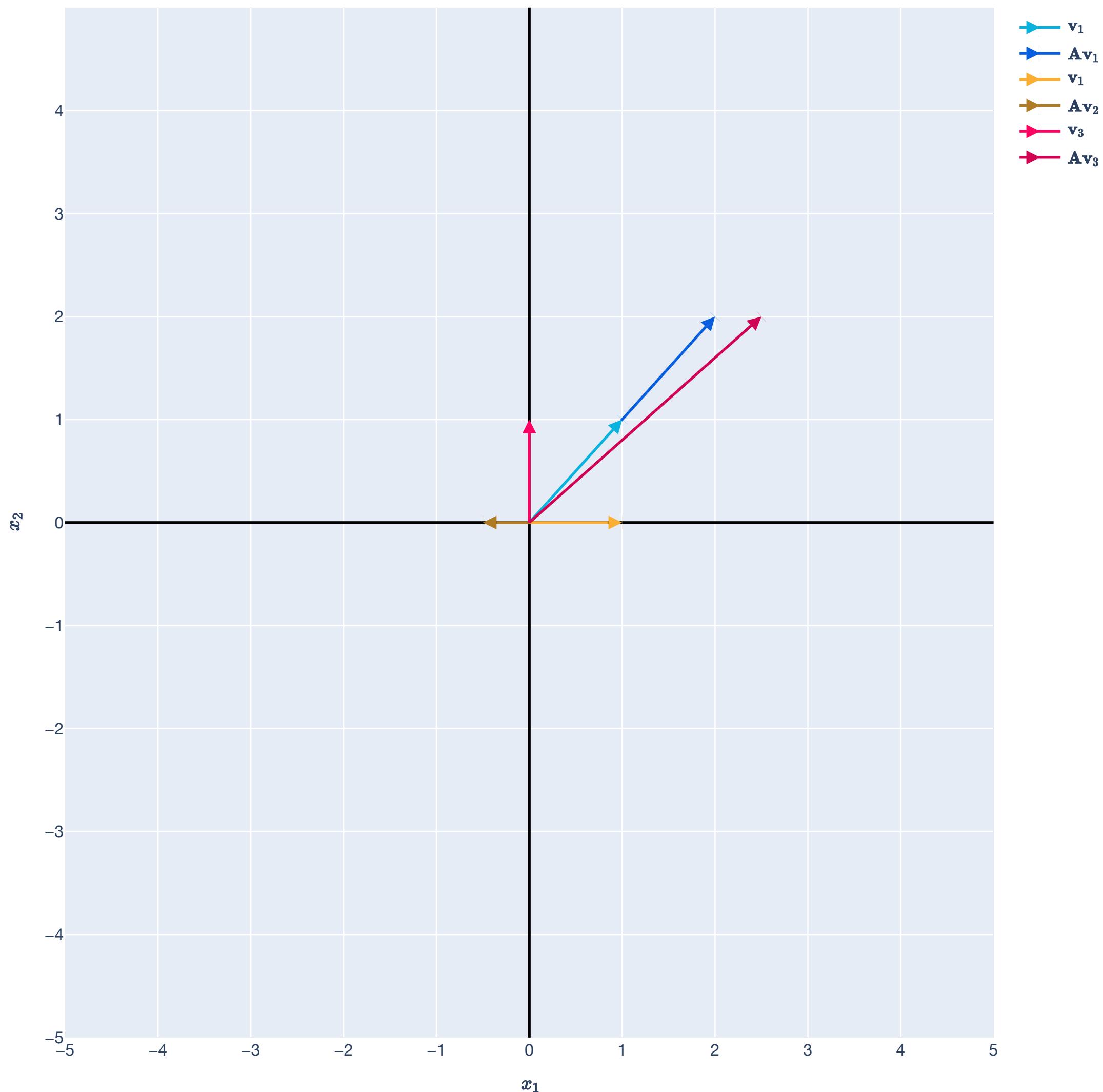
Let  $\mathbf{A} \in \mathbb{R}^{d \times d}$  be a *square* matrix.

This represents a linear transformation from  $\mathbb{R}^d$  to  $\mathbb{R}^d$ .

**Eigenvectors** are the vectors in  $\mathbb{R}^d$  that just get scaled by  $\mathbf{A}$ .

**Eigenvalues** are how much each eigenvector gets scaled.

Eigenvectors/eigenvalues are properties of square matrices!



# Eigenvectors and eigenvalues

## Definition

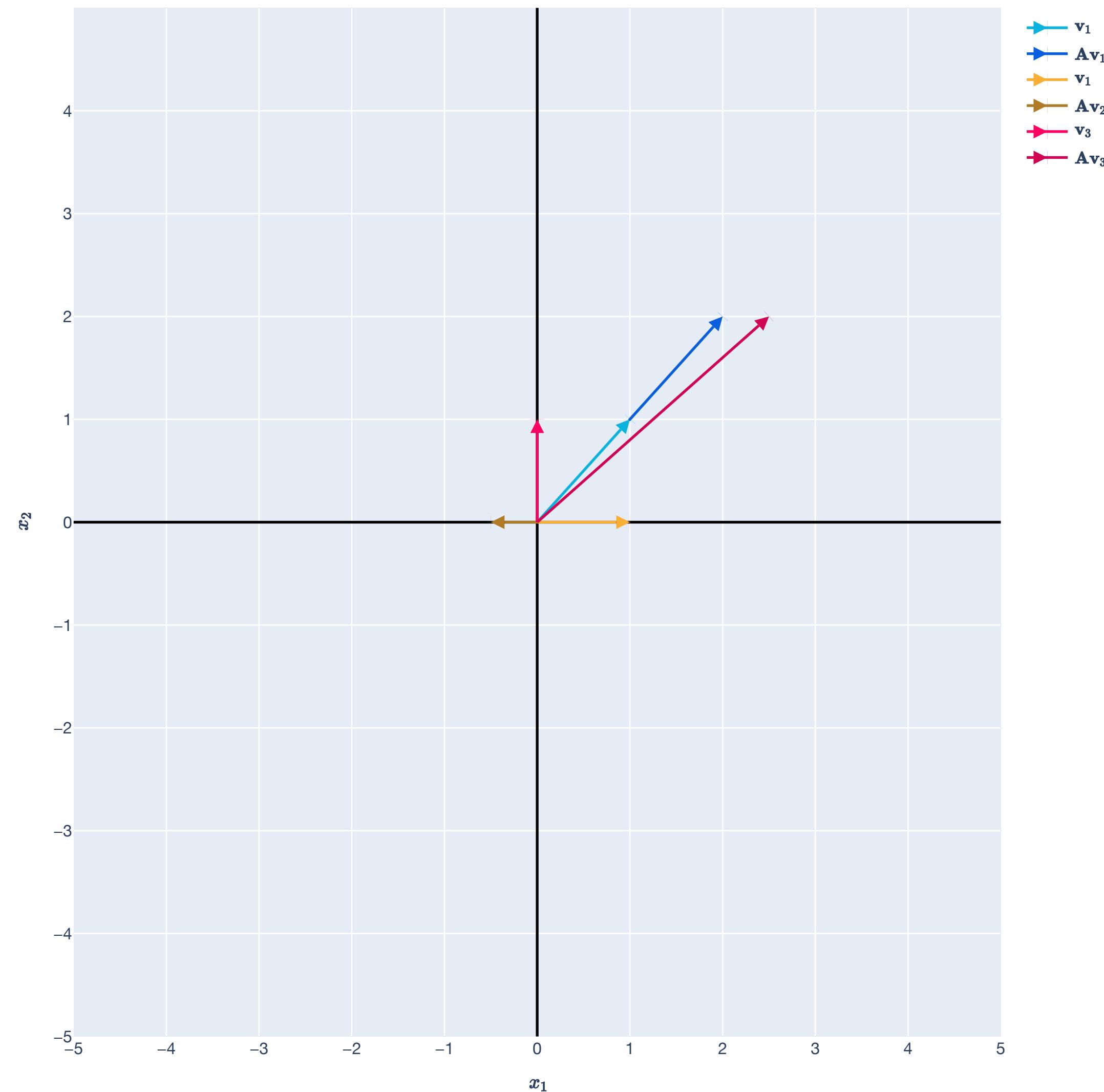
Let  $A \in \mathbb{R}^{d \times d}$  be a *square* matrix.

A nonzero vector  $v \in \mathbb{R}^d$  is an eigenvector if there exists a scalar  $\lambda \in \mathbb{R}$  such that

$$Av = \lambda v.$$

The scalar  $\lambda$  is the eigenvalue associated with the eigenvector  $v$ .

Eigenvectors/eigenvalues are properties of square matrices!



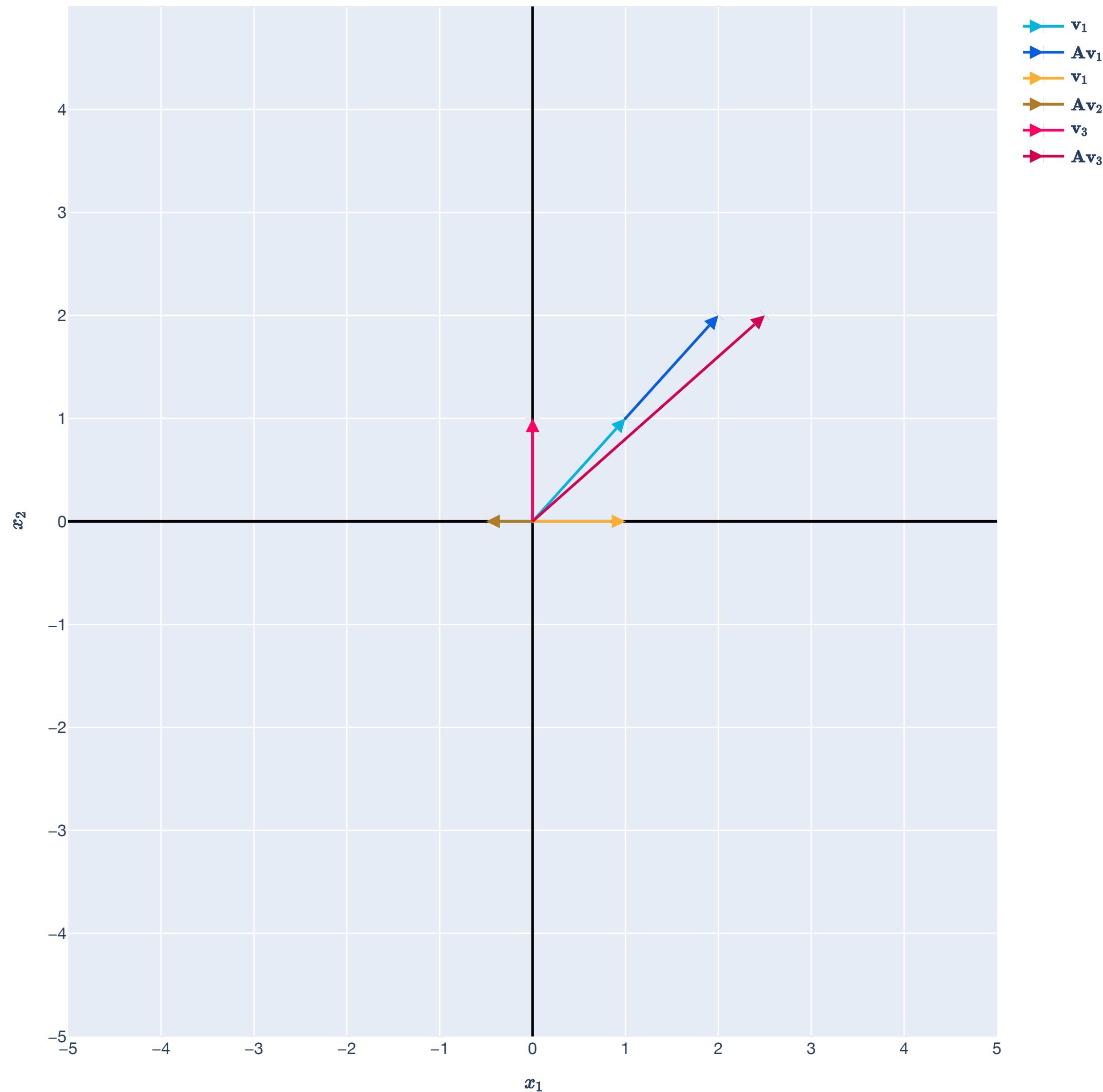
# Eigenvectors and eigenvalues

## Example

Consider the matrix  $A \in \mathbb{R}^{2 \times 2}$  given by

$$A = \begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix}.$$

What happens to the vector  $v_1 = (1,1)$ ?



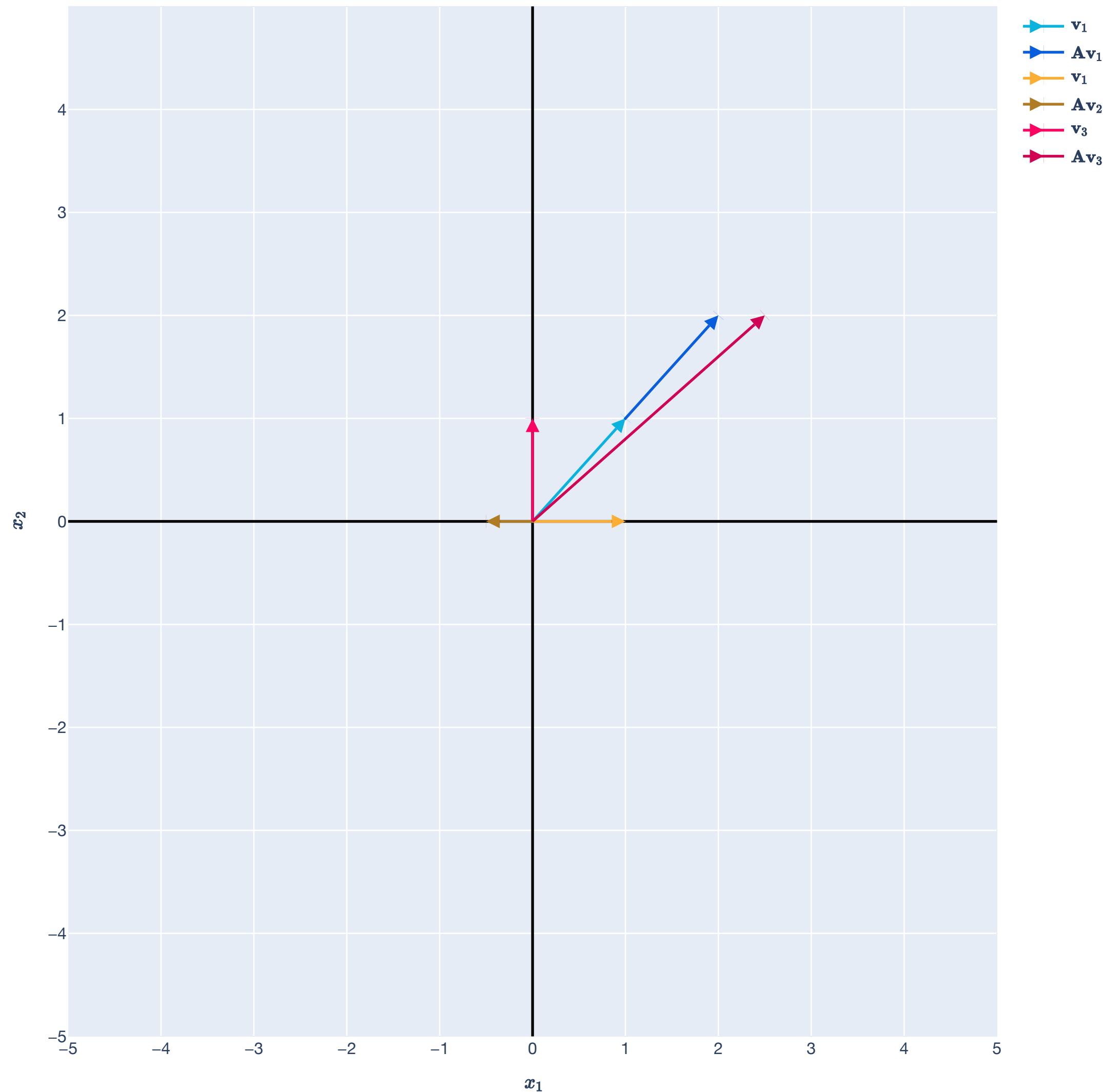
# Eigenvectors and eigenvalues

## Example

Consider the matrix  $A \in \mathbb{R}^{2 \times 2}$  given by

$$A = \begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix}.$$

What happens to the vector  $v_2 = (1,0)$ ?



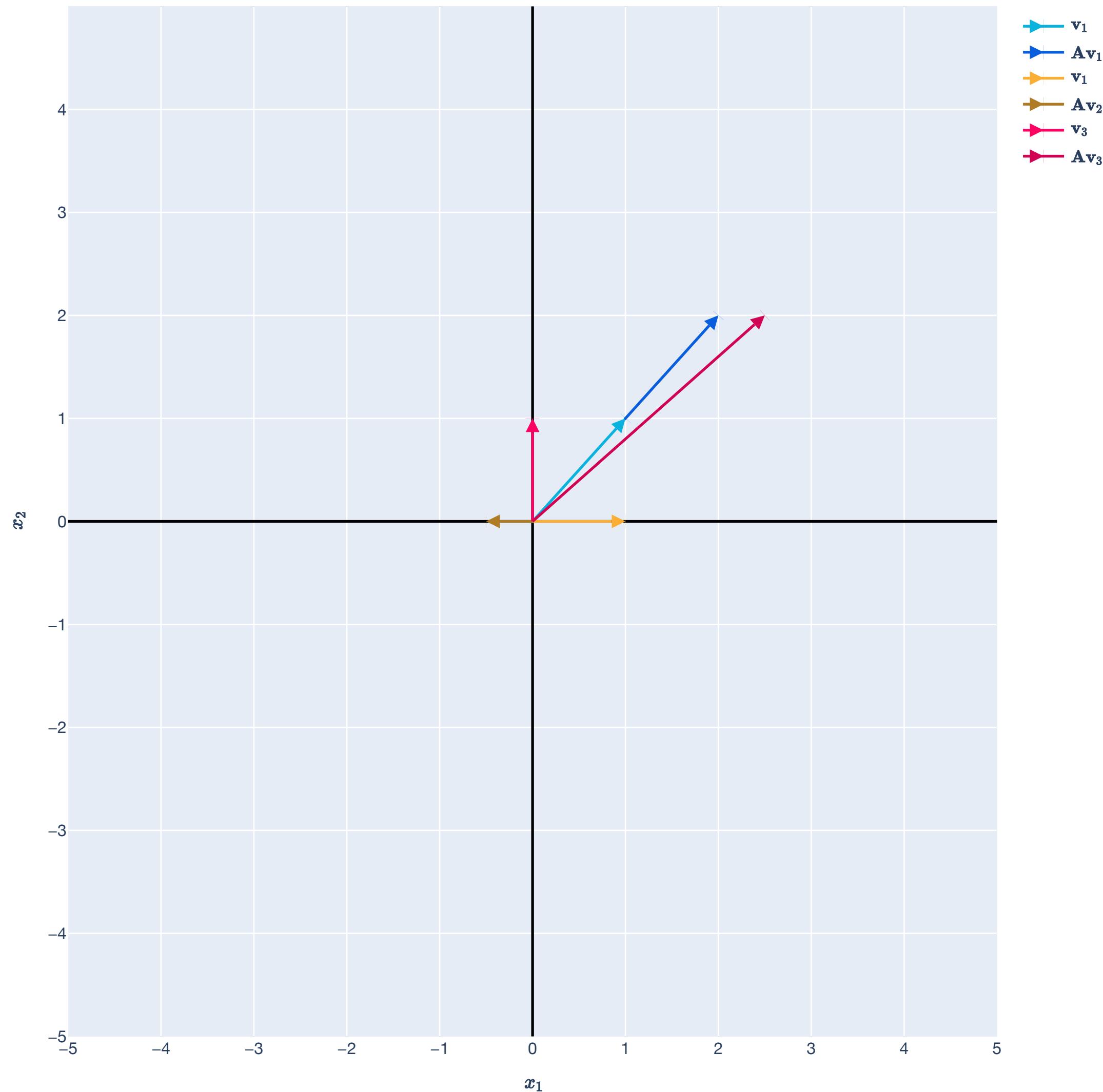
# Eigenvectors and eigenvalues

## Example

Consider the matrix  $A \in \mathbb{R}^{2 \times 2}$  given by

$$A = \begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix}.$$

What happens to the vector  $v_3 = (0, 1)$ ?



# Eigenvectors and eigenvalues

## Example

Consider the matrix  $\mathbf{A} \in \mathbb{R}^{2 \times 2}$  given by

$$\mathbf{A} = \begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix}.$$

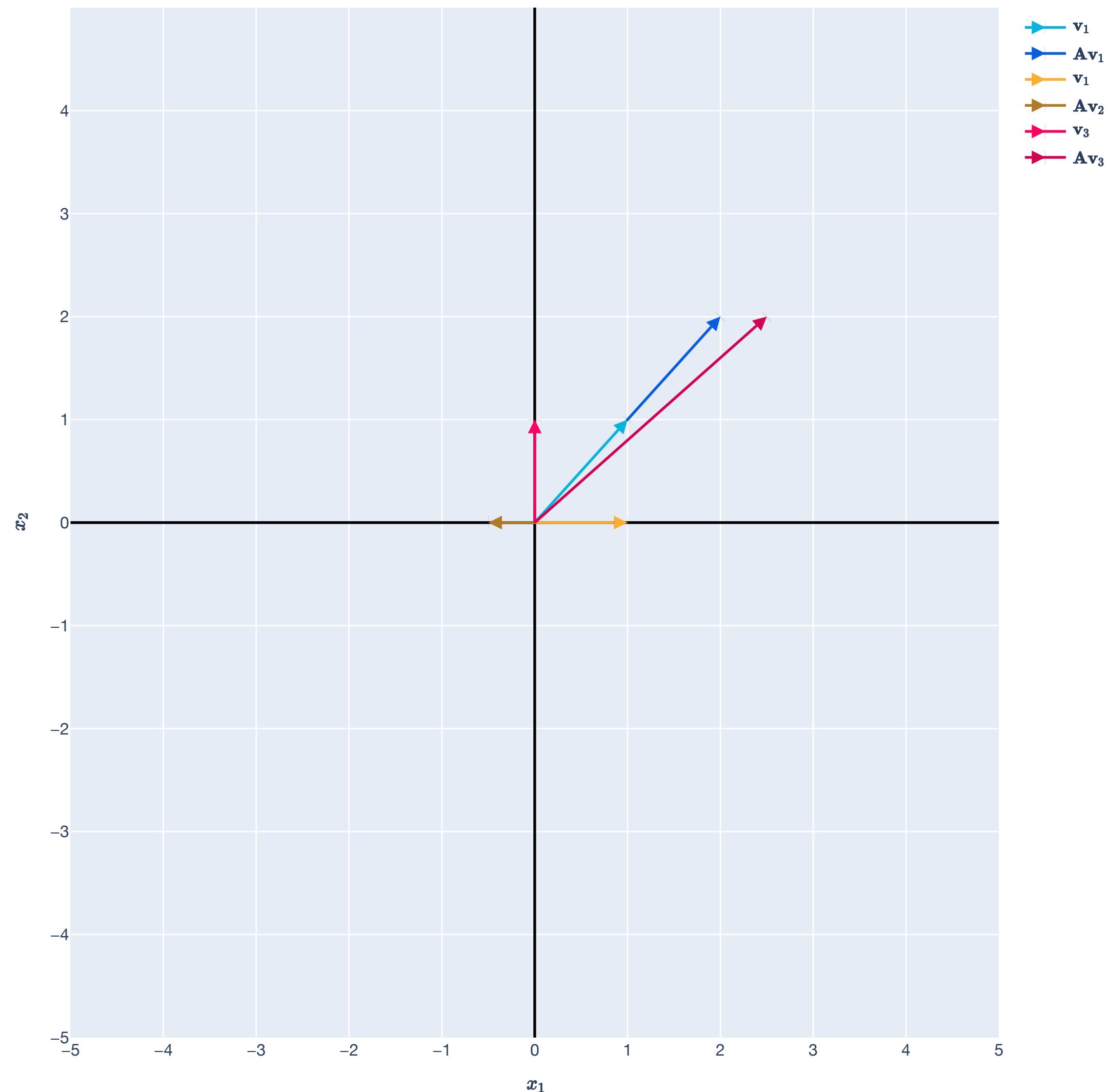
Eigenvectors (with eigenvalues  $\lambda_1 = 2$  and  $\lambda_2 = -1/2$ ):

$$\begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix} = 2 \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -1/2 \\ 0 \end{bmatrix} = -\frac{1}{2} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Not an eigenvector:

$$\begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 5/2 \\ 2 \end{bmatrix}$$



# Eigenvectors and eigenvalues

## Example

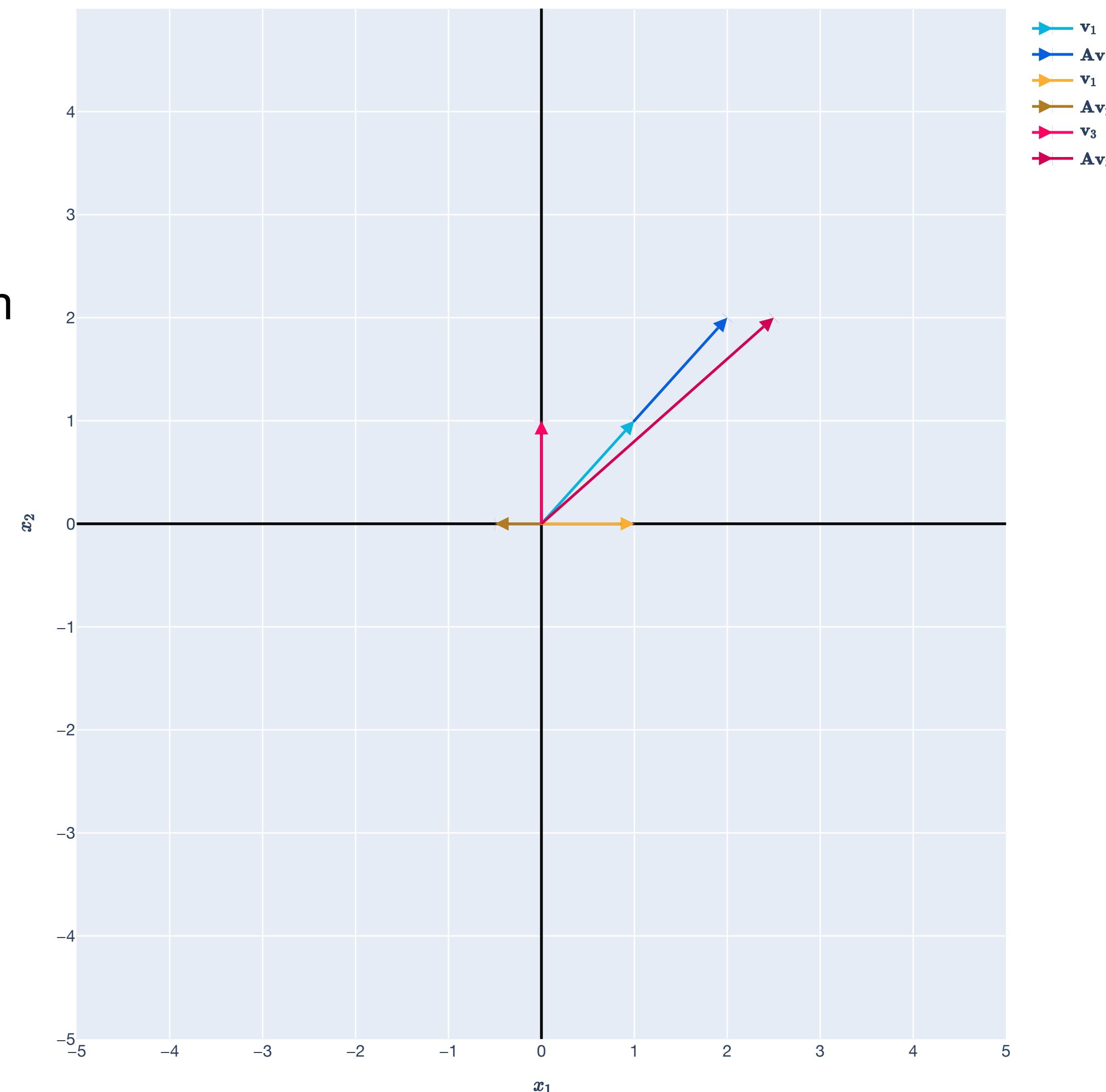
$$A = \begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix}$$

$v_1 = (1,1)$  and  $v_2 = (1,0)$  are linearly independent – they form a basis for  $\mathbb{R}^2$ .

We can write any  $x \in \mathbb{R}^2$  in terms of  $v_1$  and  $v_2$ :

$$x = a v_1 + b v_2.$$

$$x = \begin{bmatrix} \uparrow & \uparrow \\ v_1 & v_2 \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$



# Eigenvectors and eigenvalues

## Example

$$\mathbf{A} = \begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix}$$

$\mathbf{v}_1 = (1,1)$  and  $\mathbf{v}_2 = (1,0)$  are linearly independent eigenvectors – they form a basis for  $\mathbb{R}^2$ . Their eigenvalues are  $\lambda_1 = 2$  and  $\lambda_2 = -1/2$ .

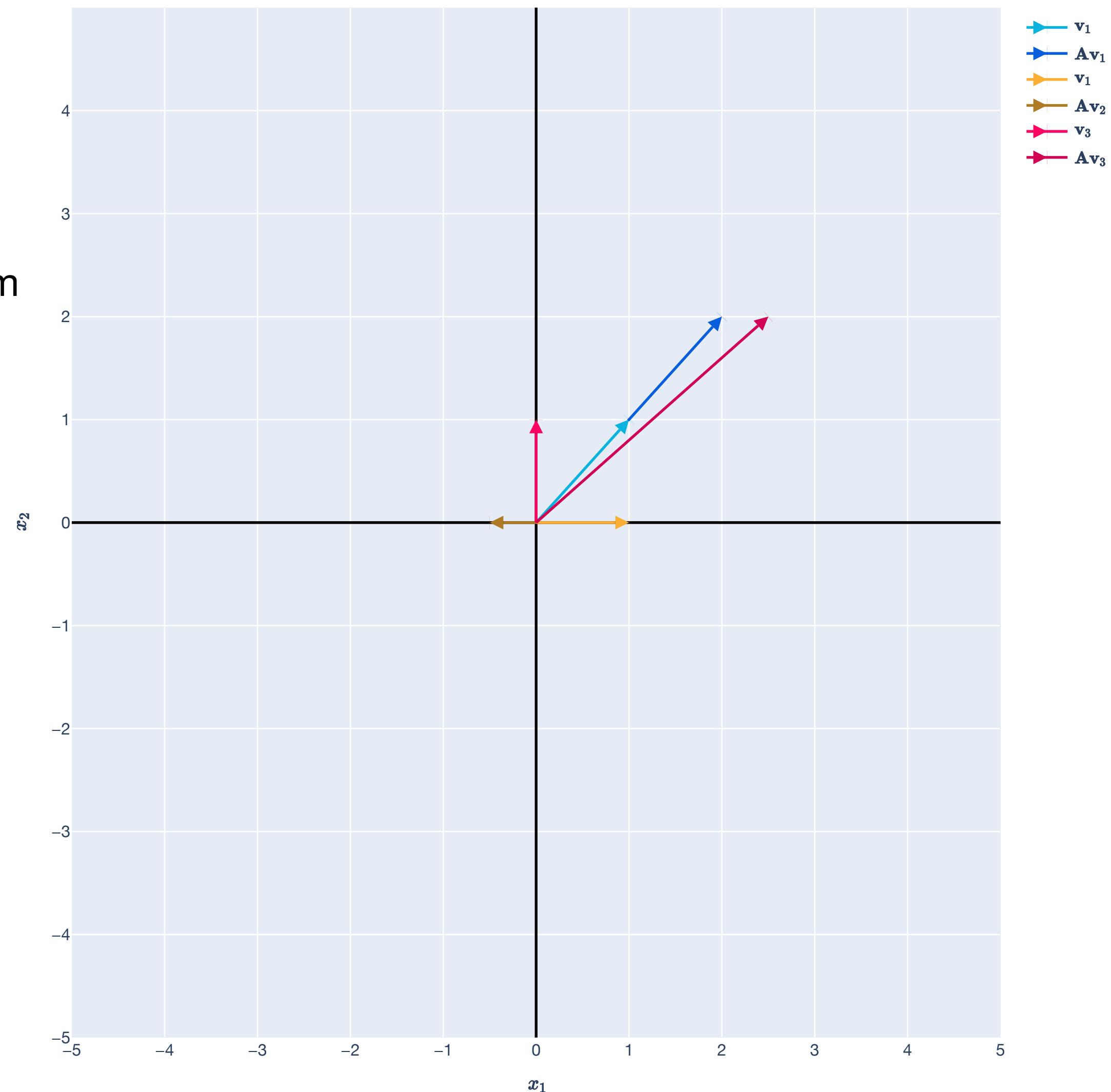
We can write any  $\mathbf{x} \in \mathbb{R}^2$  in terms of  $\mathbf{v}_1$  and  $\mathbf{v}_2$ :

$$\mathbf{x} = a\mathbf{v}_1 + b\mathbf{v}_2.$$

$$\mathbf{x} = \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{v}_1 & \mathbf{v}_2 \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$

Repeated multiplication:

$$\mathbf{A}^t \mathbf{x} = \mathbf{A}^t(a\mathbf{v}_1 + b\mathbf{v}_2) = a\mathbf{A}^t\mathbf{v}_1 + b\mathbf{A}^t\mathbf{v}_2 = a2^t\mathbf{v}_1 + b\left(-\frac{1}{2}\right)^t \mathbf{v}_2$$



# Eigenvectors and eigenvalues

## Example

$$\mathbf{A} = \begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix}$$

$\mathbf{v}_1 = (1, 1)$  and  $\mathbf{v}_2 = (1, 0)$  are linearly independent eigenvectors – they form a basis for  $\mathbb{R}^2$ . Their eigenvalues are  $\lambda_1 = 2$  and  $\lambda_2 = -1/2$ .

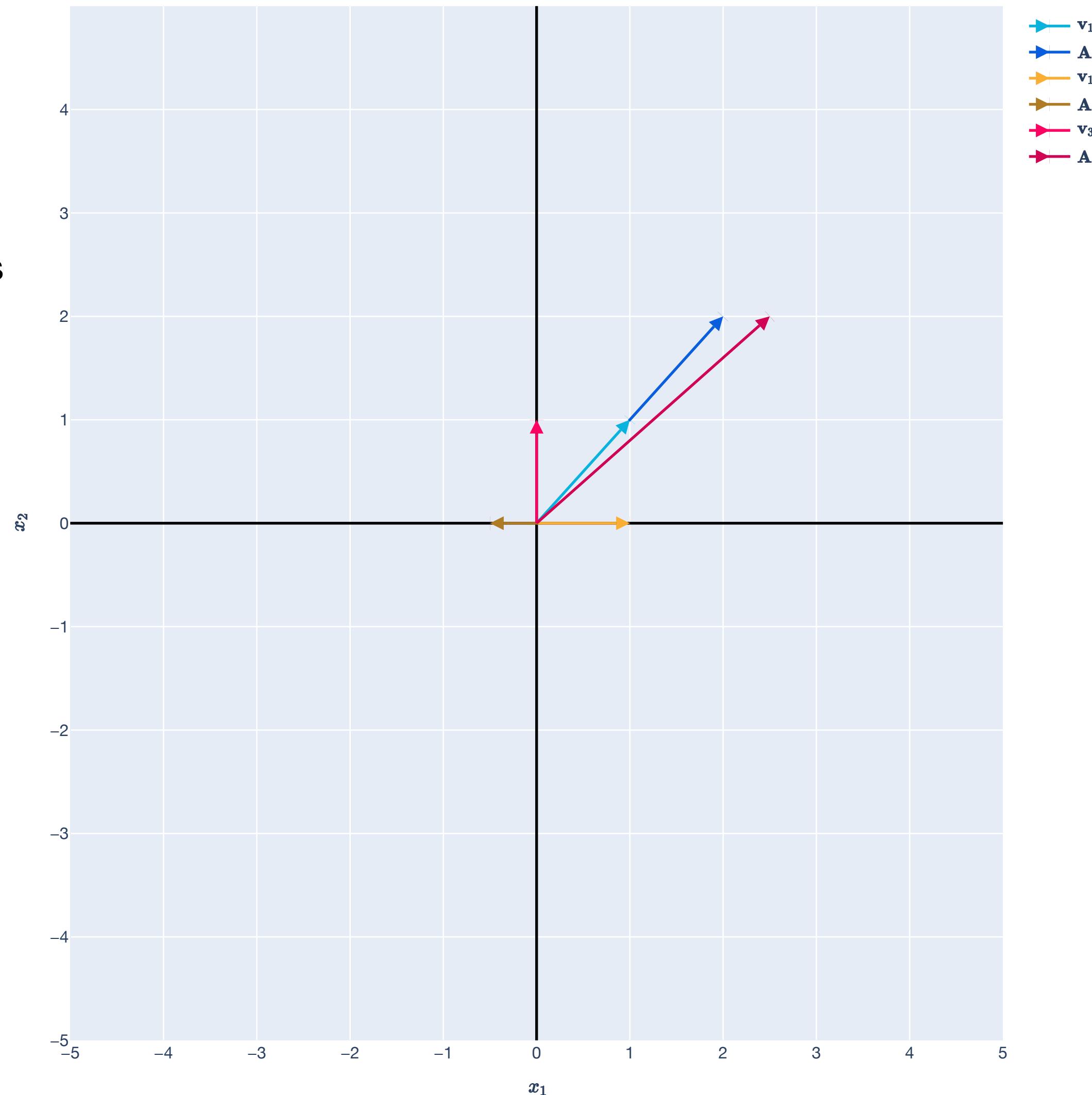
We can write any  $\mathbf{x} \in \mathbb{R}^2$  in terms of  $\mathbf{v}_1$  and  $\mathbf{v}_2$ :

$$\mathbf{x} = a\mathbf{v}_1 + b\mathbf{v}_2.$$

$$\mathbf{x} = \underbrace{\begin{bmatrix} \uparrow & \uparrow \\ \mathbf{v}_1 & \mathbf{v}_2 \\ \downarrow & \downarrow \end{bmatrix}}_{\mathbf{V}} \begin{bmatrix} a \\ b \end{bmatrix} \Rightarrow \begin{bmatrix} a \\ b \end{bmatrix} = \mathbf{V}^{-1}\mathbf{x}$$

Repeated multiplication:

$$\mathbf{A}^t \mathbf{x} = \mathbf{A}^t(a\mathbf{v}_1 + b\mathbf{v}_2) = a\mathbf{A}^t\mathbf{v}_1 + b\mathbf{A}^t\mathbf{v}_2 = a2^t\mathbf{v}_1 + b\left(-\frac{1}{2}\right)^t\mathbf{v}_2$$



# Eigenvectors and eigenvalues

## Example

$$\mathbf{A} = \begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix}$$

$\mathbf{v}_1 = (1, 1)$  and  $\mathbf{v}_2 = (1, 0)$  are linearly independent eigenvectors — they form a basis for  $\mathbb{R}^2$ . Their eigenvalues are  $\lambda_1 = 2$  and  $\lambda_2 = -1/2$ .

We can write any  $\mathbf{x} \in \mathbb{R}^2$  in terms of  $\mathbf{v}_1$  and  $\mathbf{v}_2$ :

$$\mathbf{x} = a\mathbf{v}_1 + b\mathbf{v}_2.$$
$$\mathbf{x} = \underbrace{\begin{bmatrix} \uparrow & \uparrow \\ \mathbf{v}_1 & \mathbf{v}_2 \\ \downarrow & \downarrow \end{bmatrix}}_{\mathbf{V}} \begin{bmatrix} a \\ b \end{bmatrix} \implies \begin{bmatrix} a \\ b \end{bmatrix} = \mathbf{V}^{-1}\mathbf{x}$$

Repeated multiplication:

$$\mathbf{A}^t \mathbf{x} = \mathbf{A}^t(a\mathbf{v}_1 + b\mathbf{v}_2) = a\mathbf{A}^t\mathbf{v}_1 + b\mathbf{A}^t\mathbf{v}_2 = a2^t\mathbf{v}_1 + b\left(-\frac{1}{2}\right)^t\mathbf{v}_2 \implies \mathbf{A}^t \mathbf{x} = \mathbf{V} \begin{bmatrix} 2^t & 0 \\ 0 & (-1/2)^t \end{bmatrix} \mathbf{V}^{-1} \mathbf{x}$$

# Eigenvectors and eigenvalues

## Example

Repeated multiplication:

$$\mathbf{A}^t \mathbf{x} = \mathbf{A}^t(a\mathbf{v}_1 + b\mathbf{v}_2) = a\mathbf{A}^t\mathbf{v}_1 + b\mathbf{A}^t\mathbf{v}_2 = a2^t\mathbf{v}_1 + b\left(-\frac{1}{2}\right)^t \mathbf{v}_2 \implies \mathbf{A}^t \mathbf{x} = \mathbf{V} \begin{bmatrix} 2^t & 0 \\ 0 & (-1/2)^t \end{bmatrix} \mathbf{V}^{-1} \mathbf{x}$$

Single multiplication:

$$\mathbf{A}\mathbf{x} = \mathbf{V} \begin{bmatrix} 2 & 0 \\ 0 & -1/2 \end{bmatrix} \mathbf{V}^{-1} \mathbf{x}$$

$\mathbf{A} = \mathbf{V}\Lambda\mathbf{V}^{-1}$ , where  $\Lambda \in \mathbb{R}^{2 \times 2}$  is diagonal.

# Eigendecomposition

## Definition

**Prop (Eigendecomposition of a diagonalizable matrix).** Let  $\mathbf{A} \in \mathbb{R}^{d \times d}$  be a matrix with  $d$  linearly independent eigenvectors

$$\mathbf{A}\mathbf{v}_1 = \lambda_1\mathbf{v}_1$$

⋮

$$\mathbf{A}\mathbf{v}_d = \lambda_d\mathbf{v}_d$$

Then,  $\mathbf{A}$  has the [\*eigendecomposition\*](#):

$$\mathbf{A} = \mathbf{V}\Lambda\mathbf{V}^{-1} = \begin{bmatrix} \mathbf{v}_1 & \dots & \mathbf{v}_d \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \dots & \lambda_d \end{bmatrix} \begin{bmatrix} \mathbf{v}_1 & \dots & \mathbf{v}_d \end{bmatrix}^{-1}.$$

Such a matrix is said to be [\*diagonalizable\*](#).

# Eigendecomposition

## Example

$A = \begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix}$  has the eigenvectors  $\mathbf{v}_1 = (1,1)$  and  $\mathbf{v}_2 = (1,0)$ :

$$A\mathbf{v}_1 = 2\mathbf{v}_1 \text{ and } A\mathbf{v}_2 = -\frac{1}{2}\mathbf{v}_2.$$

$\mathbf{v}_1$  and  $\mathbf{v}_2$  are *linearly independent*, so  $A$  is *diagonalizable* with *eigendecomposition*:

$$A = Q\Lambda Q^{-1}$$

$$\begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & -1/2 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & -1 \end{bmatrix}$$

# Eigendecomposition

## Example

$A = \begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix}$  has the eigenvectors  $v_1 = (1,1)$  and  $v_2 = (1,0)$ :

$$Av_1 = 2v_1 \text{ and } Av_2 = -\frac{1}{2}v_2.$$

$v_1$  and  $v_2$  are *linearly independent*, so  $A$  is *diagonalizable* with *eigendecomposition*:

$$A = Q\Lambda Q^{-1}$$

$$\begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & -1/2 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & -1 \end{bmatrix}$$

*Question:* But when do (square) matrices have a basis of eigenvectors?

# Eigendecomposition

## Connection with SVD

# Connection with SVD

## Eigendecomposition from SVD

Eigendecomposition only applies to *square* matrices  $\mathbf{A} \in \mathbb{R}^{d \times d}$ :

$$\mathbf{A} = \mathbf{Q}\Lambda\mathbf{Q}^{-1}.$$

The SVD applies to *any* matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ :

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top.$$

# Connection with SVD

## Eigendecomposition from SVD

The SVD applies to *any* matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ :

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top.$$

Consider the square matrix  $\mathbf{A} = \mathbf{X}^\top\mathbf{X} \in \mathbb{R}^{d \times d}$ . By the SVD:

$$\begin{aligned}\mathbf{A} &= \mathbf{X}^\top\mathbf{X} \\ &= \mathbf{V}\Sigma^\top\mathbf{U}^\top\mathbf{U}\Sigma\mathbf{V}^\top \\ &= \mathbf{V}\Sigma^\top\Sigma\mathbf{V}^\top\end{aligned}$$

# Connection with SVD

## Eigendecomposition from SVD

The SVD applies to *any* matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ :

$$\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top.$$

Consider the square matrix  $\mathbf{A} = \mathbf{X}^\top\mathbf{X} \in \mathbb{R}^{d \times d}$ . By the SVD:

$$\mathbf{A} = \underbrace{\mathbf{V}}_{d \times d} \underbrace{\Sigma^\top \Sigma}_{d \times d} \underbrace{\mathbf{V}^\top}_{d \times d}$$

The *eigendecomposition* of  $\mathbf{A}$  is:

$$\mathbf{A} = \underbrace{\mathbf{Q}}_{d \times d} \underbrace{\Lambda}_{d \times d} \underbrace{\mathbf{Q}^{-1}}_{d \times d}$$

# Connection with SVD

## Eigendecomposition from SVD

**Theorem (SVD and Eigendecomposition).** Let  $\mathbf{X} \in \mathbb{R}^{n \times d}$  be a matrix with  $\text{rank}(\mathbf{X}) = r$  and  $\mathbf{A} = \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{d \times d}$ . Let the SVD of  $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$  have singular values

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0,$$

and let  $\mathbf{v}_1, \dots, \mathbf{v}_d$  be the columns of  $\mathbf{V} \in \mathbb{R}^{d \times d}$ . Then, each  $\mathbf{v}_i$  is an eigenvector for  $\mathbf{A}$  with corresponding eigenvalue  $\lambda_i = \sigma_i^2$ , and the eigendecomposition of  $\mathbf{A}$  is:

$$\mathbf{A} = \mathbf{V}\Lambda\mathbf{V}^\top,$$

where  $\Lambda \in \mathbb{R}^{d \times d}$  is the diagonal matrix with entries  $\lambda_i = \sigma_i^2$  for  $i \in [d]$ .

# Connection with SVD

## Eigendecomposition from SVD

Therefore, if  $\mathbf{A} = \mathbf{X}^\top \mathbf{X}$  (for any matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$ ), we know that we have  $d$  linearly independent eigenvectors – this is a case where  $\mathbf{A}$  is diagonalizable!

Moreover, the diagonalization looks like:

$$\mathbf{A} = \mathbf{V}\Lambda\mathbf{V}^\top$$

where  $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$  is the SVD.

# **Positive Semidefinite Matrices**

## Definition and Connections

# Positive Semidefinite (PSD) Matrices

## First definition

A square matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is **positive semidefinite (PSD)** if there exists a matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  such that:

$$\mathbf{A} = \mathbf{X}^T \mathbf{X}.$$

*Note: If you've seen PSD matrices before, this isn't the usual definition (but it's equivalent, as we'll see in a bit).*

# Positive Semidefinite (PSD) Matrices

## Symmetry of PSD Matrices

A square matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is positive semidefinite (PSD) if there exists a matrix  $\mathbf{X} \in \mathbb{R}^{n \times d}$  such that:

$$\mathbf{A} = \mathbf{X}^\top \mathbf{X}.$$

Prop (Symmetry of PSD matrices). All positive semidefinite matrices are symmetric. If  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is PSD, then

$$\mathbf{A} = \mathbf{A}^\top.$$

# Positive Semidefinite (PSD) Matrices

## Example

$$A = \begin{bmatrix} 5/2 & 3/2 \\ 3/2 & 5/2 \end{bmatrix}$$
 is positive semidefinite.

# Positive Semidefinite (PSD) Matrices

## Example

$$A = \begin{bmatrix} 5/2 & 3/2 \\ 3/2 & 5/2 \end{bmatrix} \text{ is positive semidefinite.}$$

Its “square root” is the matrix

$$X = \begin{bmatrix} \frac{2}{\sqrt{2}} & \frac{2}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ 0 & 0 \end{bmatrix}.$$

To verify:

$$X^T X = \begin{bmatrix} \frac{2}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \\ \frac{2}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \frac{2}{\sqrt{2}} & \frac{2}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 5/2 & 3/2 \\ 3/2 & 5/2 \end{bmatrix} = A$$

# PSD Matrices and Eigendecomposition

## Connection to eigenvalues

By Theorem (SVD and Eigendecomposition), if  $\mathbf{A}$  is PSD with  $\mathbf{A} = \mathbf{X}^\top \mathbf{X}$  and  $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$  then

$$\mathbf{A} = \mathbf{V}\Lambda\mathbf{V}^\top,$$

with orthonormal eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_d$

and nonnegative eigenvalues  $\lambda_1 = \sigma_1^2, \dots, \lambda_d = \sigma_d^2$

The reverse direction is also true!

# PSD Matrices and Eigendecomposition

## Second definition

A square matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is **positive semidefinite (PSD)** if  $\mathbf{A}$  has  $d$  eigenvectors forming an orthonormal basis for  $\mathbb{R}^d$  with corresponding nonnegative eigenvalues  $\lambda_1, \dots, \lambda_d \geq 0$ .

# Positive Semidefinite (PSD) Matrices

## Example

$\mathbf{A} = \begin{bmatrix} 5/2 & 3/2 \\ 3/2 & 5/2 \end{bmatrix}$  is positive semidefinite.

It has the eigenvectors  $\mathbf{v}_1 = \left( \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right)$  and  $\mathbf{v}_2 = \left( \frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}} \right)$ :

$$\mathbf{A}\mathbf{v}_1 = \begin{bmatrix} 5/2 & 3/2 \\ 3/2 & 5/2 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix} = \begin{bmatrix} 4/\sqrt{2} \\ 4/\sqrt{2} \end{bmatrix} = 4 \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix} \implies \lambda_1 = 4$$

$$\mathbf{A}\mathbf{v}_2 = \begin{bmatrix} 5/2 & 3/2 \\ 3/2 & 5/2 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \end{bmatrix} = \begin{bmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \end{bmatrix} \implies \lambda_2 = 1$$

The eigenvectors are orthonormal and  $\lambda_1, \lambda_2 \geq 0$ , so  $\mathbf{A} = \mathbf{V}\Lambda\mathbf{V}^\top$ .

# Positive Semidefinite (PSD) Matrices

## Third definition

A square matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is positive semidefinite (PSD) if, for any  $\mathbf{x} \in \mathbb{R}^d$ ,

$$\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0.$$

This is often taken as the definition of PSD (but it is equivalent to the other two definitions in previous slides).

# Positive Semidefinite (PSD) Matrices

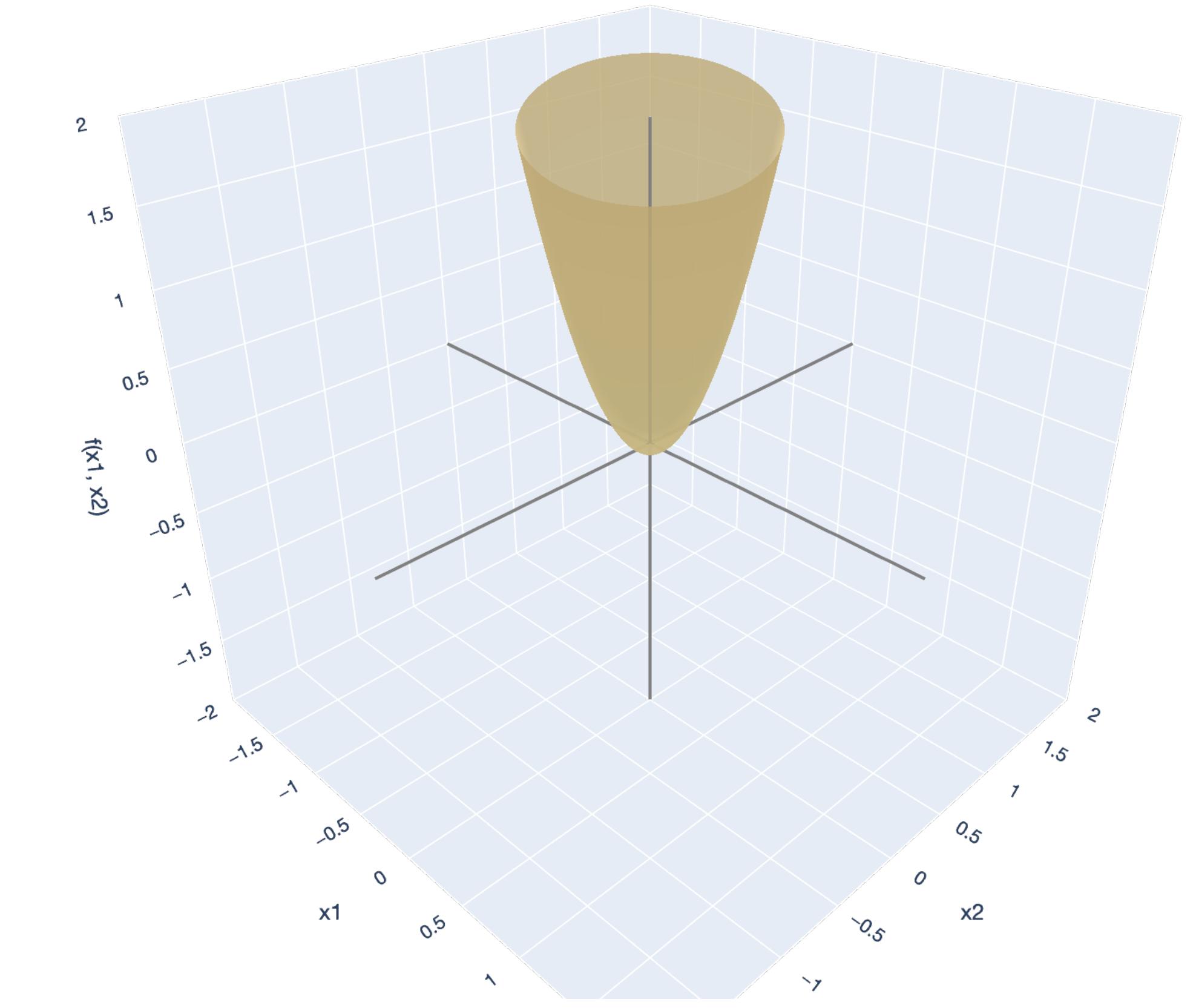
## Example

$A = \begin{bmatrix} 5/2 & 3/2 \\ 3/2 & 5/2 \end{bmatrix}$  is positive semidefinite.

Consider any vector  $\mathbf{x} = (x_1, x_2) \in \mathbb{R}^d$ .

$$\mathbf{x}^\top A \mathbf{x} = [x_1 \ x_2] \begin{bmatrix} 5/2 & 3/2 \\ 3/2 & 5/2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = [x_1 \ x_2] \begin{bmatrix} (5/2)x_1 + (3/2)x_2 \\ (3/2)x_1 + (5/2)x_2 \end{bmatrix}$$

$$\mathbf{x}^\top A \mathbf{x} = (5/2)x_1^2 + 3x_1x_2 + (5/2)x_2^2$$



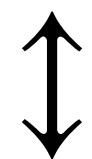
x1-axis   x2-axis   f(x1, x2)-axis

# Positive Semidefinite (PSD) Matrices

## All definitions

A square matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is positive semidefinite (PSD) if...

there exists  $\mathbf{X} \in \mathbb{R}^{n \times d}$  such that  $\mathbf{A} = \mathbf{X}^\top \mathbf{X}$ .



all eigenvalues of  $\mathbf{A}$  are nonnegative:  $\lambda_1 \geq 0, \dots, \lambda_d \geq 0$ .



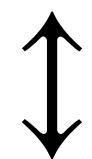
$\mathbf{x}^\top \mathbf{A} \mathbf{x} \geq 0$  for any  $\mathbf{x} \in \mathbb{R}^d$ .

# Positive Definite (PD) Matrices

## All definitions

A square matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  is **positive definite (PD)** if...

there exists an *invertible matrix*  $\mathbf{X} \in \mathbb{R}^{d \times d}$  such that  $\mathbf{A} = \mathbf{X}^\top \mathbf{X}$ .



all eigenvalues of  $\mathbf{A}$  are positive:  $\lambda_1 > 0, \dots, \lambda_d > 0$ .



$\mathbf{x}^\top \mathbf{A} \mathbf{x} > 0$  for any  $\mathbf{x} \in \mathbb{R}^d$ .

# Spectral Theorem

## Statement

*Question: But when does a square matrix  $A \in \mathbb{R}^{d \times d}$  have a basis of eigenvectors (and, hence, is diagonalizable)?*

A: When  $A$  is positive semidefinite!

But even more generally...

# Spectral Theorem

## Statement

**Theorem (Spectral Theorem).** Let  $\mathbf{A} \in \mathbb{R}^{d \times d}$  be a square, symmetric matrix (i.e.  $\mathbf{A}^\top = \mathbf{A}$ ). Then,  $\mathbf{A}$  is diagonalizable:  $\mathbf{A}$  has an orthonormal basis of  $d$  eigenvectors and an eigendecomposition

$$\mathbf{A} = \mathbf{Q}\Lambda\mathbf{Q}^\top.$$

But, in this generality,  $\lambda_i$  can be negative!

# **Principal Components Analysis**

## Application of Eigendecomposition

# Principal Components Analysis

Example: “Eigenfaces” and facial recognition

Observed: Matrix of *training images*  $\mathbf{X} \in \mathbb{R}^{n \times d}$ :

$$\mathbf{X} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{x}_1 & \dots & \mathbf{x}_d \\ \downarrow & & \downarrow \end{bmatrix} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix}.$$

Each row is a “flattened” image vector. Typically, each pixel is in  $[0, 255]$  for grayscale images.

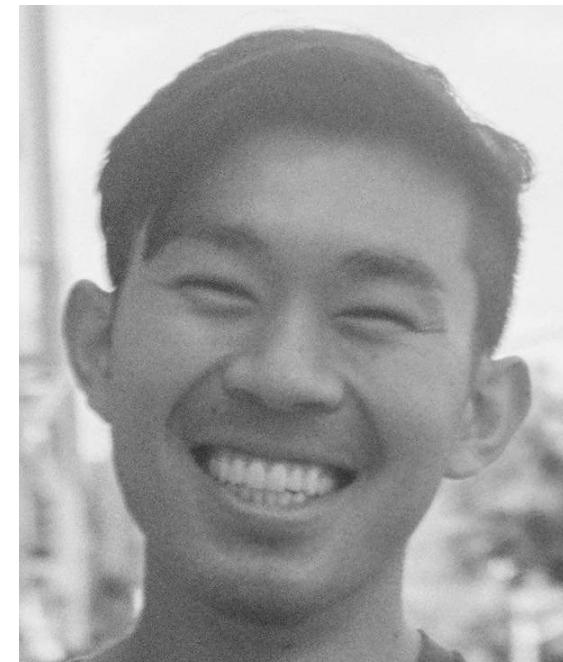
Images are very high-dimensional:  $d = \text{width in pixels} \times \text{height in pixels}$  (e.g.  $d = 1080 \times 1080 = 1,166,400$ ).

# Principal Components Analysis

**Example: “Eigenfaces” and facial recognition**

Consider a dataset of 1,000 grayscale face images  $\mathbf{x}_1, \dots, \mathbf{x}_{1000} \in \mathbb{R}^{1080 \times 1080}$ ...

e.g.  $\mathbf{x}_1 =$



*Naive facial recognition:* Get a new face, linear search over 1,000 faces for the “closest” face (perhaps in Euclidean norm  $\|\mathbf{x} - \mathbf{x}_i\|$ ).

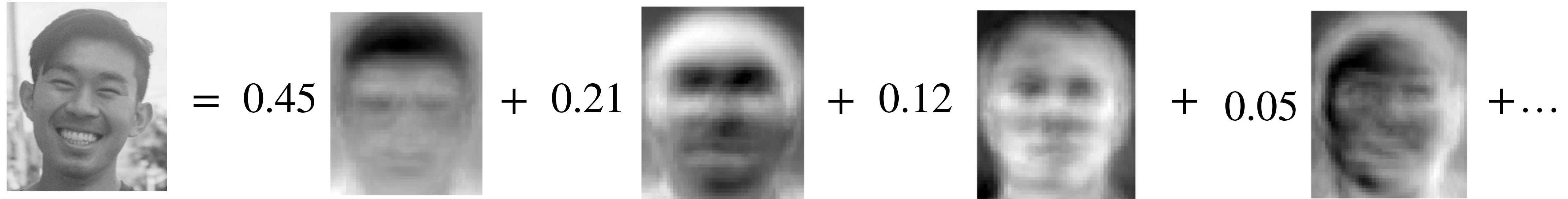
Storage: 1166400 integers  $\times$  1000 images  $\approx$  1 GB.

# Principal Components Analysis

## Example: “Eigenfaces” and facial recognition

Suppose we can find a “basis” of representative faces:  $\mathbf{v}_1, \dots, \mathbf{v}_k$  where  $k \ll n$ .

Then, we can represent any face as a linear combination of the basis faces!



A black and white photograph of a smiling man's face is shown on the left. To its right is an equals sign (=). Following the equals sign are five terms, each consisting of a scalar value and a grayscale image of a face. The scalars are 0.45, 0.21, 0.12, 0.05, and a plus sign (...). The images are progressively more abstract, representing the projection of the original face onto a lower-dimensional subspace defined by the eigenfaces.

$$\text{Smiling Face} = 0.45 \begin{matrix} \text{Eigenface 1} \end{matrix} + 0.21 \begin{matrix} \text{Eigenface 2} \end{matrix} + 0.12 \begin{matrix} \text{Eigenface 3} \end{matrix} + 0.05 \begin{matrix} \text{Eigenface 4} \end{matrix} + \dots$$

*Improved facial recognition:* Store  $k$  “eigenfaces.” Given a new face  $\mathbf{x}_0$ , project the face onto the subspace spanned by the eigenfaces to get  $\Pi(\mathbf{x}_0)$ . Compare  $\Pi(\mathbf{x}_0)$  to each face’s projection in the database in Euclidean norm  $\|\Pi(\mathbf{x}_0) - \Pi(\mathbf{x}_i)\|$ .

# Principal Components Analysis

## Example: PCA in 2D

**Observed:** Matrix of *training points*  $\mathbf{X} \in \mathbb{R}^{n \times 2}$ :

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{bmatrix}.$$

Want to find the directions that most explain the “variance” of the data.

# Principal Components Analysis

## Example: PCA in 2D

**Observed:** Matrix of *training points*  $\mathbf{X} \in \mathbb{R}^{n \times 2}$ :

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{bmatrix}.$$

Want to find the directions that most explain the “variance” of the data.

The matrix  $\mathbf{C} = \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{2 \times 2}$  is the *covariance matrix* of the data.

# Principal Components Analysis

## Example: PCA in 2D

**Observed:** Matrix of *training points*  $\mathbf{X} \in \mathbb{R}^{n \times 2}$ :

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{bmatrix} = \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{x}_1 & \mathbf{x}_2 \\ \downarrow & \downarrow \end{bmatrix}$$

The matrix  $\mathbf{C} = \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{2 \times 2}$  is the *covariance matrix* of the data.

$$\mathbf{C} = \begin{bmatrix} \mathbf{x}_1^\top \mathbf{x}_1 & \mathbf{x}_1^\top \mathbf{x}_2 \\ \mathbf{x}_2^\top \mathbf{x}_1 & \mathbf{x}_2^\top \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \|\mathbf{x}_1\|^2 & \mathbf{x}_1^\top \mathbf{x}_2 \\ \mathbf{x}_1^\top \mathbf{x}_2 & \|\mathbf{x}_2\|^2 \end{bmatrix}$$

# Principal Components Analysis

## Example: PCA in 2D

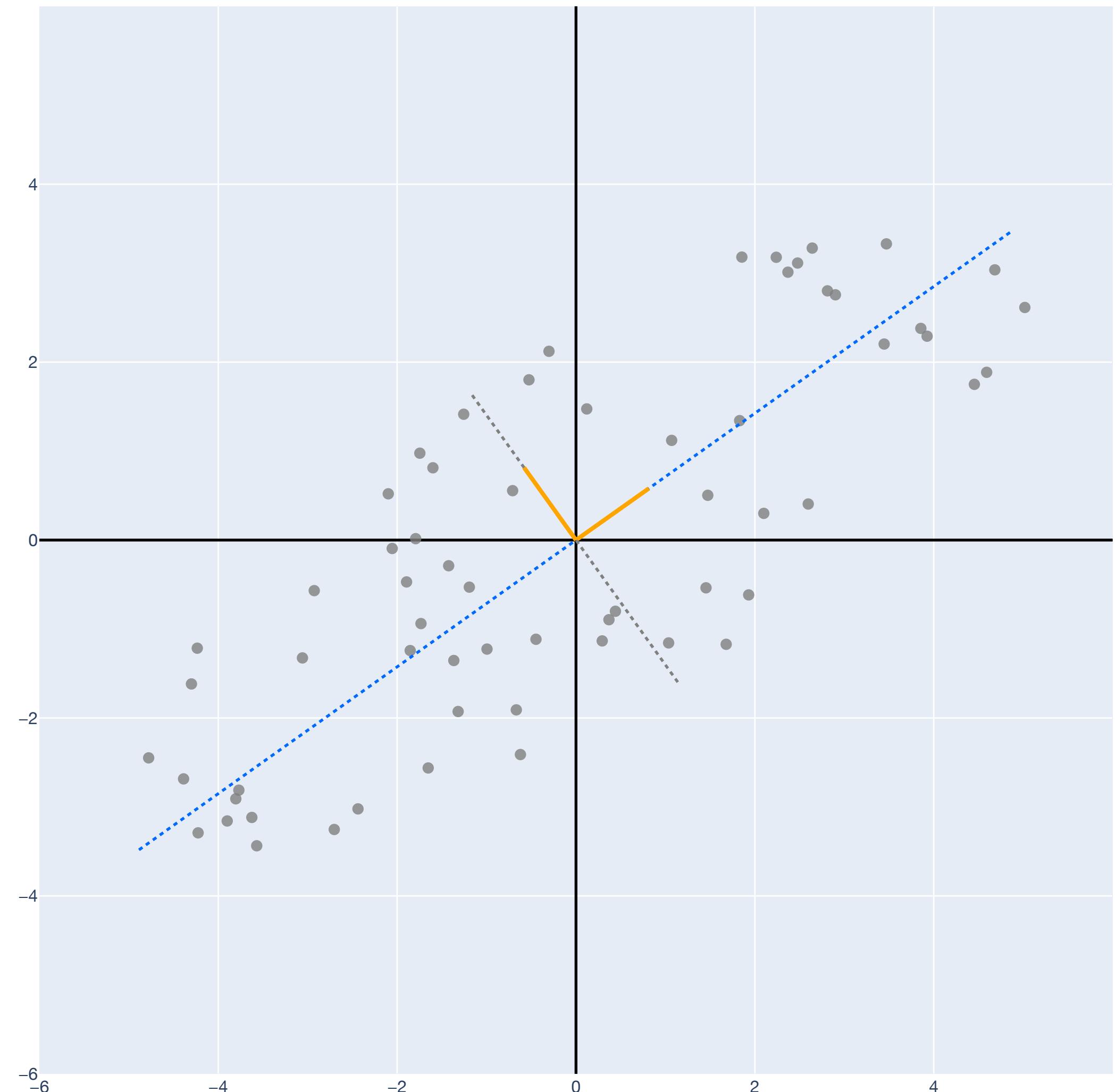
Observed: Matrix of *training points*  $\mathbf{X} \in \mathbb{R}^{n \times 2}$ :

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ \vdots & \vdots \\ x_{n1} & x_{n2} \end{bmatrix} = \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{x}_1 & \mathbf{x}_2 \\ \downarrow & \downarrow \end{bmatrix}$$

The matrix  $\mathbf{C} = \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{2 \times 2}$  is the *covariance matrix* of the data.

$$\mathbf{C} = \begin{bmatrix} \mathbf{x}_1^\top \mathbf{x}_1 & \mathbf{x}_1^\top \mathbf{x}_2 \\ \mathbf{x}_2^\top \mathbf{x}_1 & \mathbf{x}_2^\top \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \|\mathbf{x}_1\|^2 & \mathbf{x}_1^\top \mathbf{x}_2 \\ \mathbf{x}_1^\top \mathbf{x}_2 & \|\mathbf{x}_2\|^2 \end{bmatrix}$$

PCA: Find the ordered set of vectors  $\mathbf{v}_1, \dots, \mathbf{v}_d \in \mathbb{R}^d$  that explain the most variance to least variance in the data.



# Derivation of PCA

## Eigendecomposition and PCA

*PCA = Eigendecomposition of the covariance matrix!*

Consider a (column-centered) dataset  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and construct its covariance matrix  $\mathbf{C} = \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{d \times d}$ . By definition,  $\mathbf{C}$  is positive semidefinite.

Therefore, it is diagonalizable with eigendecomposition:

$$\mathbf{C} = \mathbf{X}^\top \mathbf{X} = \mathbf{V} \Lambda \mathbf{V}^\top, \text{ with eigenvectors } \mathbf{v}_1, \dots, \mathbf{v}_d.$$

With eigenvectors ordered  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ , choose a cutoff point  $k \ll d$ , and keep eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_k$ .

The eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_k$  give an orthonormal basis for a  $k$ -dimensional subspace.

# Derivation of PCA

## Eigendecomposition and PCA

*PCA = Eigendecomposition of the covariance matrix!*

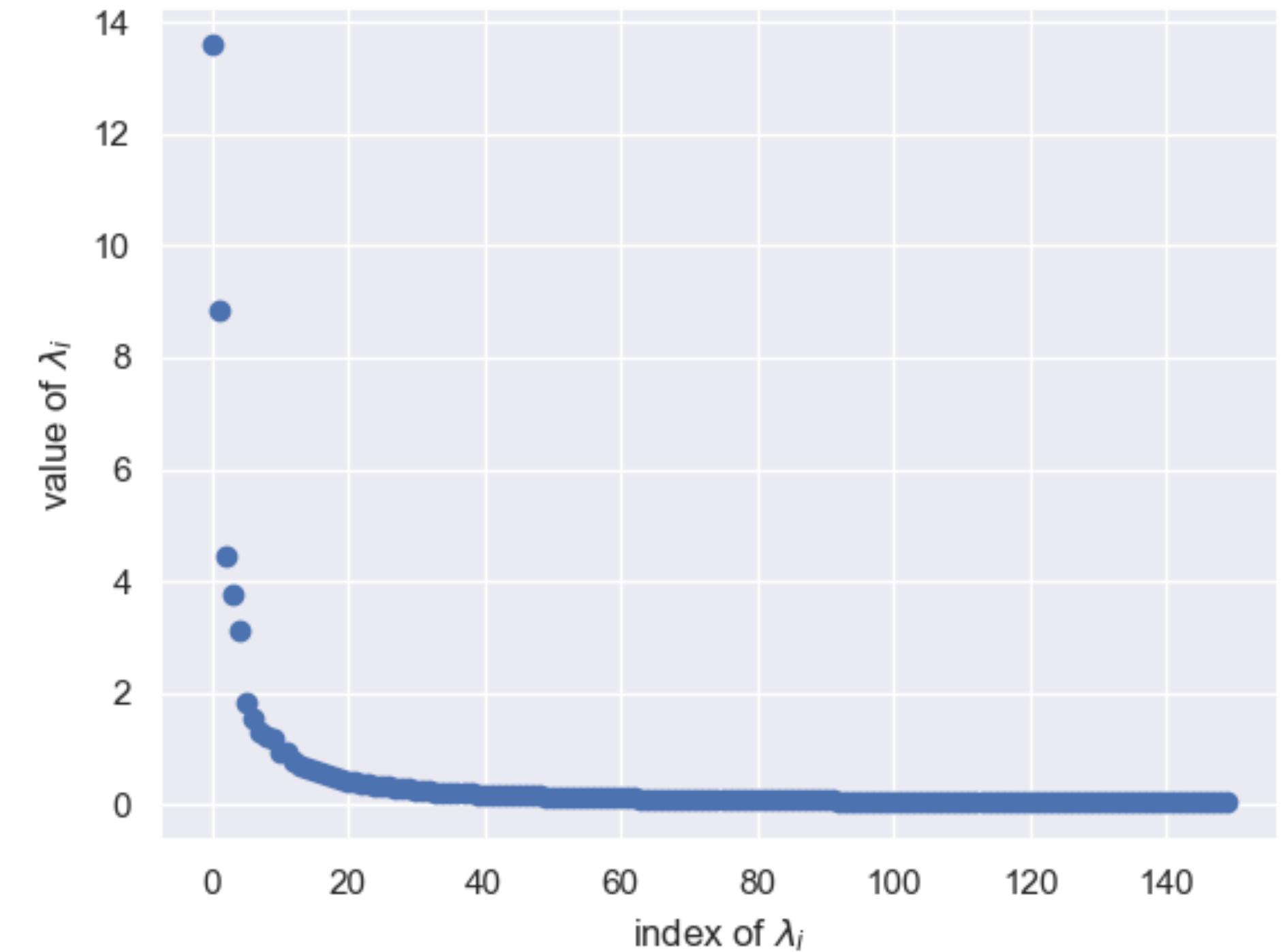
Consider a (column-centered) dataset  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and construct its covariance matrix  $\mathbf{C} = \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{d \times d}$ . By definition,  $\mathbf{C}$  is positive semidefinite.

Therefore, it is diagonalizable with eigendecomposition:

$$\mathbf{C} = \mathbf{X}^\top \mathbf{X} = \mathbf{V} \Lambda \mathbf{V}^\top, \text{ with eigenvectors } \mathbf{v}_1, \dots, \mathbf{v}_d.$$

With eigenvectors ordered  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$ , choose a cutoff point  $k \ll d$ , and keep eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_k$ .

The eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_k$  give an orthonormal basis for a  $k$ -dimensional subspace.



# Derivation of PCA

## Eigendecomposition and PCA

*PCA = Eigendecomposition of the covariance matrix!*

Consider a (column-centered) dataset  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and construct its covariance matrix  $\mathbf{C} = \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{d \times d}$ . By definition,  $\mathbf{C}$  is positive semidefinite.

Therefore, it is diagonalizable with eigendecomposition:

$$\mathbf{C} = \mathbf{X}^\top \mathbf{X} = \mathbf{V} \Lambda \mathbf{V}^\top.$$

*(Could have also just taken the right singular vectors of  $\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^\top$  if we have efficient algorithm to find the SVD – true in practice).*

# Least Squares

## Interpretation of Eigenvalues

# Regression Setup

**Observed:** Matrix of *training samples*  $\mathbf{X} \in \mathbb{R}^{n \times d}$  and vector of *training labels*  $\mathbf{y} \in \mathbb{R}^d$ .

$$\mathbf{X} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{x}_1 & \dots & \mathbf{x}_d \\ \downarrow & & \downarrow \end{bmatrix} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix}.$$

**Unknown:** *Weight vector*  $\mathbf{w} \in \mathbb{R}^d$  with weights  $w_1, \dots, w_d$ .

**Goal:** For each  $i \in [n]$ , we predict:  $\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i = w_1 x_{i1} + \dots + w_d x_{id} \in \mathbb{R}$ .

Choose a weight vector that “fits the training data”:  $\mathbf{w} \in \mathbb{R}^d$  such that  $y_i \approx \hat{y}_i$  for  $i \in [n]$ , or:

$$\mathbf{X}\mathbf{w} = \hat{\mathbf{y}} \approx \mathbf{y}.$$

# Regression Setup

**Goal:** For each  $i \in [n]$ , we predict:  $\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i = w_1 x_{i1} + \dots + w_d x_{id} \in \mathbb{R}$ .

Choose a weight vector that “fits the training data”:  $\hat{\mathbf{w}} \in \mathbb{R}^d$  such that  $y_i \approx \hat{y}_i$  for  $i \in [n]$ , or:

$$\mathbf{X}\hat{\mathbf{w}} = \hat{\mathbf{y}} \approx \mathbf{y}.$$

To find  $\hat{\mathbf{w}}$ , we follow the *principle of least squares*.

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

# Error in Regression

## Error using least squares model

Choose a weight vector that “fits the training data”:  $\hat{\mathbf{w}} \in \mathbb{R}^d$  such that  $y_i \approx \hat{y}_i$  for  $i \in [n]$ , or:

$$\mathbf{X}\hat{\mathbf{w}} = \hat{\mathbf{y}} \approx \mathbf{y}.$$

But  $\hat{\mathbf{y}}$  might not be a perfect fit to  $\mathbf{y}$ !

Model this using a *true weight vector*  $\mathbf{w}^* \in \mathbb{R}^d$  and an *error term*  $\epsilon = (\epsilon_1, \dots, \epsilon_n) \in \mathbb{R}^n$ .

$$y_i = \mathbf{x}_i^\top \mathbf{w}^* + \epsilon_i \text{ for all } i \in [n]$$

$$\mathbf{y} = \mathbf{X}\mathbf{w}^* + \epsilon$$

# Error in Regression

## Error using least squares model

True labels:  $\mathbf{y} = \mathbf{X}\mathbf{w}^* + \epsilon$ .

What happens when we use the least squares weights  $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ ?

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\mathbf{w}^* + \epsilon) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\mathbf{w}^* + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon \\ &= \mathbf{w}^* + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon\end{aligned}$$

# Error in Regression

## Error using least squares model

True labels:  $\mathbf{y} = \mathbf{X}\mathbf{w}^* + \epsilon$ .

What happens when we use the least squares weights  $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ ?

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\mathbf{w}^* + \epsilon) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\mathbf{w}^* + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon \\ &= \mathbf{w}^* + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon\end{aligned}$$

When  $\epsilon = 0$  ( $\mathbf{y}$  is linearly related to  $\mathbf{X}$ ), this is perfect:  $\hat{\mathbf{w}} = \mathbf{w}^*$ !

# Error in Regression

## Error using least squares model

True labels:  $\mathbf{y} = \mathbf{X}\mathbf{w}^* + \epsilon$ .

What happens when we use the least squares weights  $\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ ?

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X}\mathbf{w}^* + \epsilon) \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}\mathbf{w}^* + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon \\ &= \mathbf{w}^* + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon\end{aligned}$$

When  $\epsilon \neq 0$ , we have an error of  $\hat{\mathbf{w}} - \mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon$ .

# Error in Regression

## Eigendecomposition perspective

Weight vector's error:  $\hat{\mathbf{w}} - \mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \epsilon$ .

We know that  $\mathbf{X}^\top \mathbf{X}$  (the *covariance matrix*) is PSD, so it is diagonalizable:

$$\mathbf{X}^\top \mathbf{X} = \mathbf{V} \Lambda \mathbf{V}^\top \implies (\mathbf{X}^\top \mathbf{X})^{-1} = \mathbf{V}^\top \Lambda^{-1} \mathbf{V}.$$

The inverse of the diagonal matrix  $\Lambda^{-1}$ :

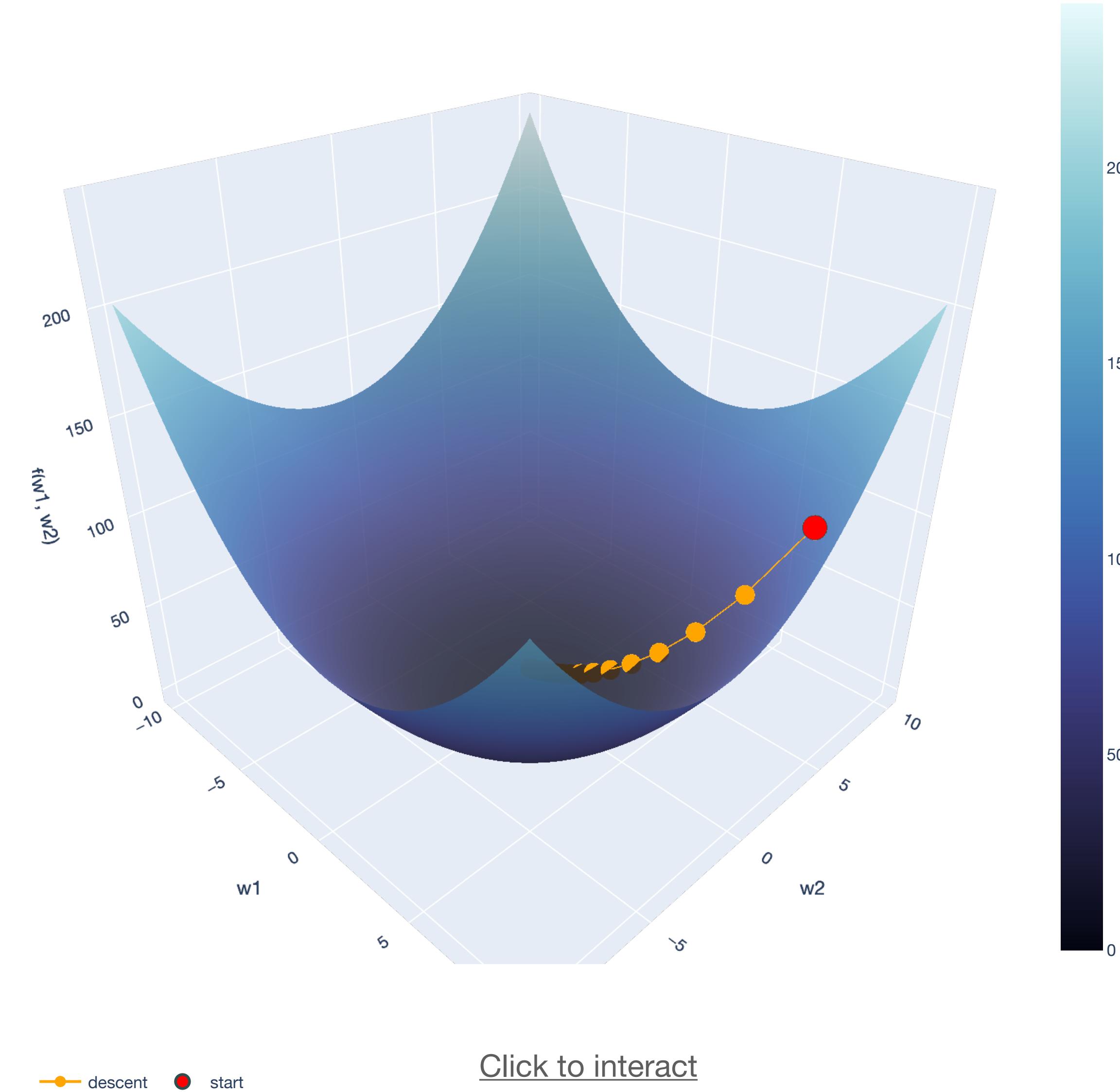
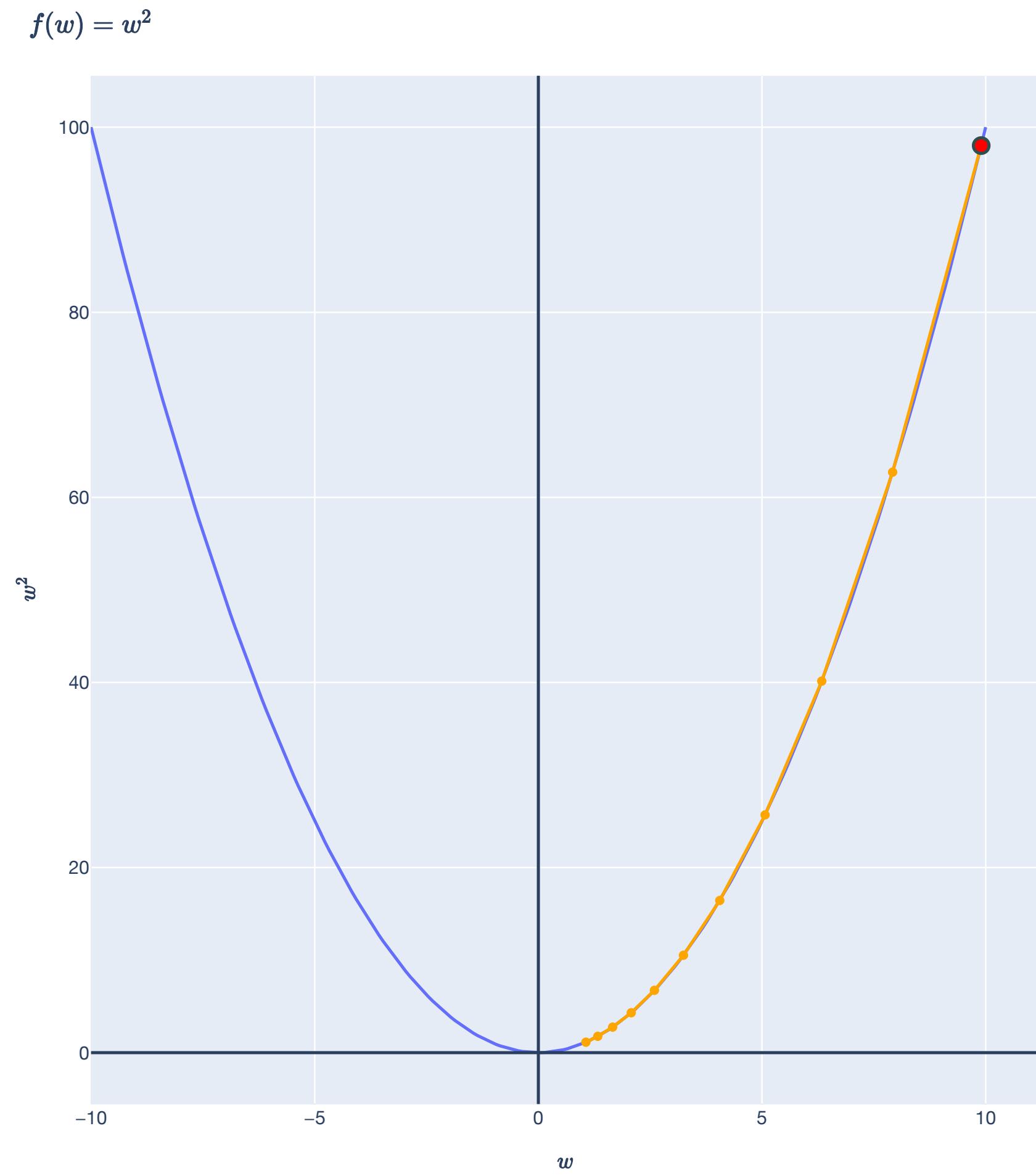
$$\Lambda^{-1} = \begin{bmatrix} 1/\lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 1/\lambda_d \end{bmatrix}, \text{ so if } \lambda_i \text{ is small, the entries of } \hat{\mathbf{w}} \text{ blow up!}$$

# **Gradient Descent**

## Positive Semidefinite Matrices and Convexity

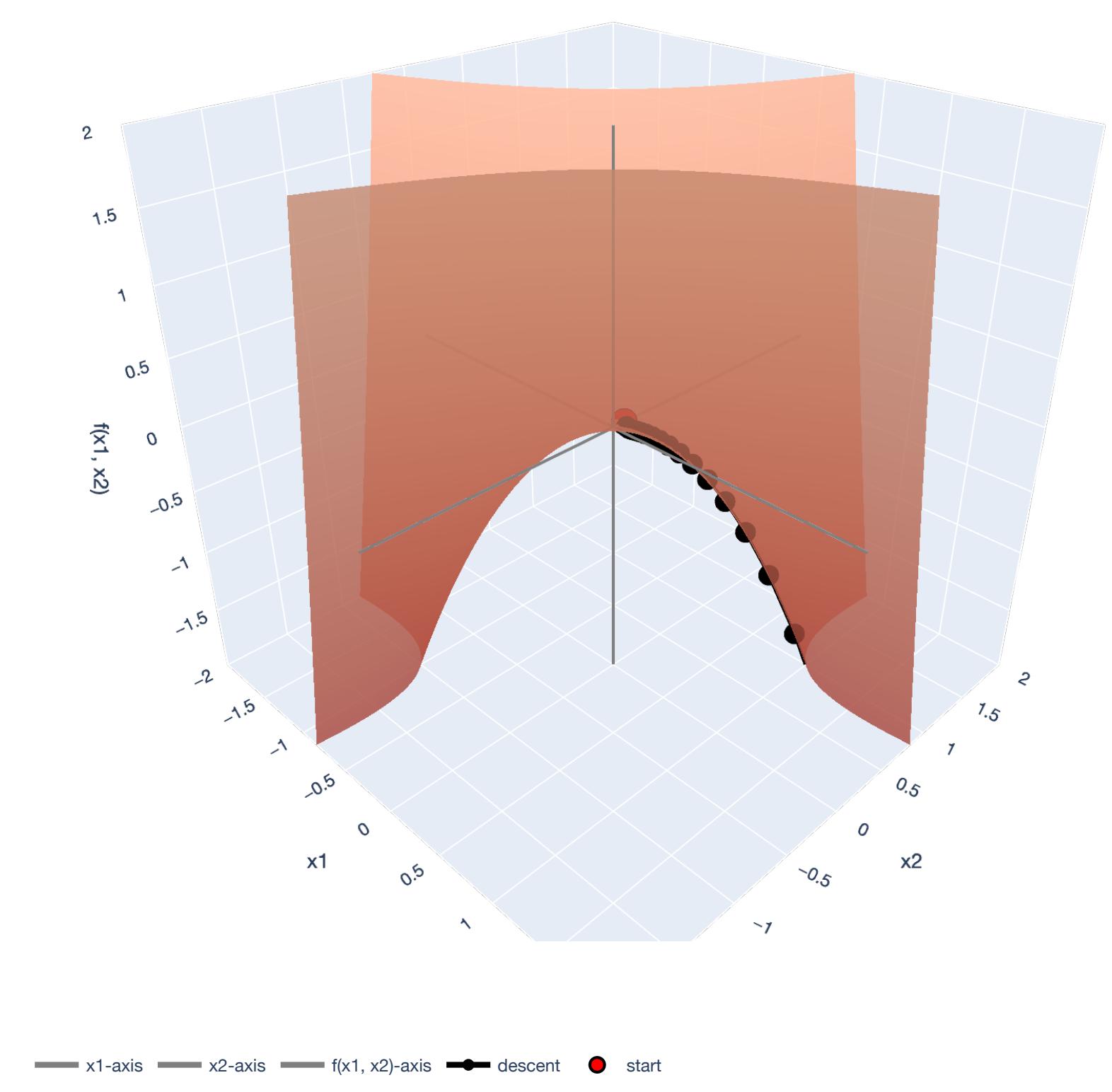
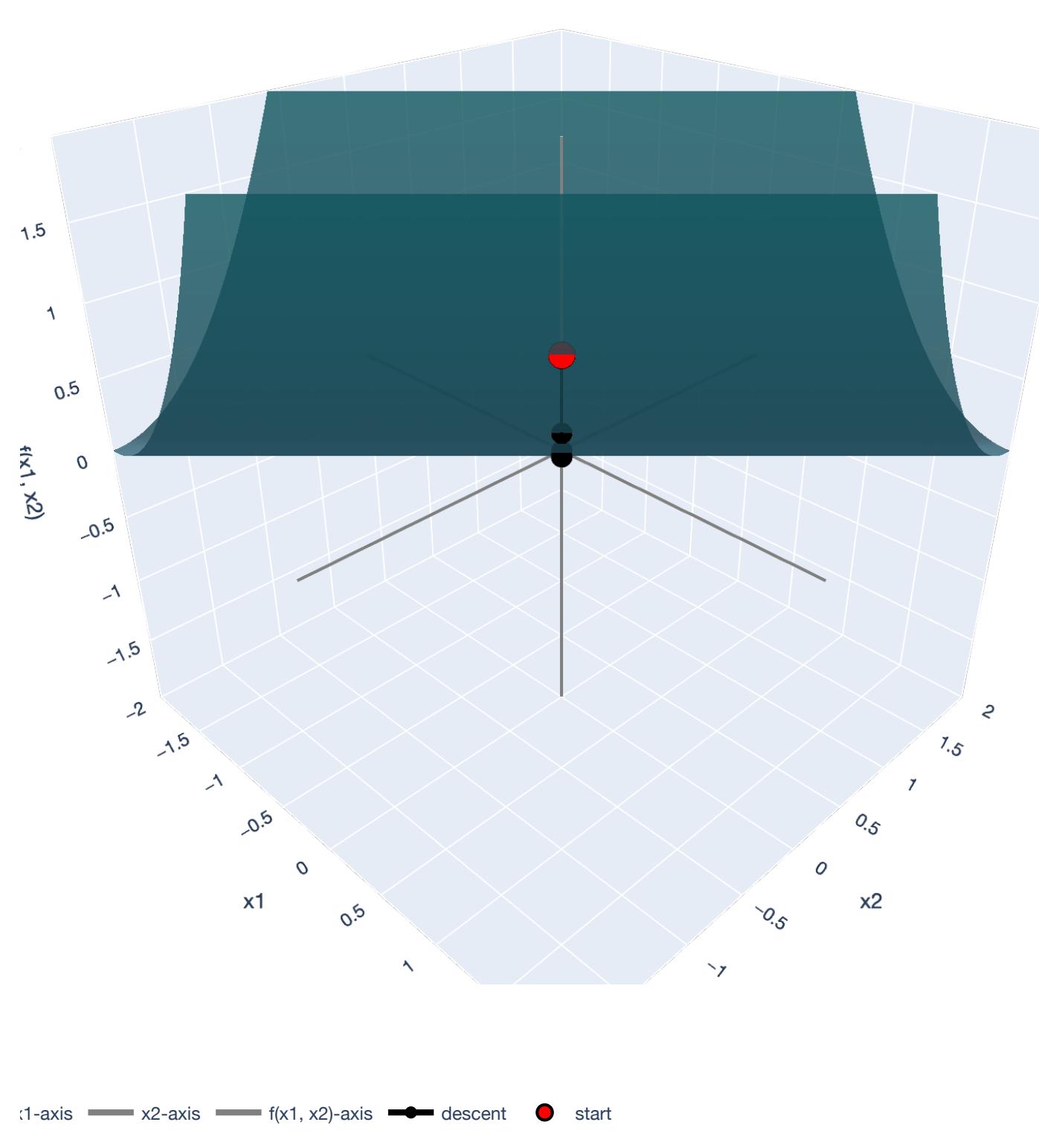
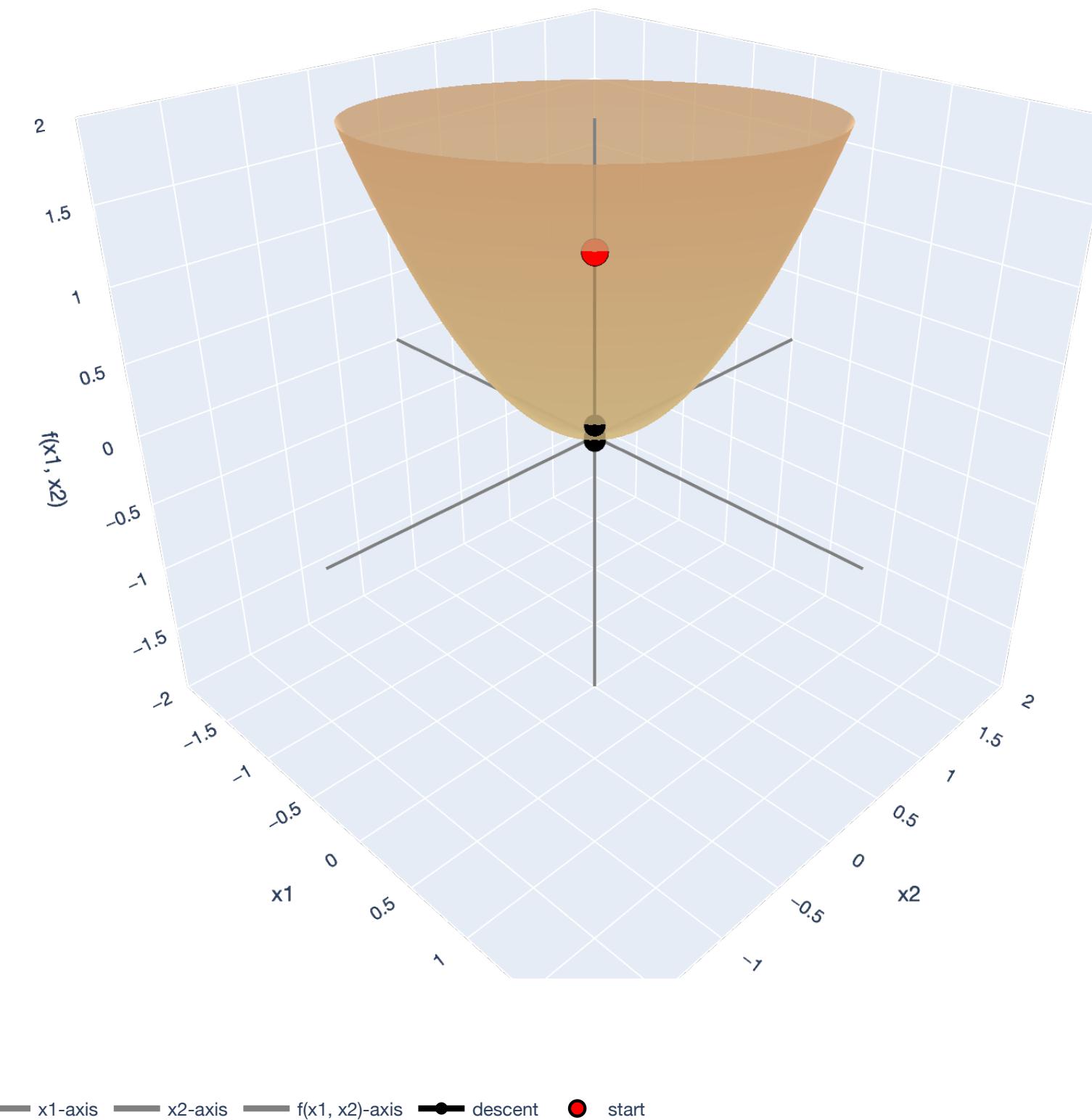
# Lesson Overview

## Big Picture: Gradient Descent



# Lesson Overview

## Big Picture: Gradient Descent



# Quadratic Forms

## 2D Example

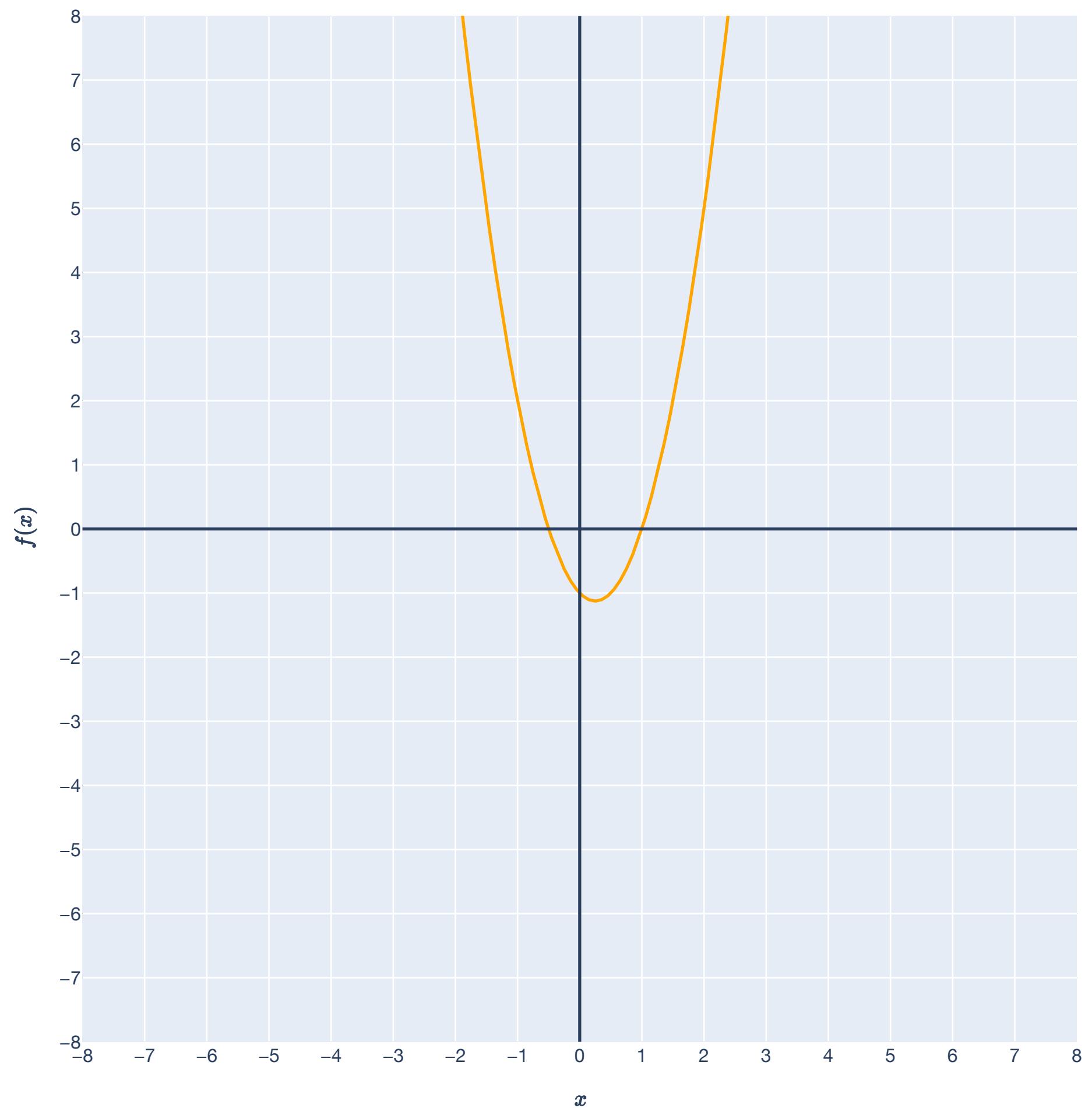
A *quadratic function*  $f: \mathbb{R} \rightarrow \mathbb{R}$  has the form

$$f(x) = ax^2 + bx + c,$$

where  $a, b, c \in \mathbb{R}$ .

**Example:**  $f(x) = 2x^2 - x - 1$

$$f(x) = 2x^2 - x - 1$$



# Quadratic Forms

## 2D Example

A *quadratic function*  $f: \mathbb{R} \rightarrow \mathbb{R}$  has the form

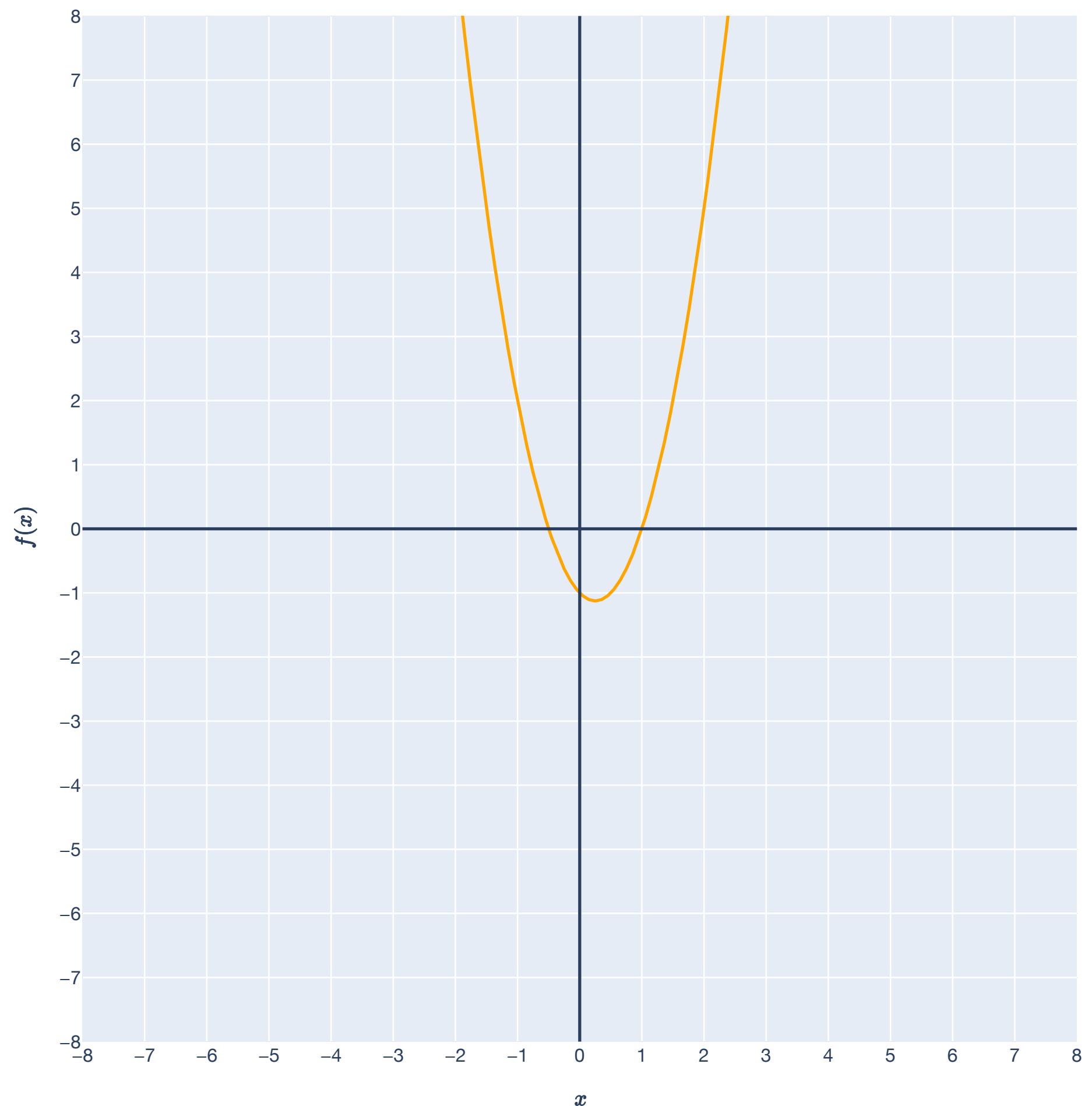
$$f(x) = ax^2 + bx + c,$$

where  $a, b, c \in \mathbb{R}$  are constants.

**Example:**  $f(x) = 2x^2 - x - 1$

We will be concerned about finding *minima* of quadratic functions.

$$f(x) = 2x^2 - x - 1$$



# Quadratic Forms

## 3D Example

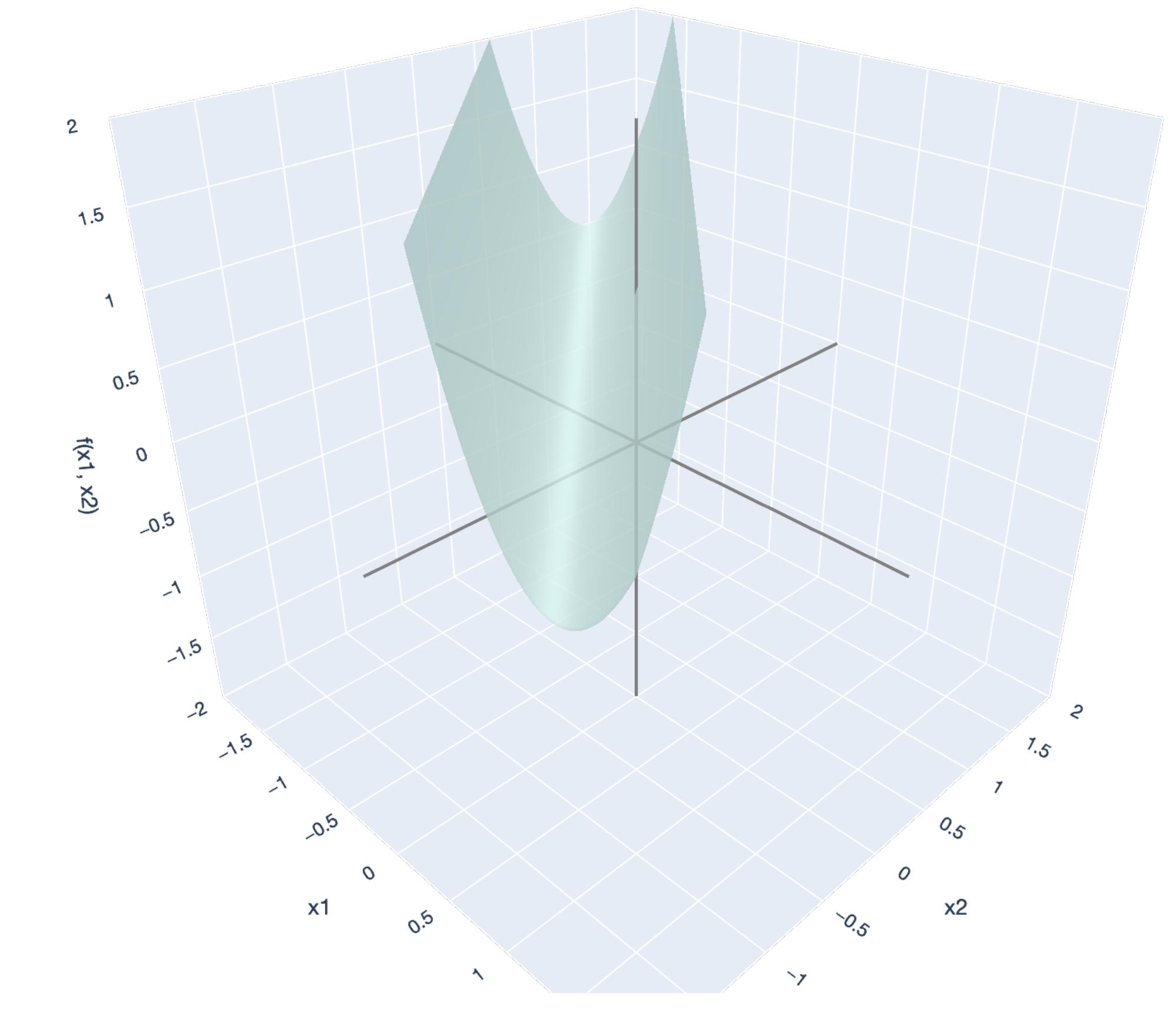
In 3D, a *quadratic function*  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  has the form

$$f(x) = ax^2 + 2bxy + cy^2 + dx + ey + f,$$

where  $a, b, c, d, e, f \in \mathbb{R}$  are all constants.

### Example:

$$f(x) = 2x^2 + 4xy + 2y^2 + 2x + 2y + 1$$

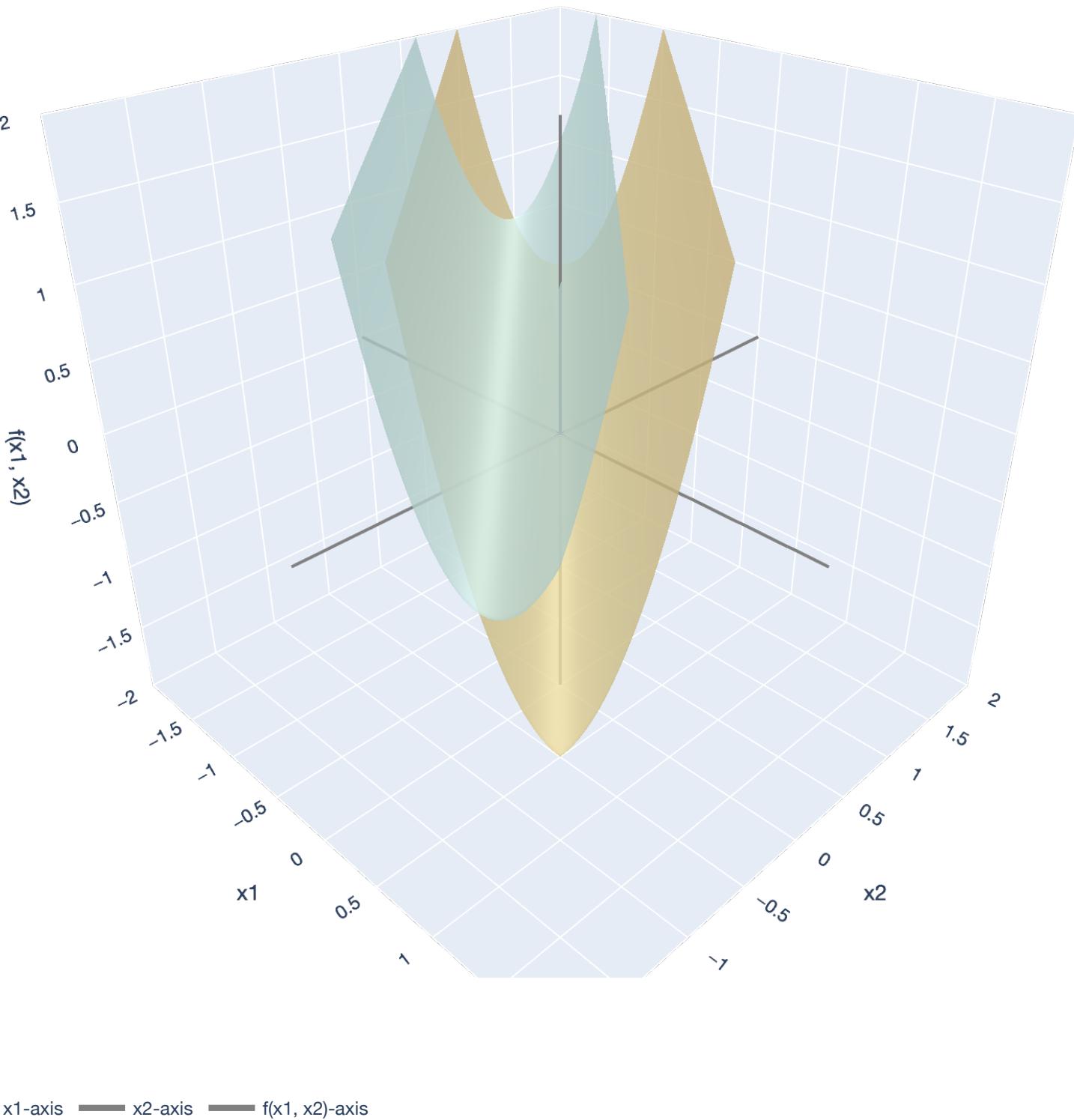


— x1-axis — x2-axis — f(x1, x2)-axis

# Quadratic Forms

## 3D Example

$$f(x) = 2x^2 + 4xy + 2y^2 + 2x + 2y + 1 \text{ vs. } f(x) = 2x^2 + 4xy + 2y^2$$



# Quadratic Forms

## 3D Example

In 3D, a *quadratic function*  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  has the form

$$f(x) = \underbrace{ax^2 + 2bxy + cy^2}_{\text{quadratic}} + \underbrace{dx + ey}_{\text{linear}} + \underbrace{f}_{\text{constant}}.$$

Let's only examine the quadratic part!

$$f(x) = ax^2 + 2bxy + cy^2.$$

# Quadratic Forms

## Relationship with matrices and eigenvalues

A function  $f: \mathbb{R}^2 \rightarrow \mathbb{R}$  is a **quadratic form** if it is a polynomial with terms of all degree two:

$$f(x) = ax^2 + 2bxy + cy^2.$$

We can rewrite this in matrix form:

$$f(x, y) = [x \ y] \begin{bmatrix} a & b \\ b & c \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$$

# Quadratic Forms

## Relationship with matrices and eigenvalues

Consider a quadratic form:

$$f(x, y) = [x \ y] \begin{bmatrix} a & b \\ b & c \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$$

The matrix  $\mathbf{A} \in \mathbb{R}^{2 \times 2}$  is always symmetric, so it is diagonalizable!

$\mathbf{A} = \mathbf{Q} \Lambda \mathbf{Q}^T$ , where  $\Lambda \in \mathbb{R}^{d \times d}$  is diagonal.

# Quadratic Forms

## Relationship with matrices and eigenvalues

The matrix  $A \in \mathbb{R}^{2 \times 2}$  is always symmetric, so it is diagonalizable!

$$A = Q\Lambda Q^T, \text{ where } \Lambda \in \mathbb{R}^{d \times d} \text{ is diagonal.}$$

$$\implies f(\mathbf{x}) = \mathbf{x}^T A \mathbf{x} = \mathbf{x}^T Q \Lambda Q^T \mathbf{x}$$

$$\implies \bar{\mathbf{x}}^T \Lambda \bar{\mathbf{x}}, \text{ where } \bar{\mathbf{x}} = Q^T \mathbf{x}.$$

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$

# Quadratic Forms

## Relationship with matrices and eigenvalues

$\mathbf{A} = \mathbf{Q}\Lambda\mathbf{Q}^\top$ , where  $\Lambda \in \mathbb{R}^{d \times d}$  is diagonal.

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$$

There are three possibilities:

1.  $\lambda_1$  and  $\lambda_2$  are *both* positive (*positive definite*).
2.  $\lambda_1$  or  $\lambda_2$  is zero, and the other is positive (*positive semidefinite*).
3.  $\lambda_1$  or  $\lambda_2$  is negative (*indefinite*).

# Quadratic Forms

Example: positive definite

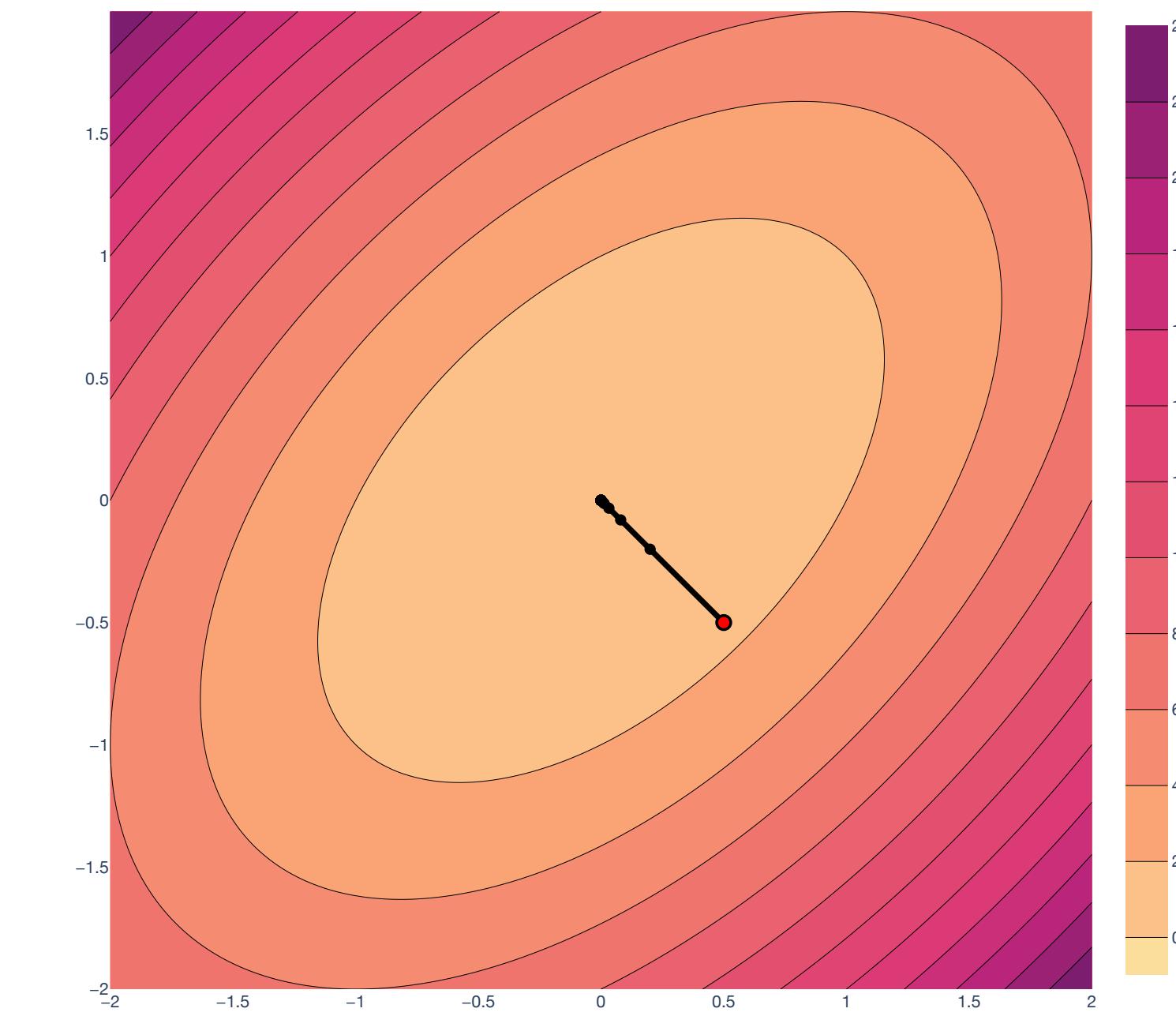
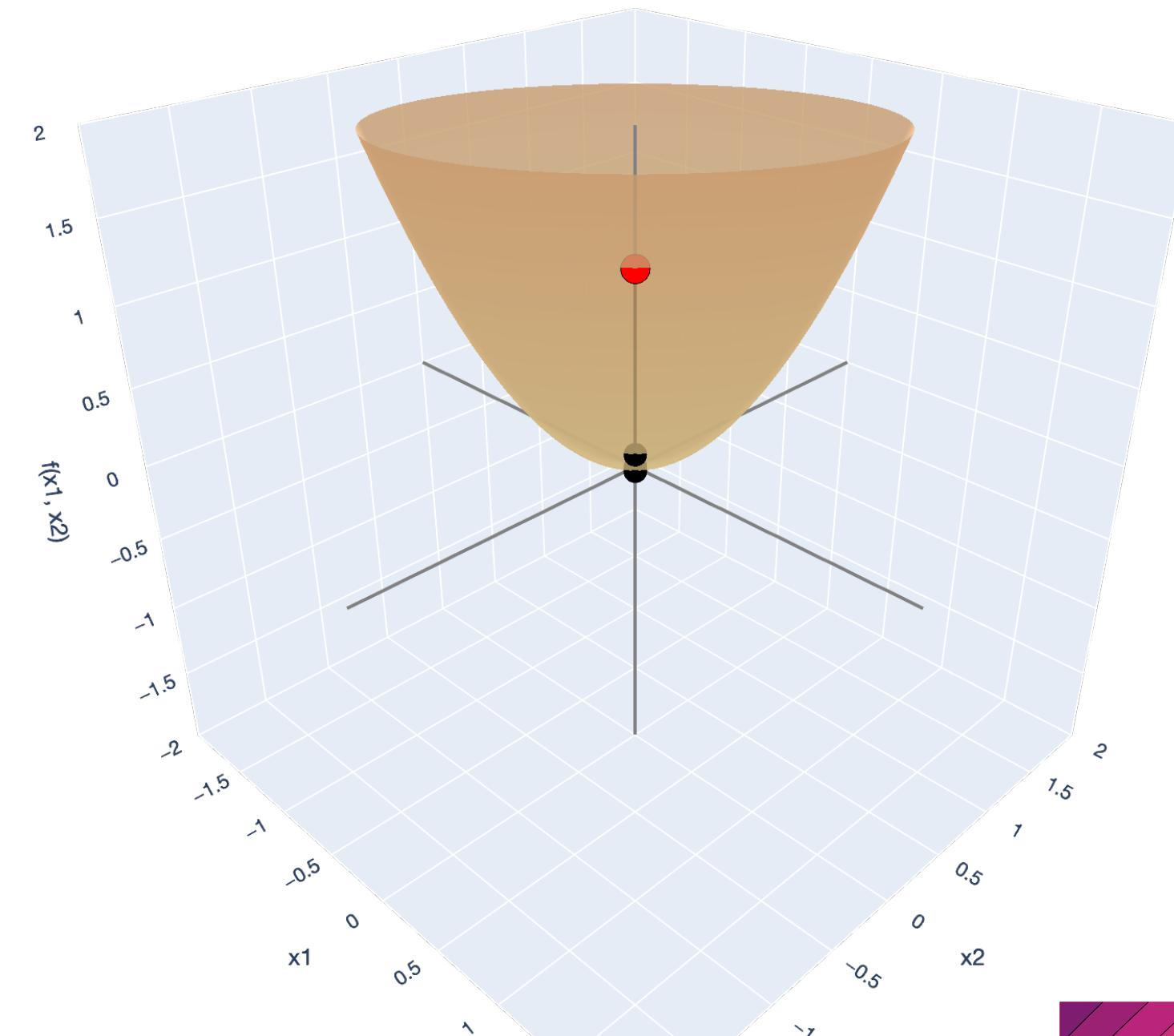
Example:

$$f(x, y) = [x \ y] \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

Eigendecomposition:

$$\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} -1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

$$\text{so } \Lambda = \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix}.$$



# Quadratic Forms

## Example: positive semidefinite

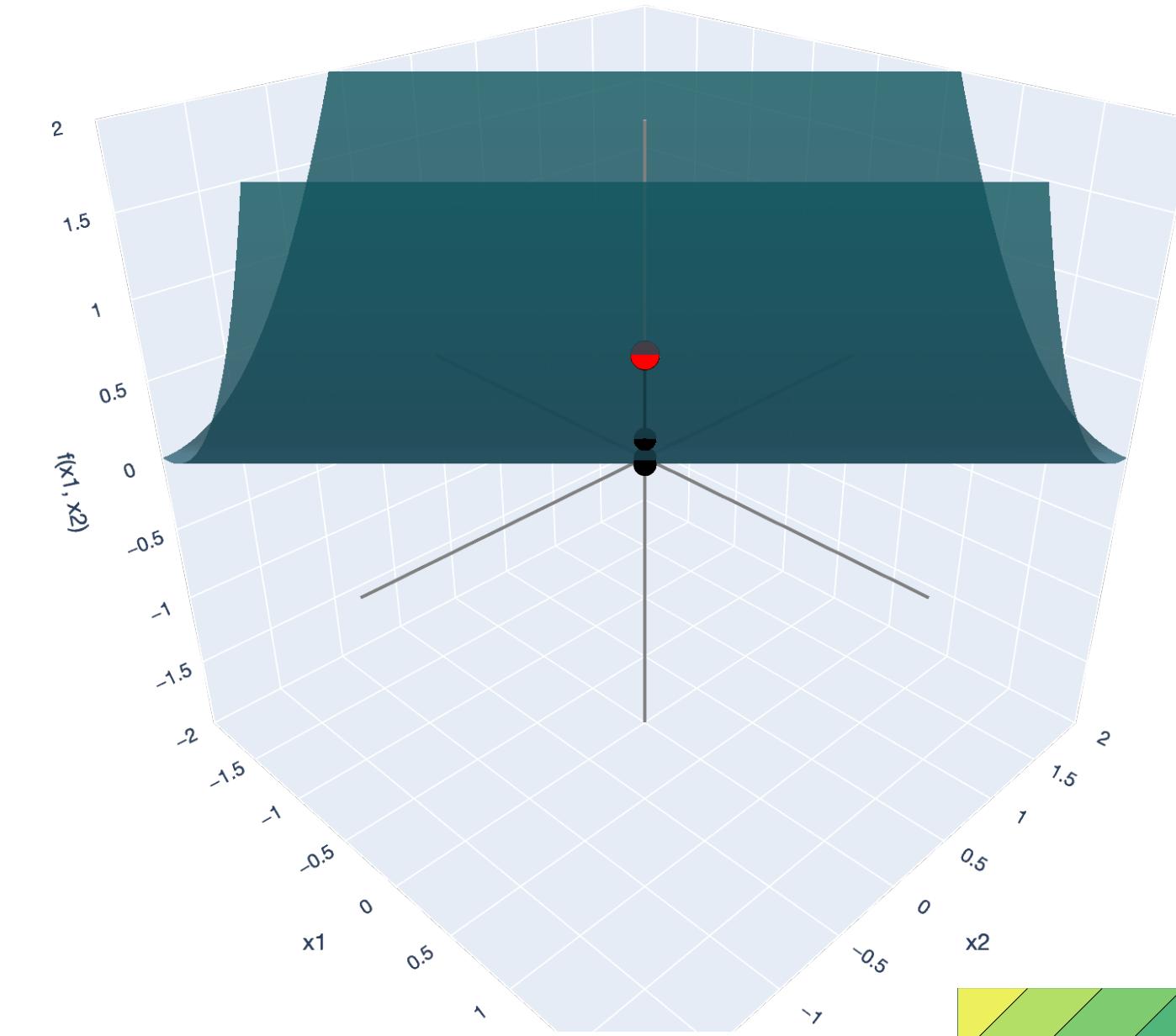
Example:

$$f(x, y) = [x \ y] \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

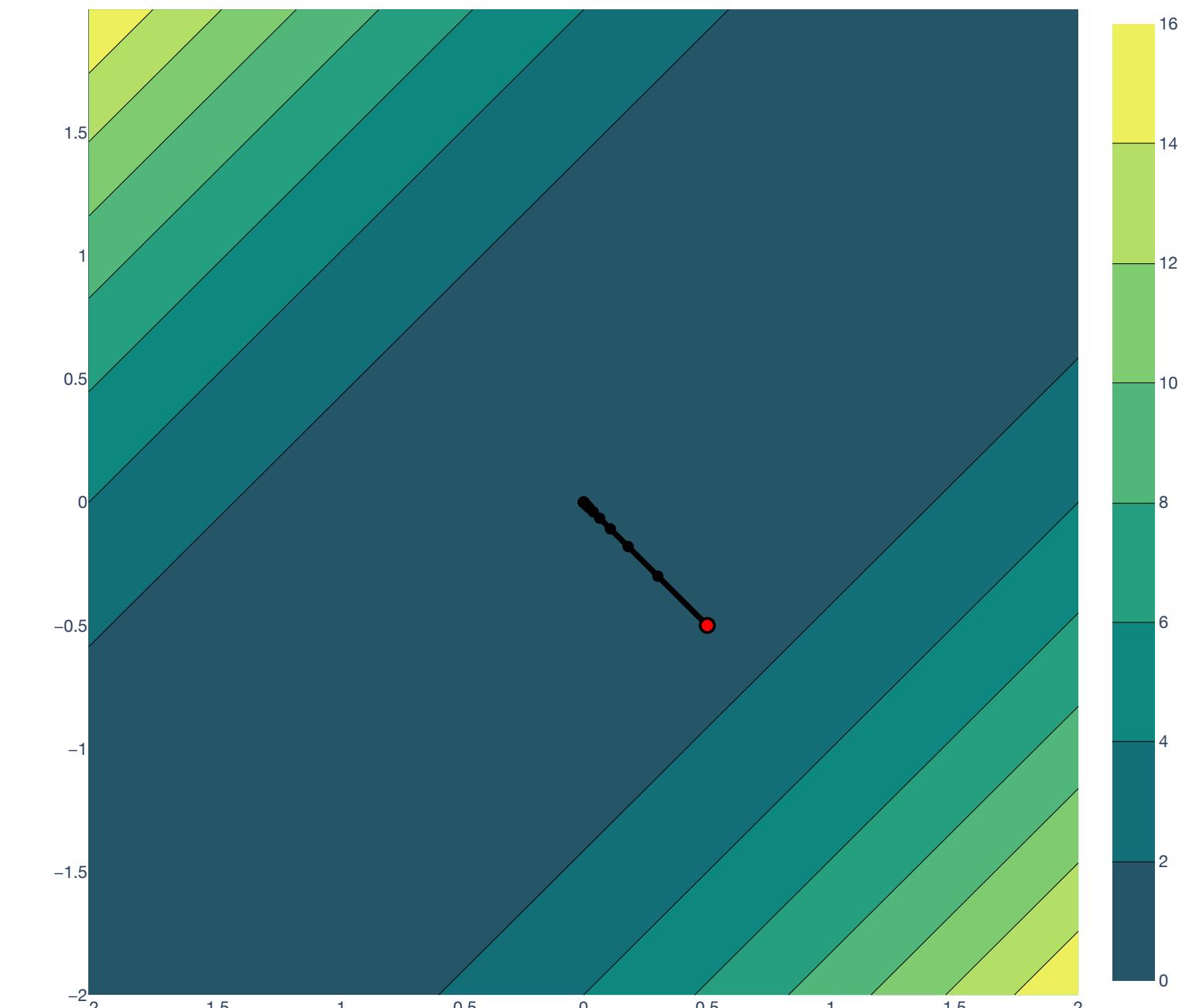
Eigendecomposition:

$$\begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} = \begin{bmatrix} -1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} -1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

$$\text{so } \Lambda = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}.$$



— x1-axis — x2-axis —  $f(x_1, x_2)$ -axis • descent ● start



• descent ● start

# Quadratic Forms

## Example: indefinite

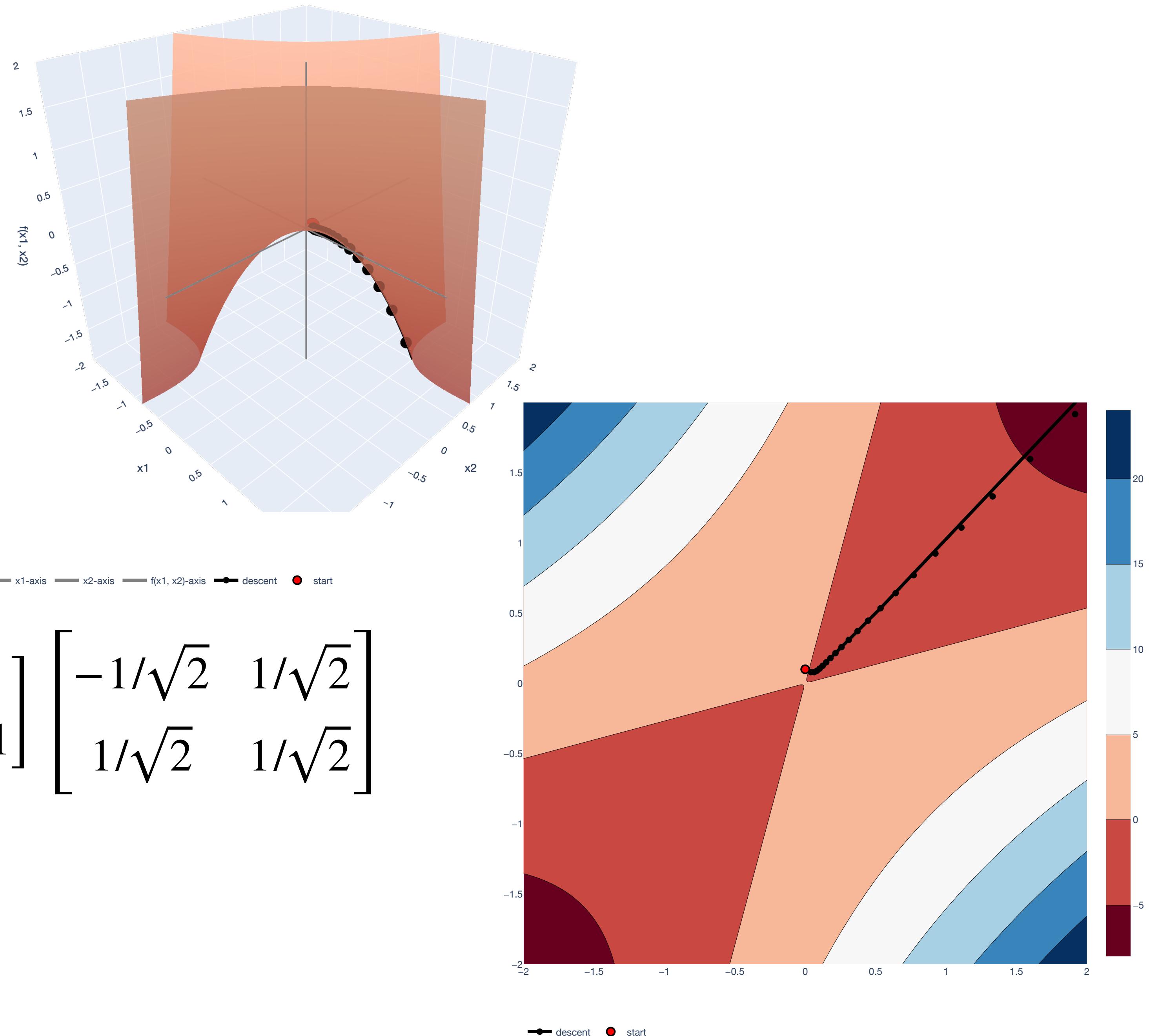
Example:

$$f(x, y) = [x \ y] \begin{bmatrix} 1 & -2 \\ -2 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

Eigendecomposition:

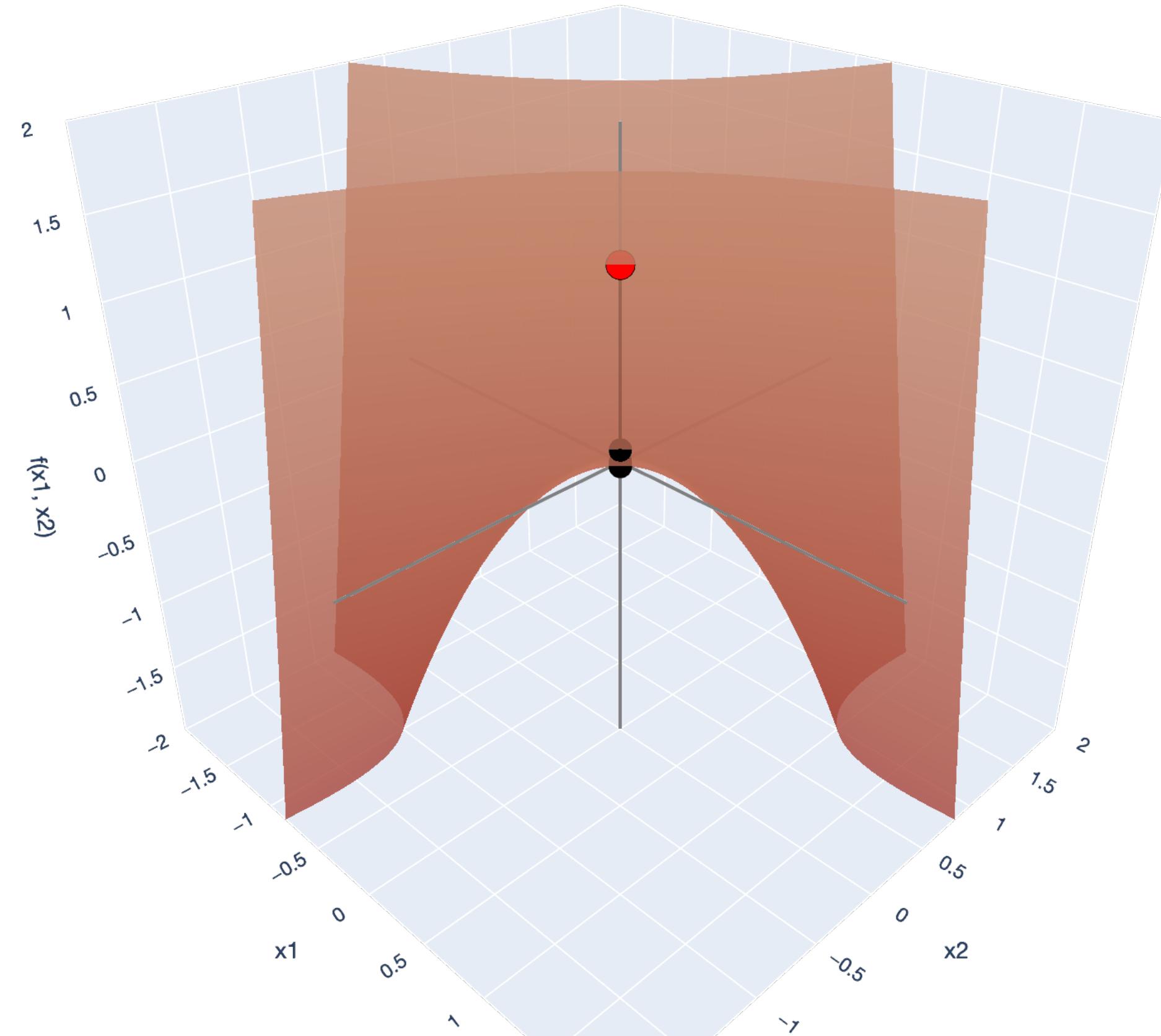
$$\begin{bmatrix} 1 & -2 \\ -2 & 1 \end{bmatrix} = \begin{bmatrix} -1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 3 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} -1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

$$\text{so } \Lambda = \begin{bmatrix} 3 & 0 \\ 0 & -1 \end{bmatrix}.$$

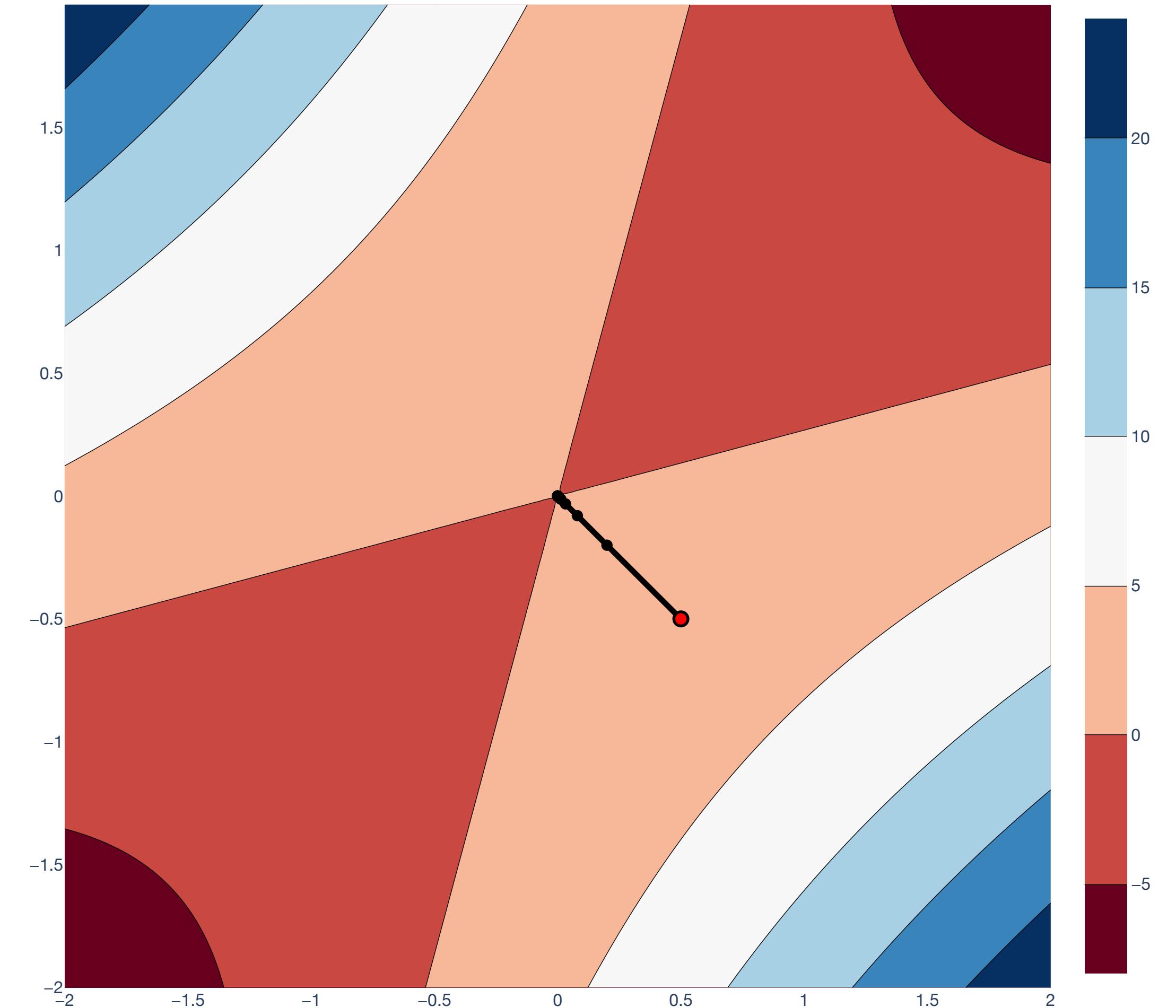


# Quadratic Forms

## Example: indefinite



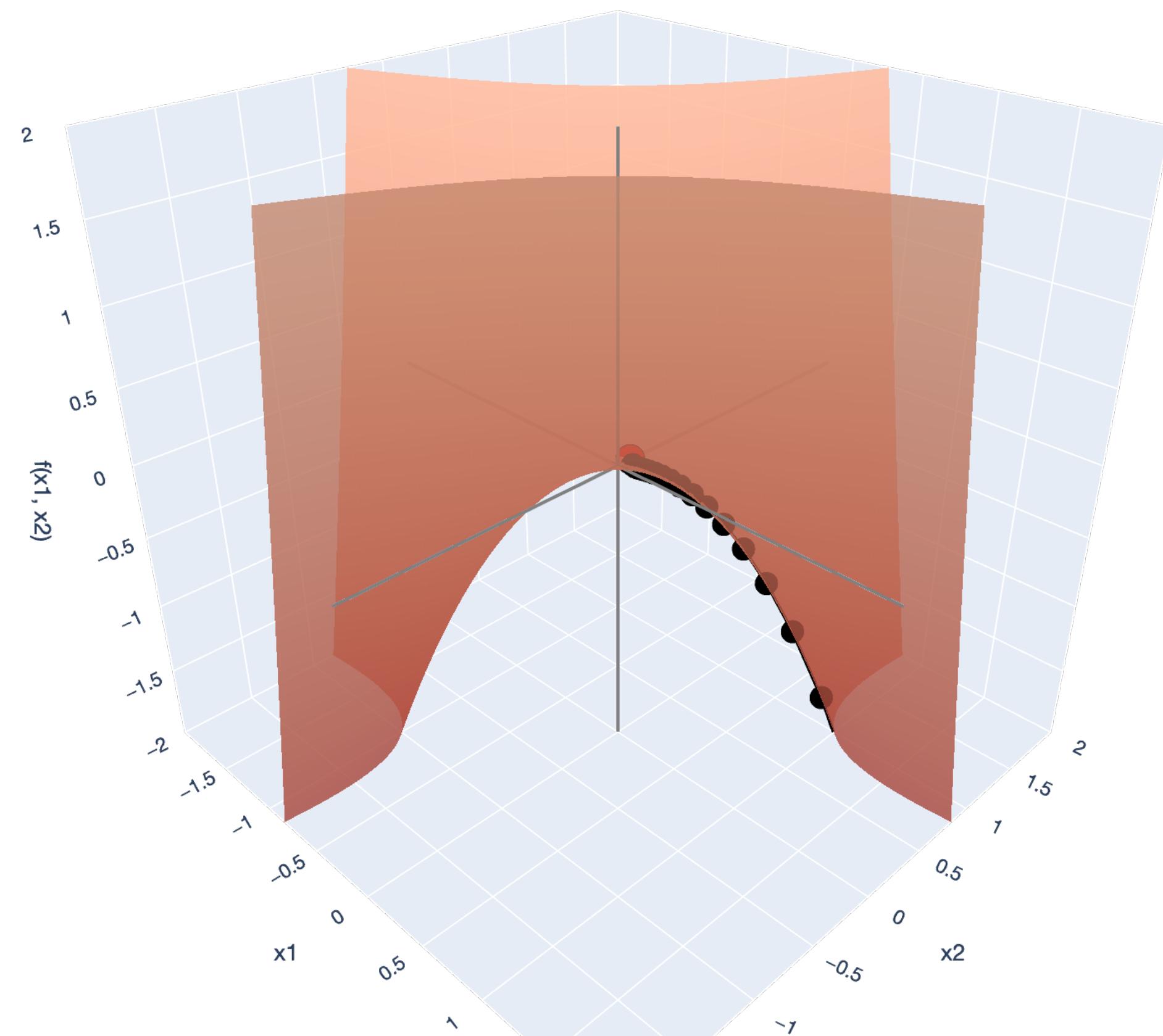
— x1-axis — x2-axis — f( $x_1, x_2$ )-axis ● descent ● start



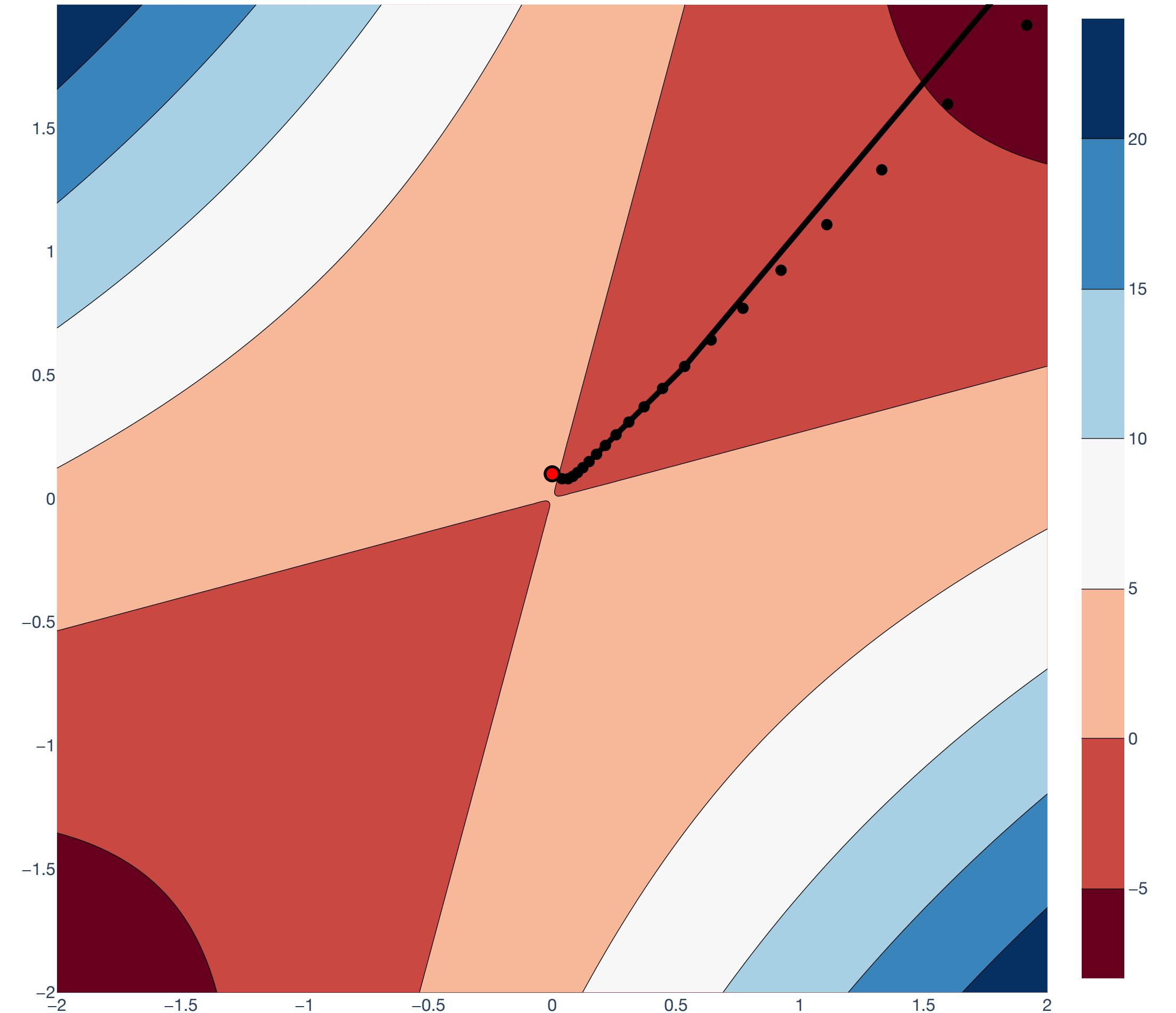
● descent ● start

# Quadratic Forms

## Example: indefinite



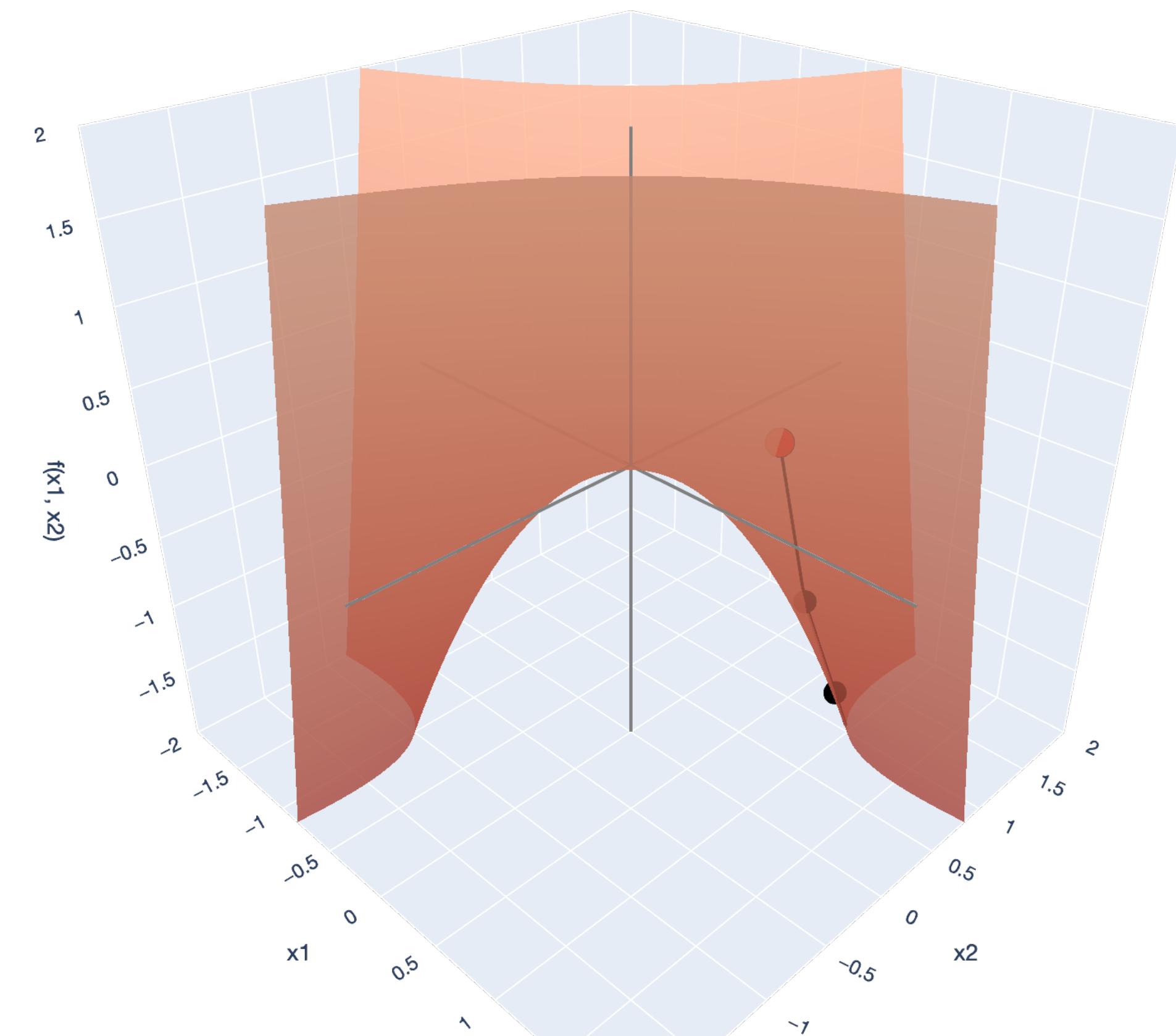
—  $x_1$ -axis —  $x_2$ -axis —  $f(x_1, x_2)$ -axis ● descent ● start



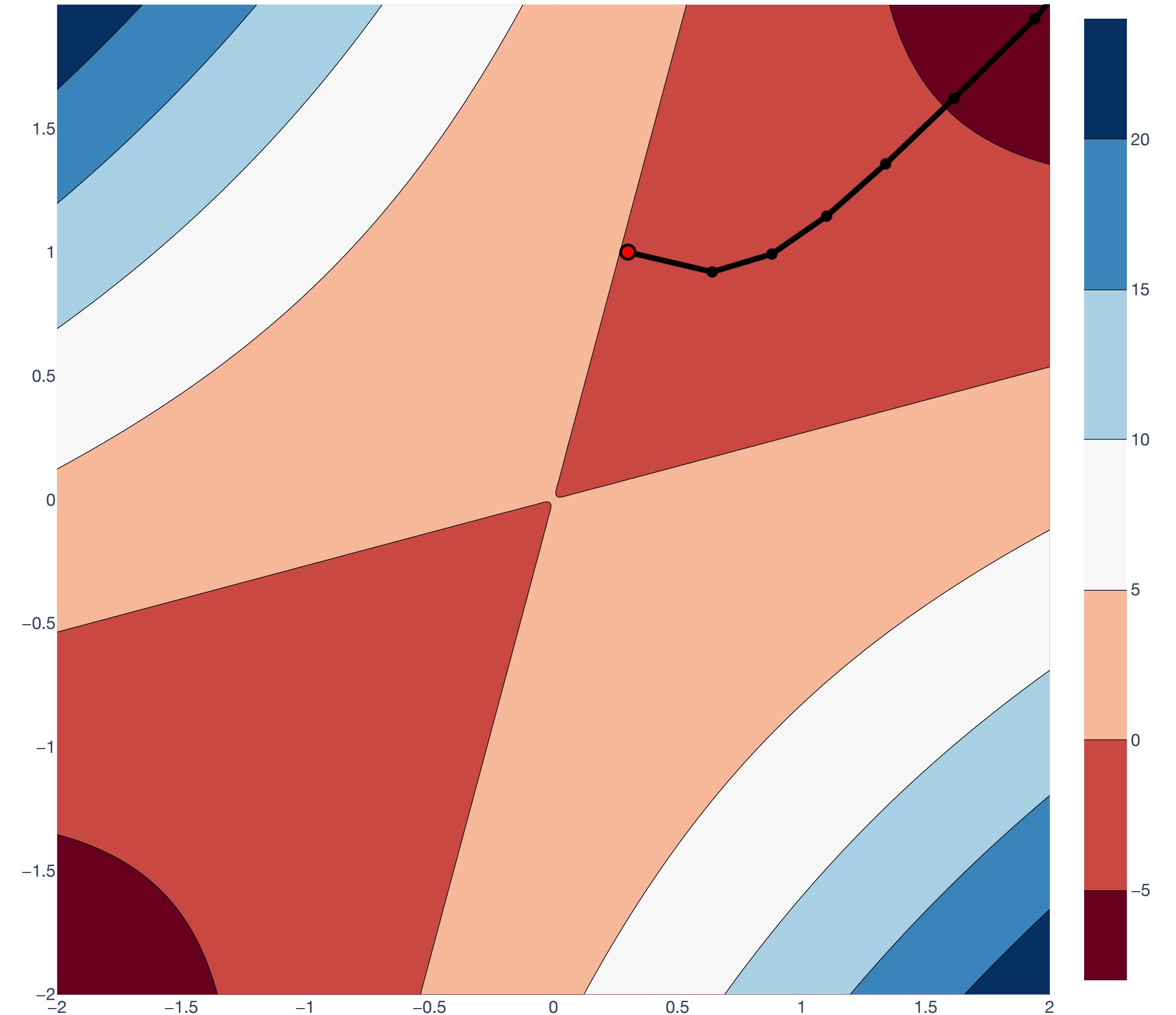
● descent ● start

# Quadratic Forms

## Example: indefinite



— x1-axis — x2-axis —  $f(x_1, x_2)$ -axis ● descent ● start



● descent ● start

# Least Squares

## Example of quadratic form

Consider the familiar function we've been thinking about:

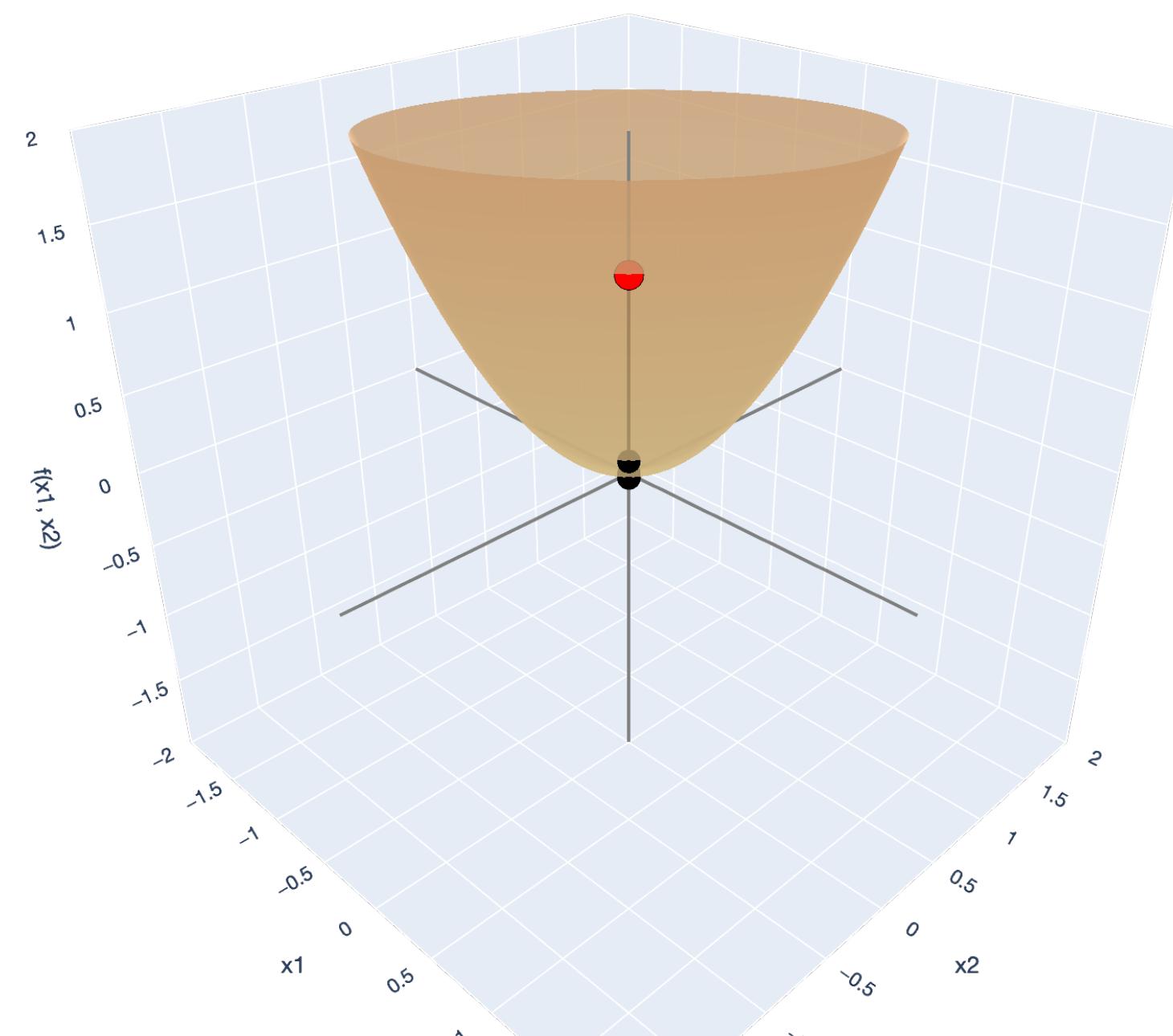
$$f(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2$$

$$(\mathbf{X}\mathbf{w} - \mathbf{y})^\top(\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{w}^\top(\mathbf{X}^\top\mathbf{X})\mathbf{w} - 2\mathbf{w}^\top(\mathbf{X}^\top\mathbf{y}) + \mathbf{y}^\top\mathbf{y}.$$

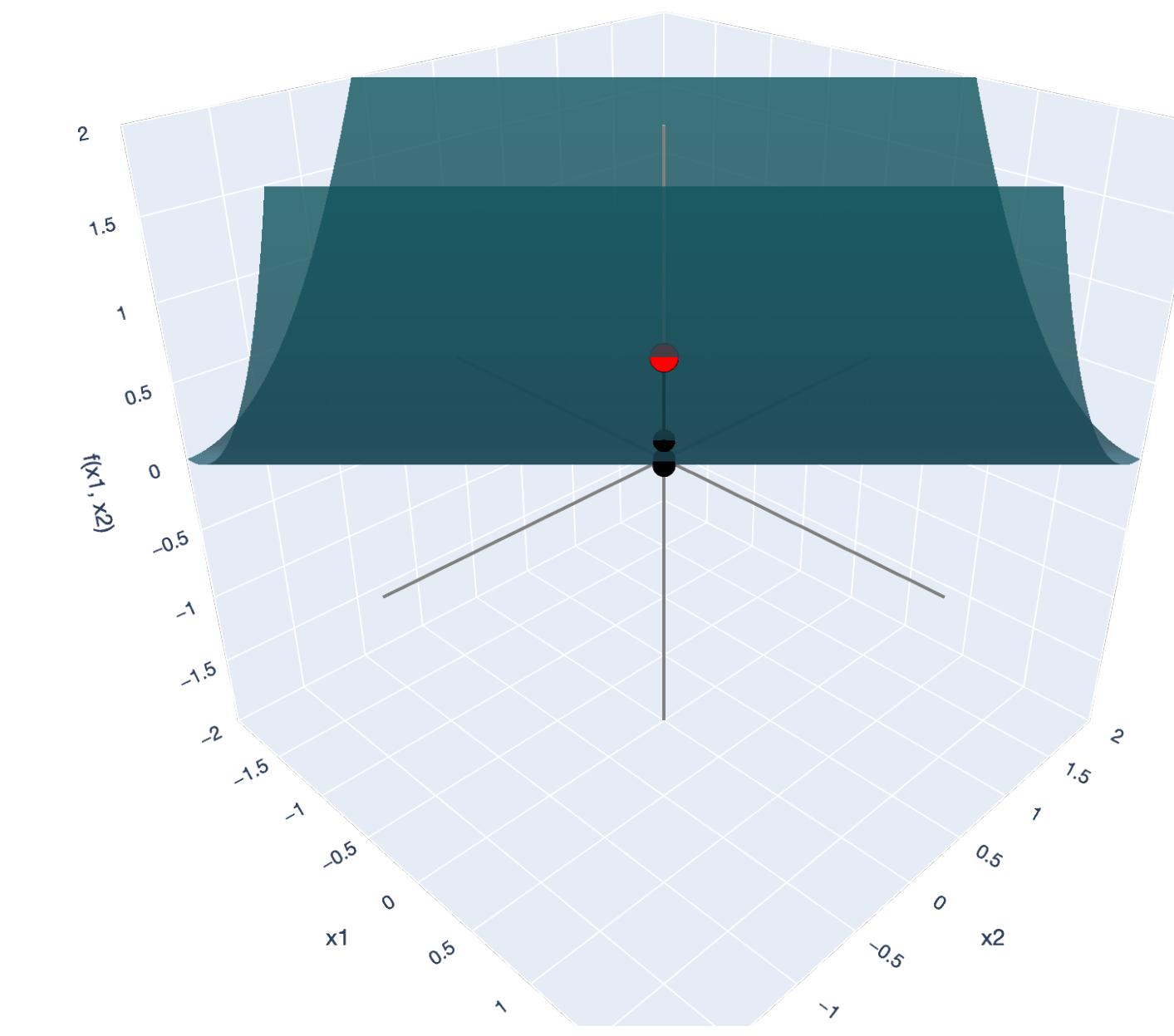
The quadratic form  $\mathbf{w}^\top(\mathbf{X}^\top\mathbf{X})\mathbf{w}$  is positive semidefinite!

# Gradient Descent

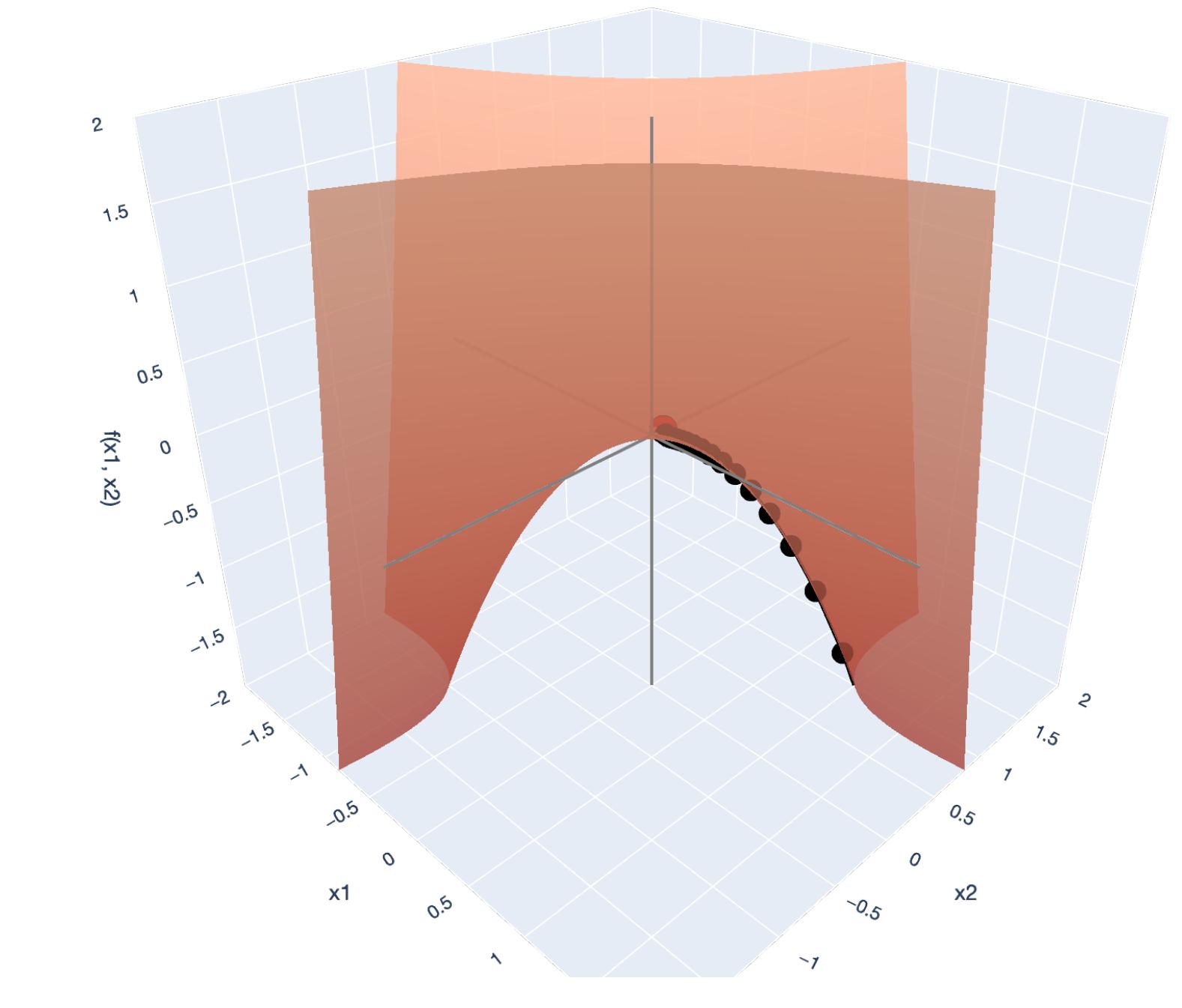
## Preview



— x1-axis — x2-axis —  $f(x_1, x_2)$ -axis ● descent ● start



— x1-axis — x2-axis —  $f(x_1, x_2)$ -axis ● descent ● start



— x1-axis — x2-axis —  $f(x_1, x_2)$ -axis ● descent ● start

$$\Lambda = \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\Lambda = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix}$$

$$\Lambda = \begin{bmatrix} 3 & 0 \\ 0 & -1 \end{bmatrix}$$

# Recap

# Lesson Overview

**Linear dynamical systems example.** Motivation for eigendecomposition as a way to make repeated matrix multiplication easier.

**Eigendecomposition.** Definition of eigenvectors, eigenvalues.

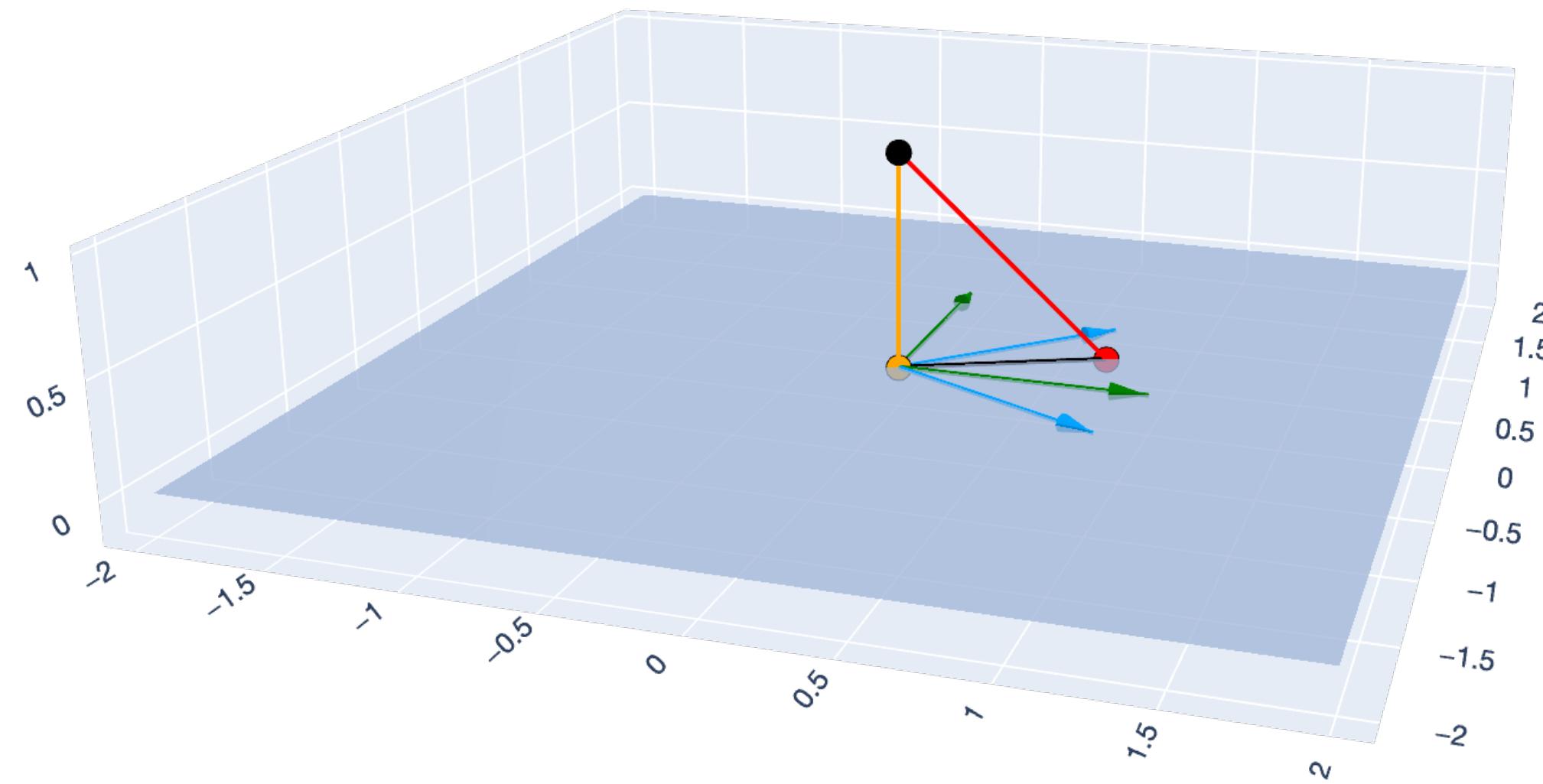
**Eigendecomposition and SVD.** The eigendecomposition drops out of the SVD.

**Spectral Theorem.** Symmetric matrices are always diagonalizable.

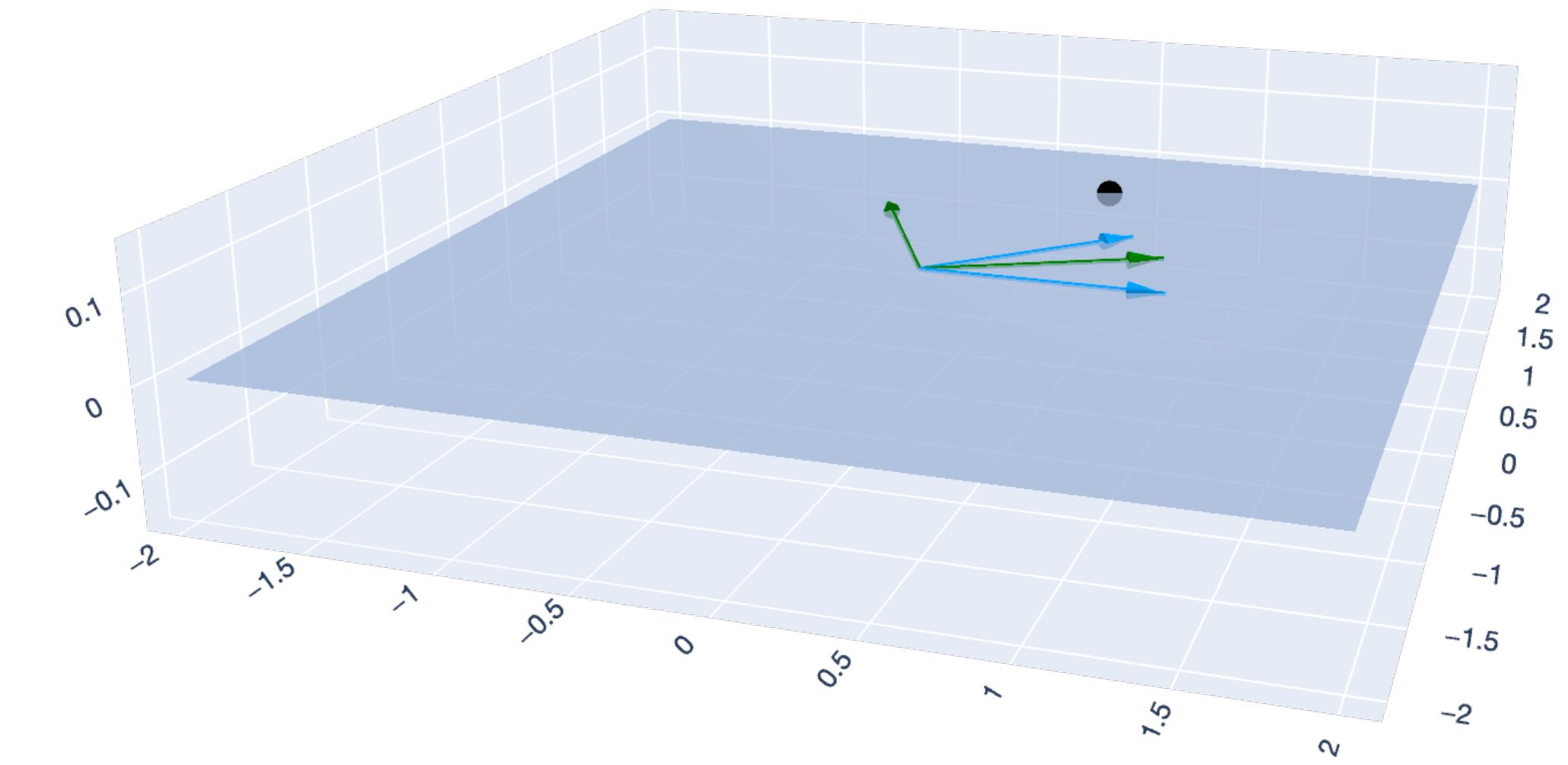
**Positive semidefinite matrices/positive definite matrices.** Definition and some visual examples through the corresponding quadratic forms.

# Lesson Overview

## Big Picture: Least Squares



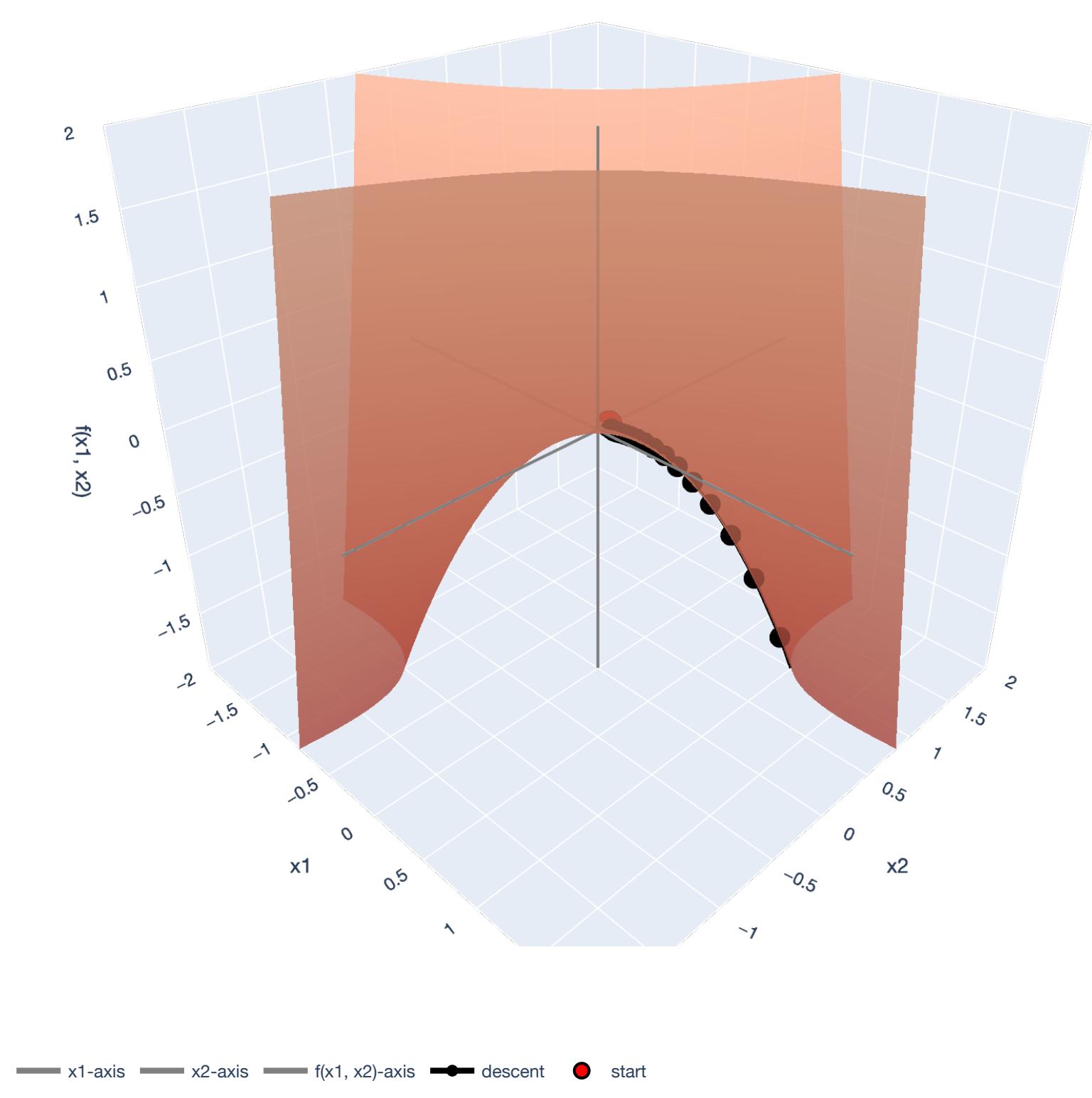
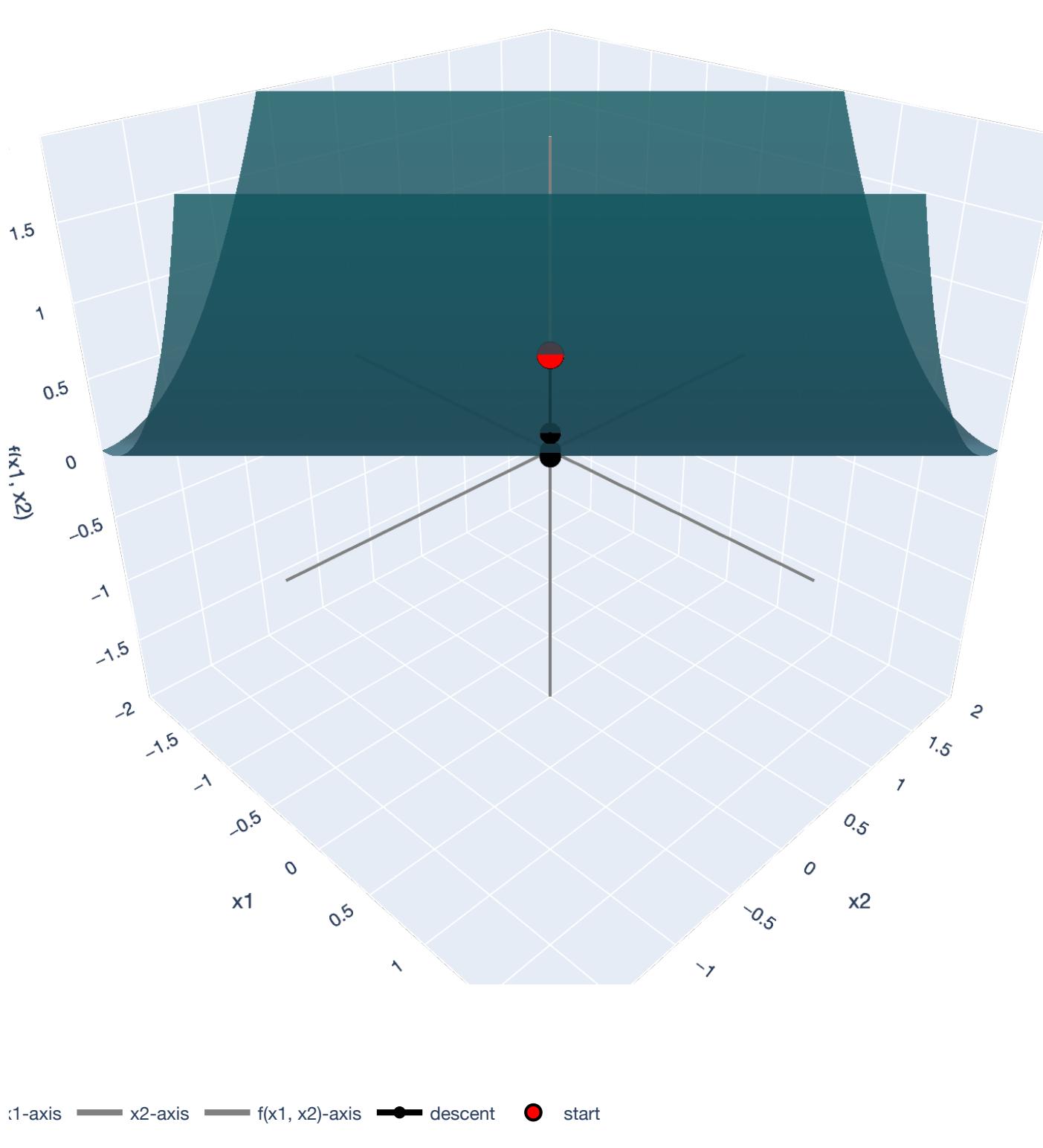
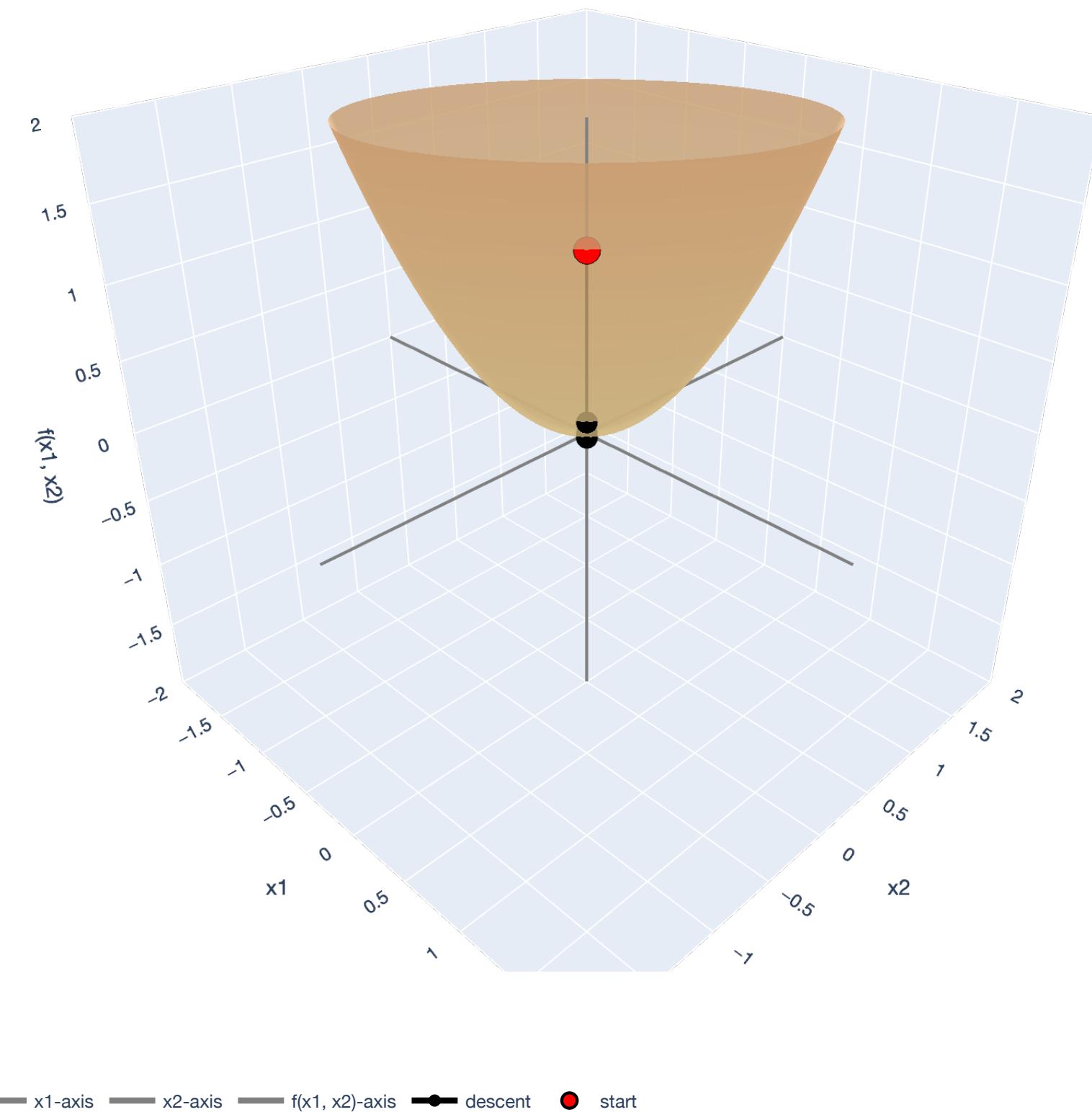
—  $x_1$  —  $x_2$  —  $u_1$  —  $u_2$  —  $y - \hat{y}$  —  $\hat{y} - y$  ●  $y$  ○  $\hat{y}$  ●  $\sim y$



—  $x_1$  —  $x_2$  —  $u_1$  —  $u_2$  ●  $y$

# Lesson Overview

## Big Picture: Gradient Descent



# References

*Mathematics for Machine Learning.* Marc Pieter Deisenroth, A. Aldo Faisal, Cheng Soon Ong.

*Vector Calculus, Linear Algebra, and Differential Forms: A Unified Approach.* John H. Hubbard and Barbara Burke Hubbard.

*Computational Linear Algebra Lecture Notes: Eigenvalues and eigenvectors.* Daniel Hsu.