

Math for Machine Learning

Week 2.2: Eigendecomposition and PSD Matrices

By: Samuel Deng

Logistics & Announcements

- HW ②: Due: July 18, next Thurs.
- HW ①: Due: July 11, tomorrow.
- OFFICE HOURS 3PM - 5PM (Zoom).
- BREAKS: & minutes = 16 min. total.

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.]

→ PCA

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

Lesson Overview

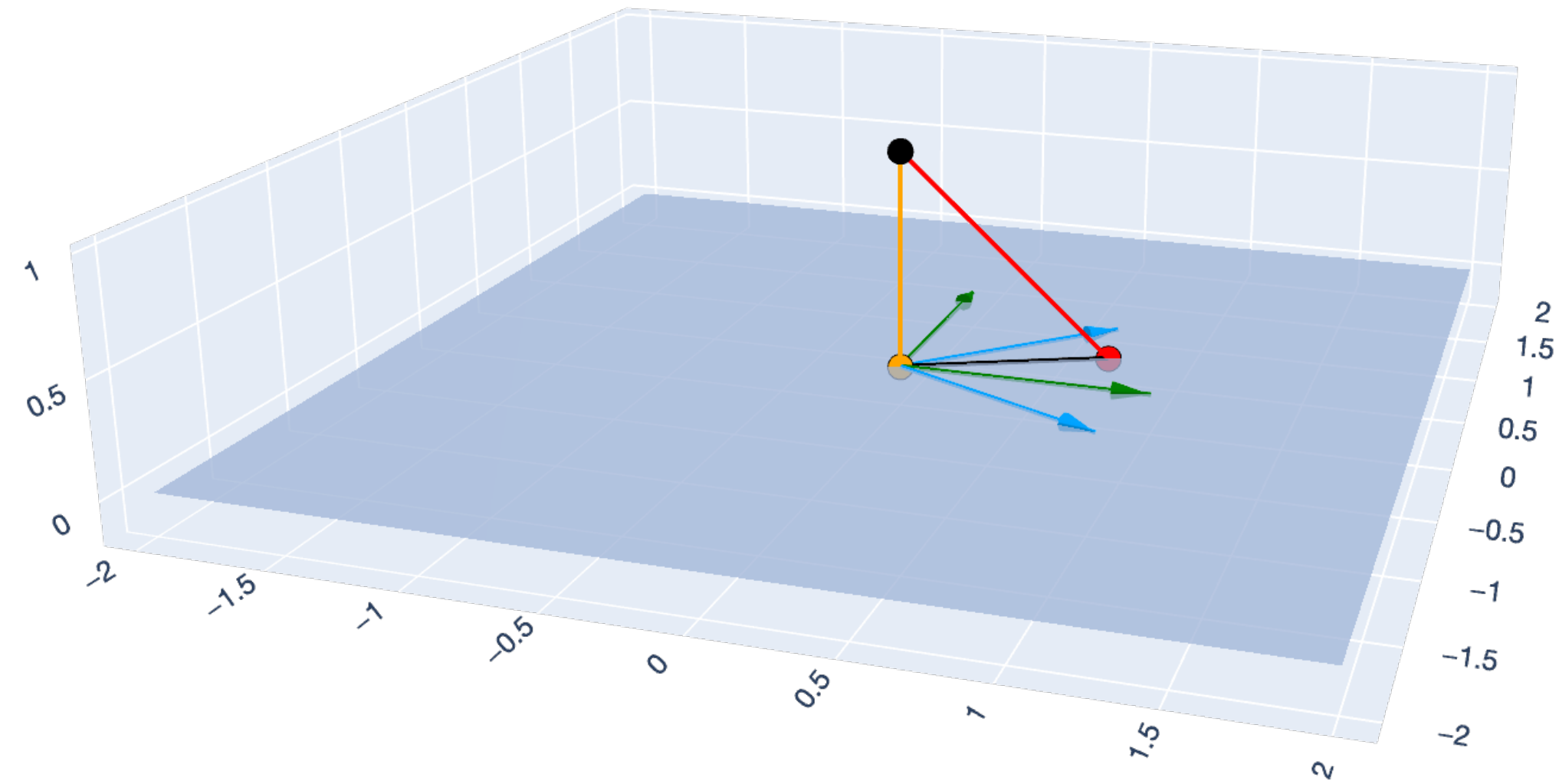
Big Picture: Least Squares

$n > d$
 $Xw \approx y$
 ↓

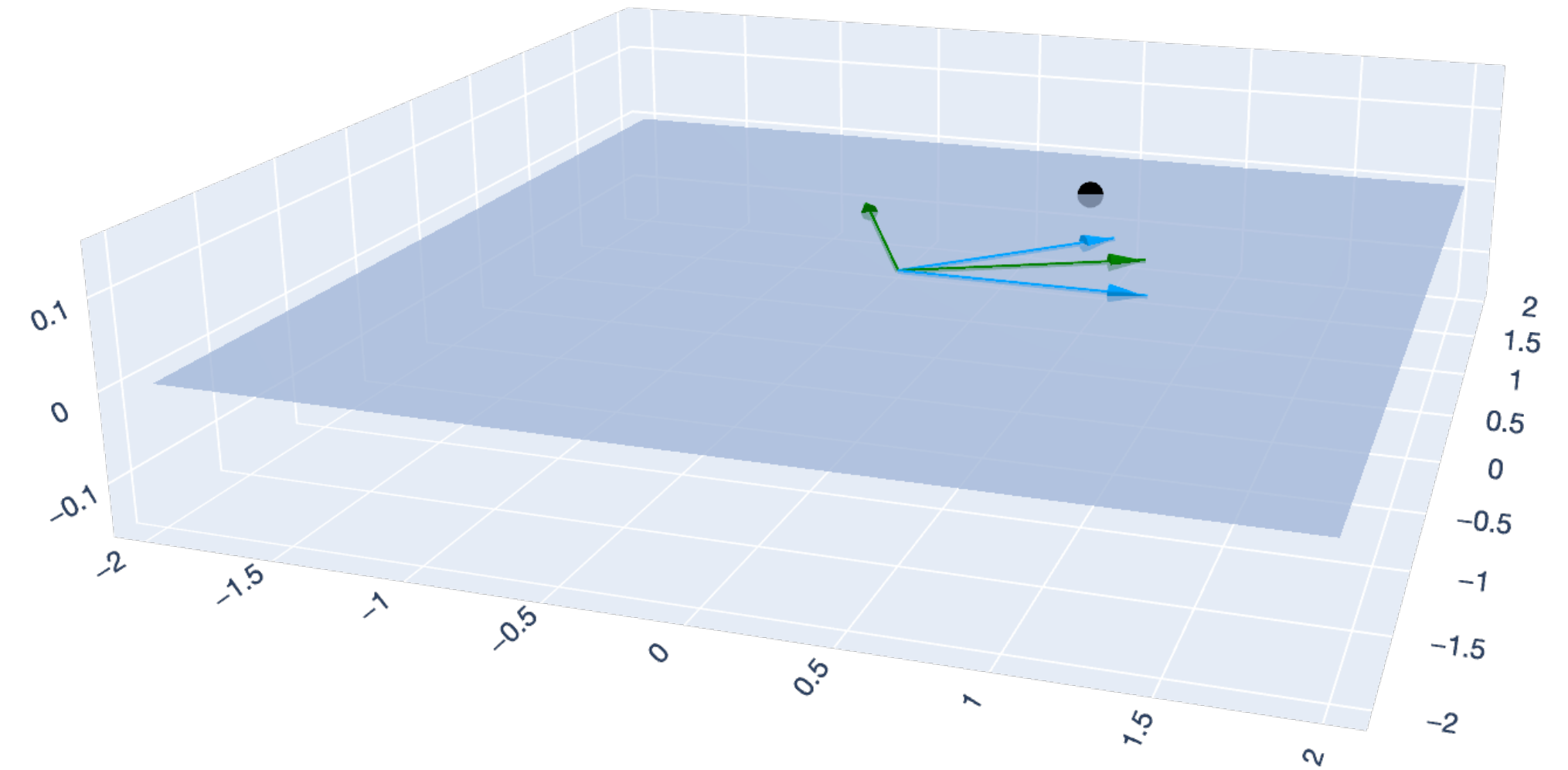
$\hat{w} = X^+ y$

$\hat{w} = X^+ y$
 $d > n$
 $Xw = y$
 ↓

more unknowns than equations.



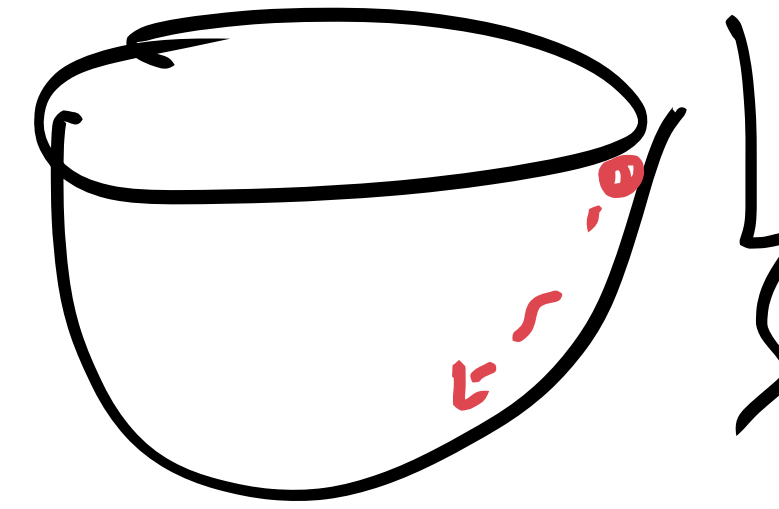
— x1 — x2 — u1 — u2 — y - ^y — ~y - ^y — ~y - y • y • ^y • ~y



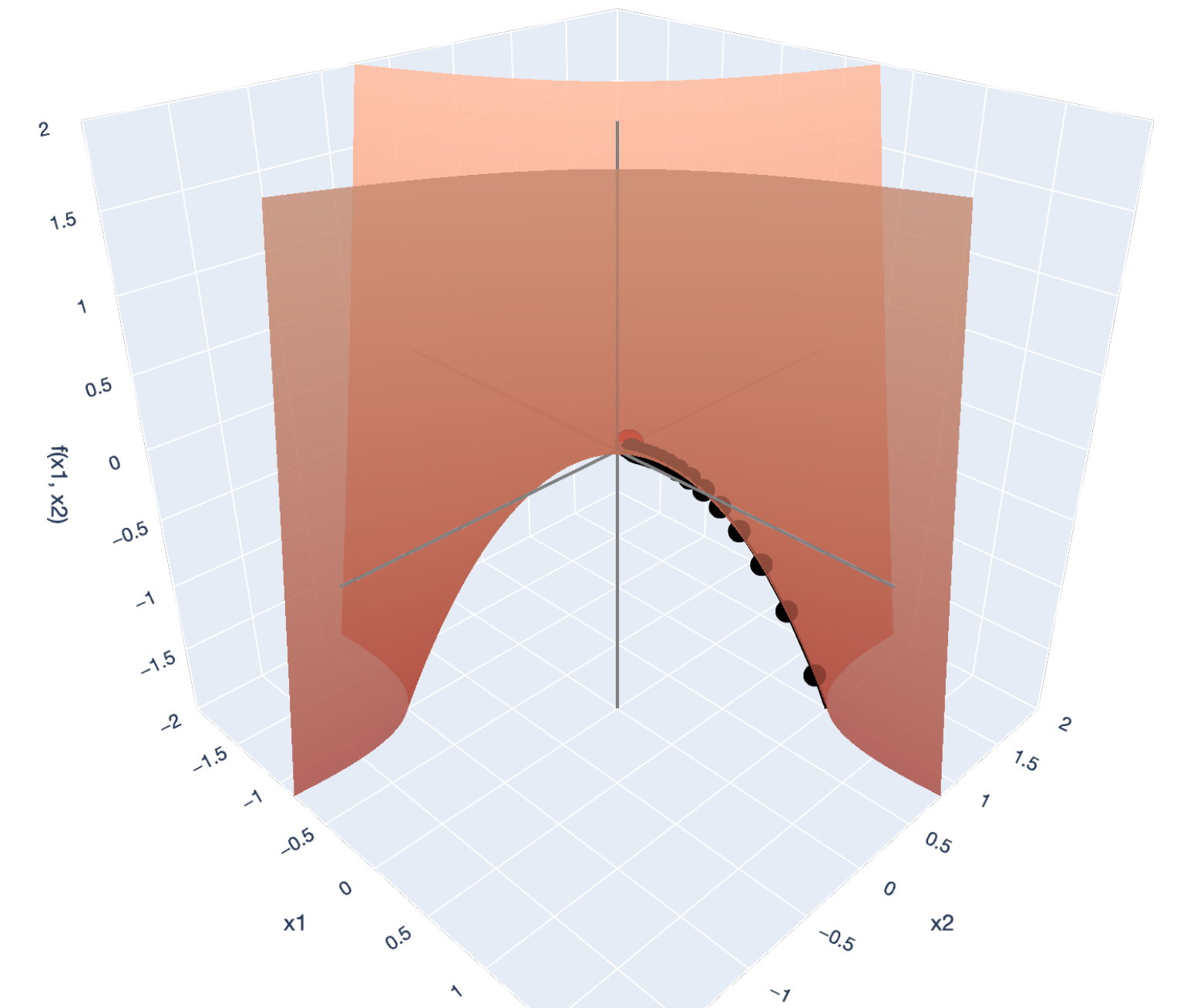
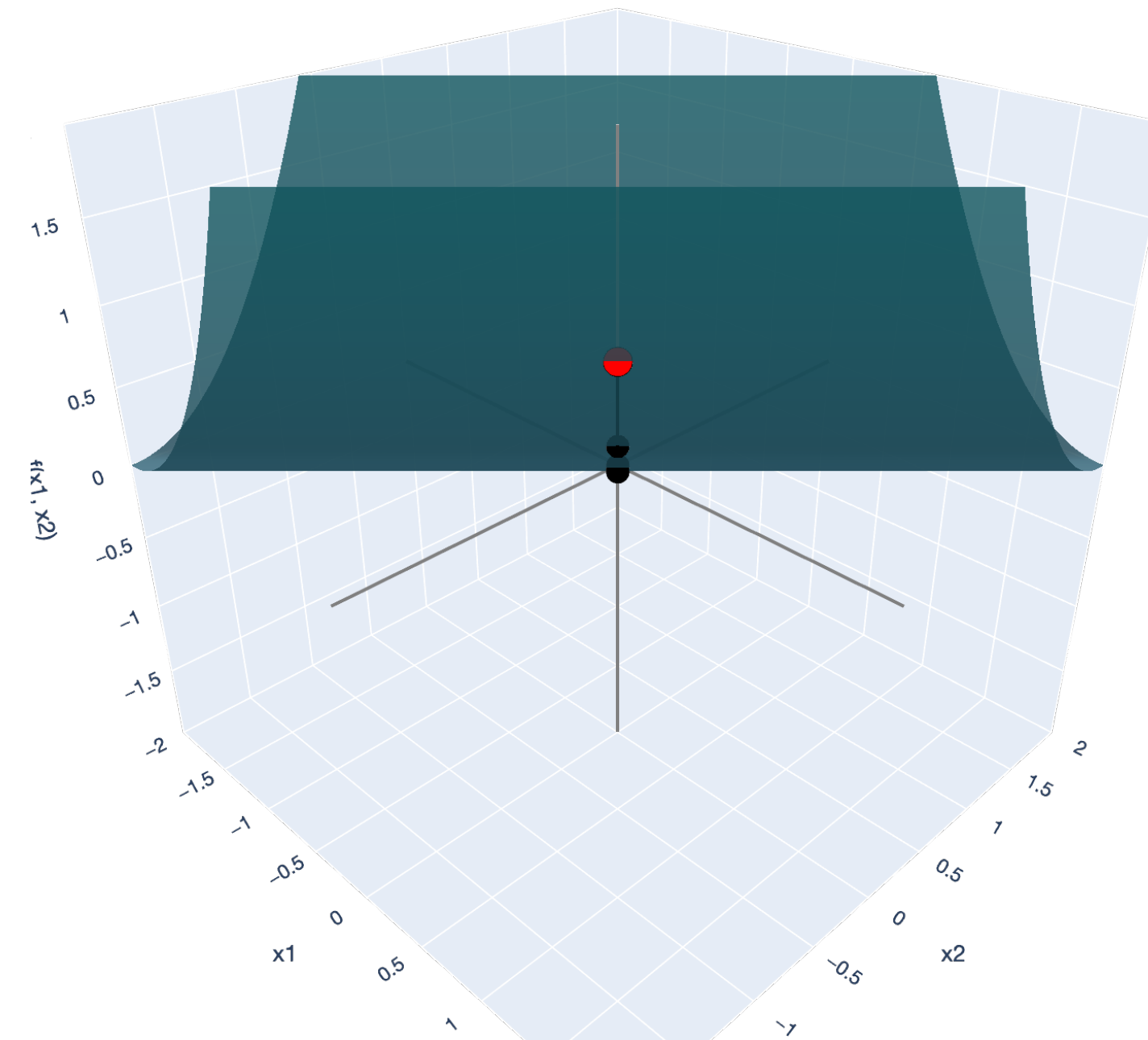
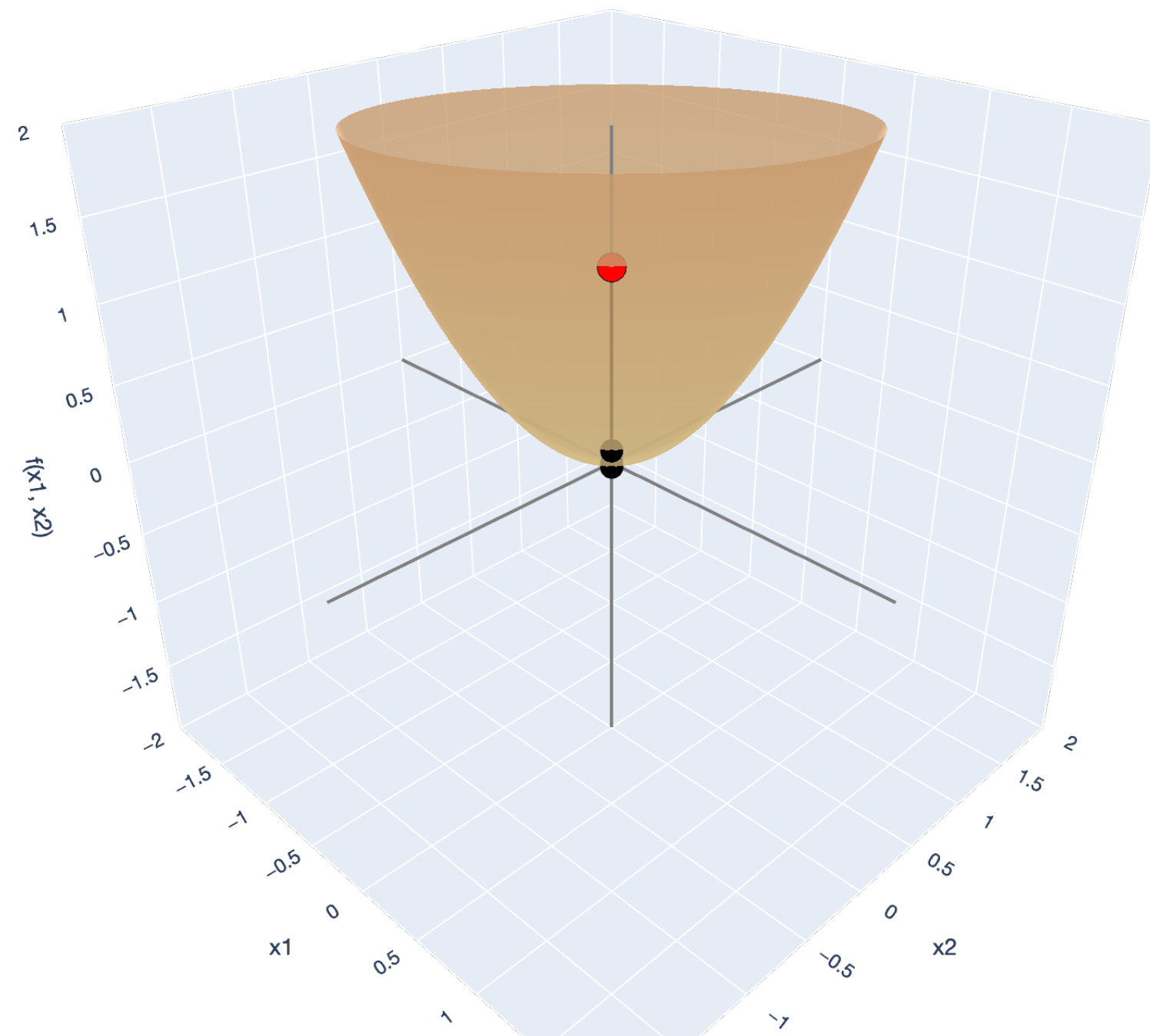
— x1 — x2 — u1 — u2 • y

Lesson Overview

Big Picture: Gradient Descent



$$\|Xw - y\|_2^2 = f(w)$$



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

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

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

$$X^T A X$$

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* $\mathbf{y} \in \mathbb{R}^d$.

$$\mathbf{X} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{x}_1 & \dots & \mathbf{x}_d \\ \downarrow & & \downarrow \end{bmatrix} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ \vdots & & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix}. \quad \left. \vphantom{\begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ \vdots & & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix}} \right\} \mathbf{X} \in \mathbb{R}^{n \times d}$$

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

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 $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ be its full SVD.

$$\mathbf{X} = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^{-1} = (\mathbf{V}^T)^{-1}\mathbf{\Sigma}^{-1}\mathbf{U}^{-1} \\ = \mathbf{V}(\mathbf{\Sigma}^{-1})\mathbf{U}^T$$

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

← left inverse: $\mathbf{\Sigma}^+\mathbf{\Sigma} = (\mathbf{\Sigma}^T\mathbf{\Sigma})^{-1}\mathbf{\Sigma}^T\mathbf{\Sigma} = \mathbf{I}$.

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

← right inverse:

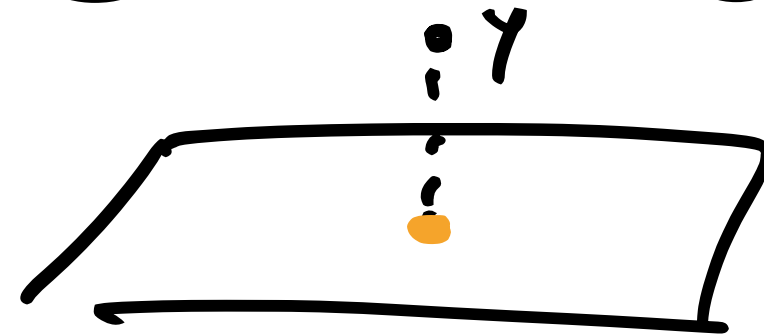
$$\mathbf{\Sigma}\mathbf{\Sigma}^+ = \mathbf{\Sigma}\mathbf{\Sigma}^T(\mathbf{\Sigma}\mathbf{\Sigma}^T)^{-1} \\ = \mathbf{I}$$

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

$$\mathbf{X}^+ := \mathbf{V}\mathbf{\Sigma}^+\mathbf{U}^T$$

Least Squares: SVD Perspective

Unified Picture



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

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

We can solve exactly.

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

which is an exact solution.

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

We approximate by least squares:

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

Choose

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

the best approximate solution:

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

$\neq 0$

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

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

Choose

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

the minimum norm solution:

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

Least Squares: SVD Perspective

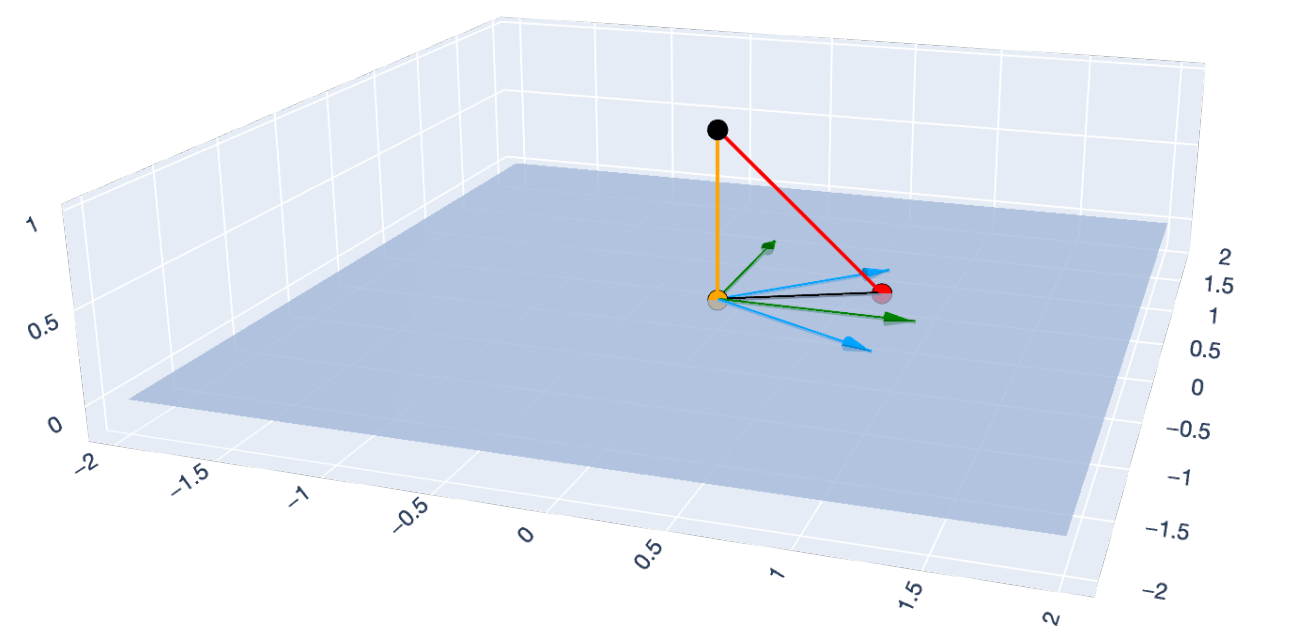
Unified Picture

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

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

We approximate by least squares:

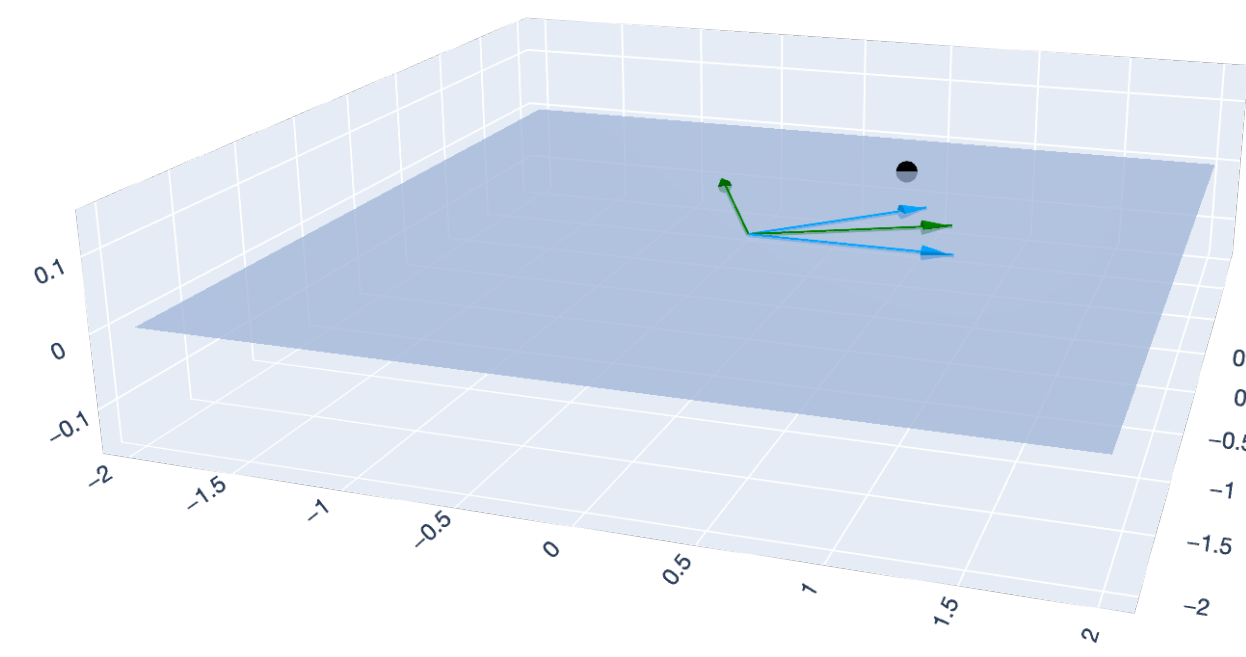
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$



— x1 — x2 — u1 — u2 — y - ^y — -y - ^y — -y - y • y • ^y • -y

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

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



— x1 — x2 — u1 — u2 • y

$y \in \text{col}(\mathbf{X})$.

$n = 3$
 $d = 2$

Singular Value Decomposition (SVD)

Matrix Decompositions

IT APPLIES TO ANY
MATRIX.

$$\boxed{\begin{array}{cccc} \underline{\mathbf{X}} & = & \underline{\mathbf{U}} & \underline{\mathbf{\Sigma}} & \underline{\mathbf{V}^T} \\ n \times d & & n \times n & n \times d & d \times d \end{array}}$$

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

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

$\mathbf{\Sigma} \in \mathbb{R}^{n \times d}$ 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 $\sigma_i > 0$.

$r = 5$ $\sigma_5 \approx 0$

***What other matrix
decompositions are out there?***

Eigendecomposition

Motivation: Linear Dynamical System

Population Change

Example of a linear dynamical system

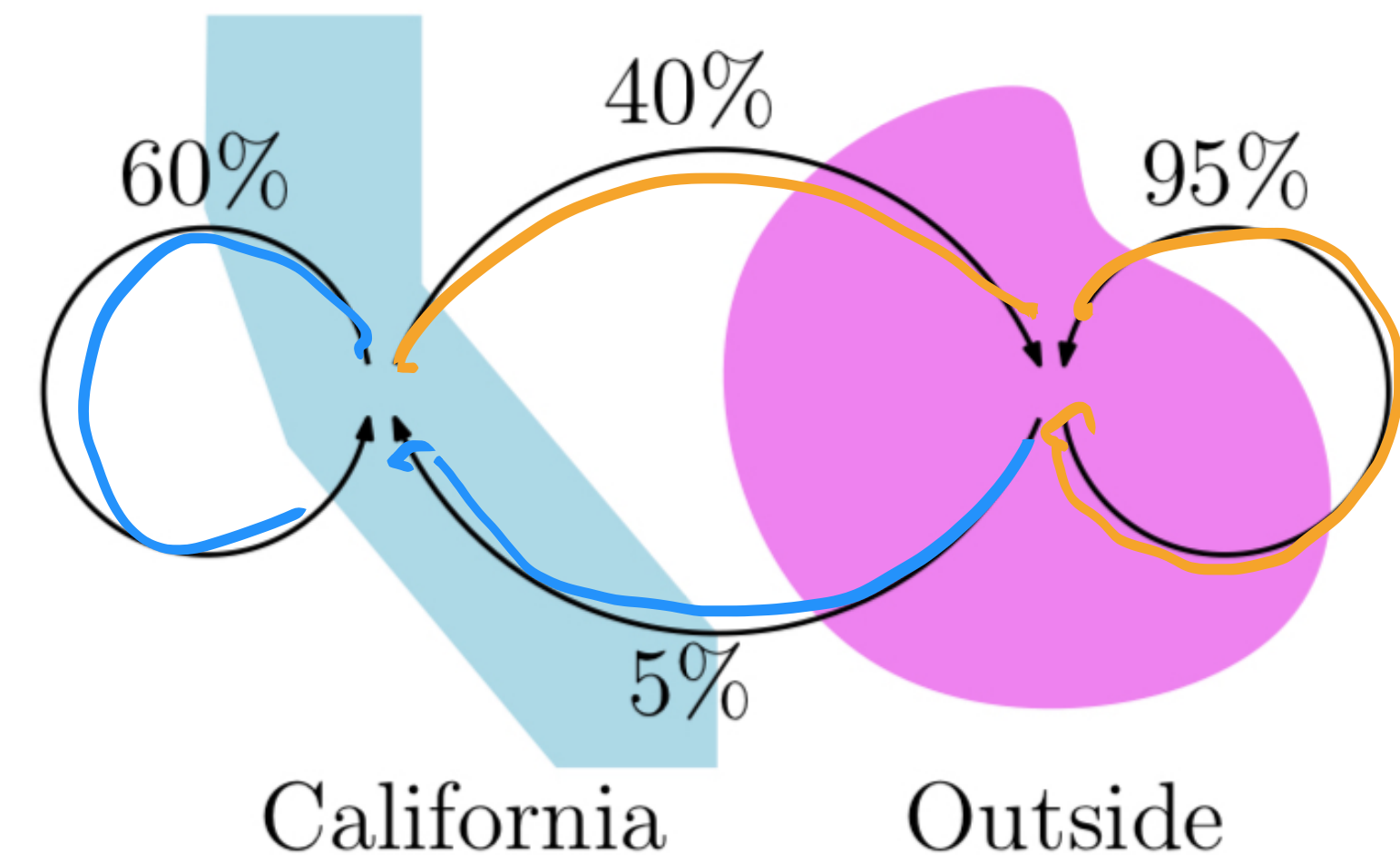
Consider the following example.

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:

$$\begin{aligned}\# \text{ inside} &= 0.6x_{in} + 0.05x_{out} \\ \# \text{ outside} &= 0.4x_{in} + 0.95x_{out}\end{aligned}$$



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.

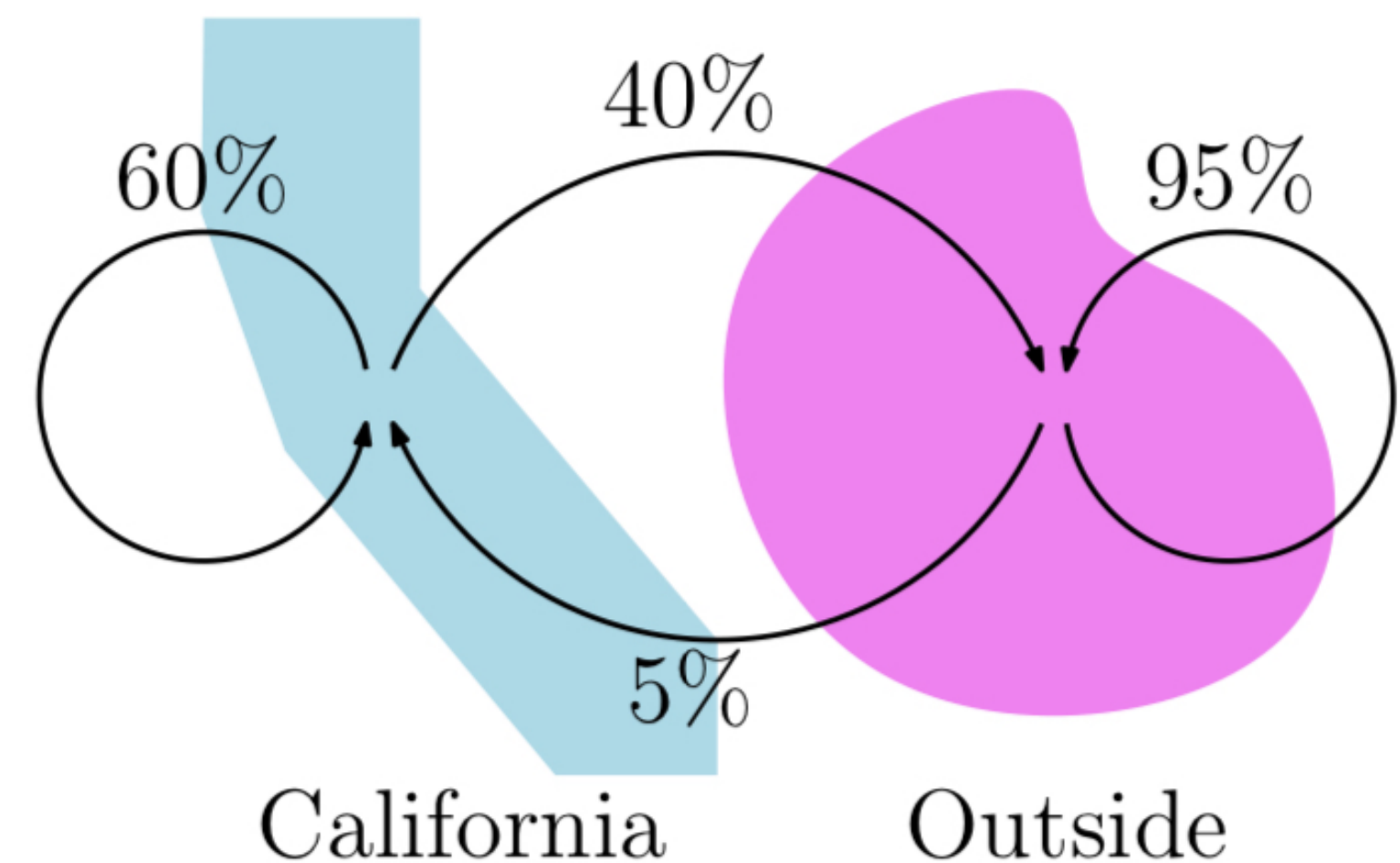
We can model this with a *transition matrix*

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

and a system of linear equations:

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

↑ # of people in/out at the end of year.



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

Population Change

Modeling with a transition matrix

Consider the transition matrix

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

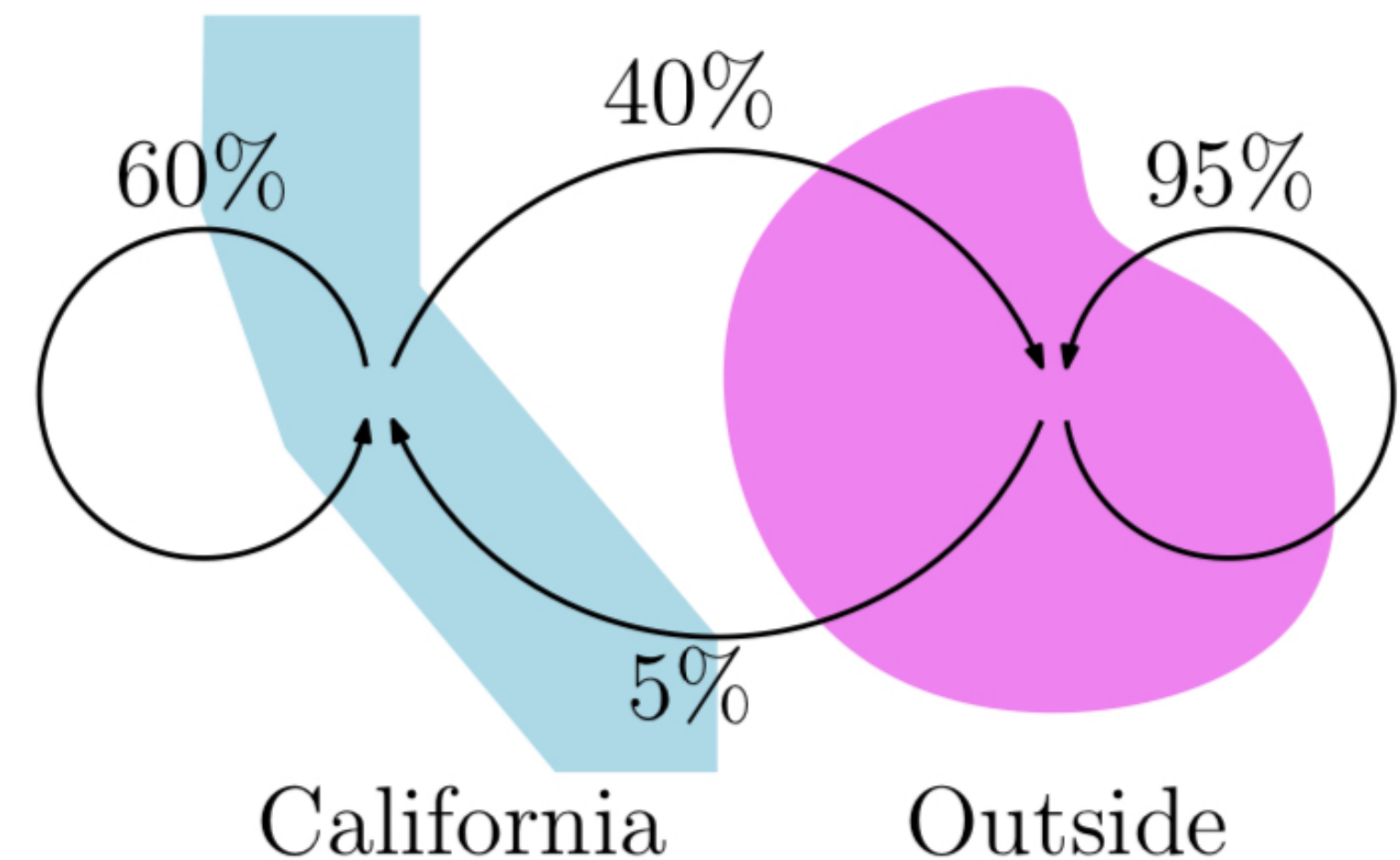
with a corresponding system of linear equations:

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

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

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

$$\begin{aligned} \mathbf{Ax} &\leftarrow \text{after 1 year.} \\ \mathbf{A}(\mathbf{Ax}) &\leftarrow \text{after 2 years } \mathbf{A}^2 \mathbf{x}. \end{aligned} \quad \left. \vphantom{\begin{aligned} \mathbf{Ax} \\ \mathbf{A}(\mathbf{Ax}) \end{aligned}} \right\} \mathbf{A}^t \mathbf{x}$$



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

Population Change

Modeling with a transition matrix

Consider the transition matrix

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

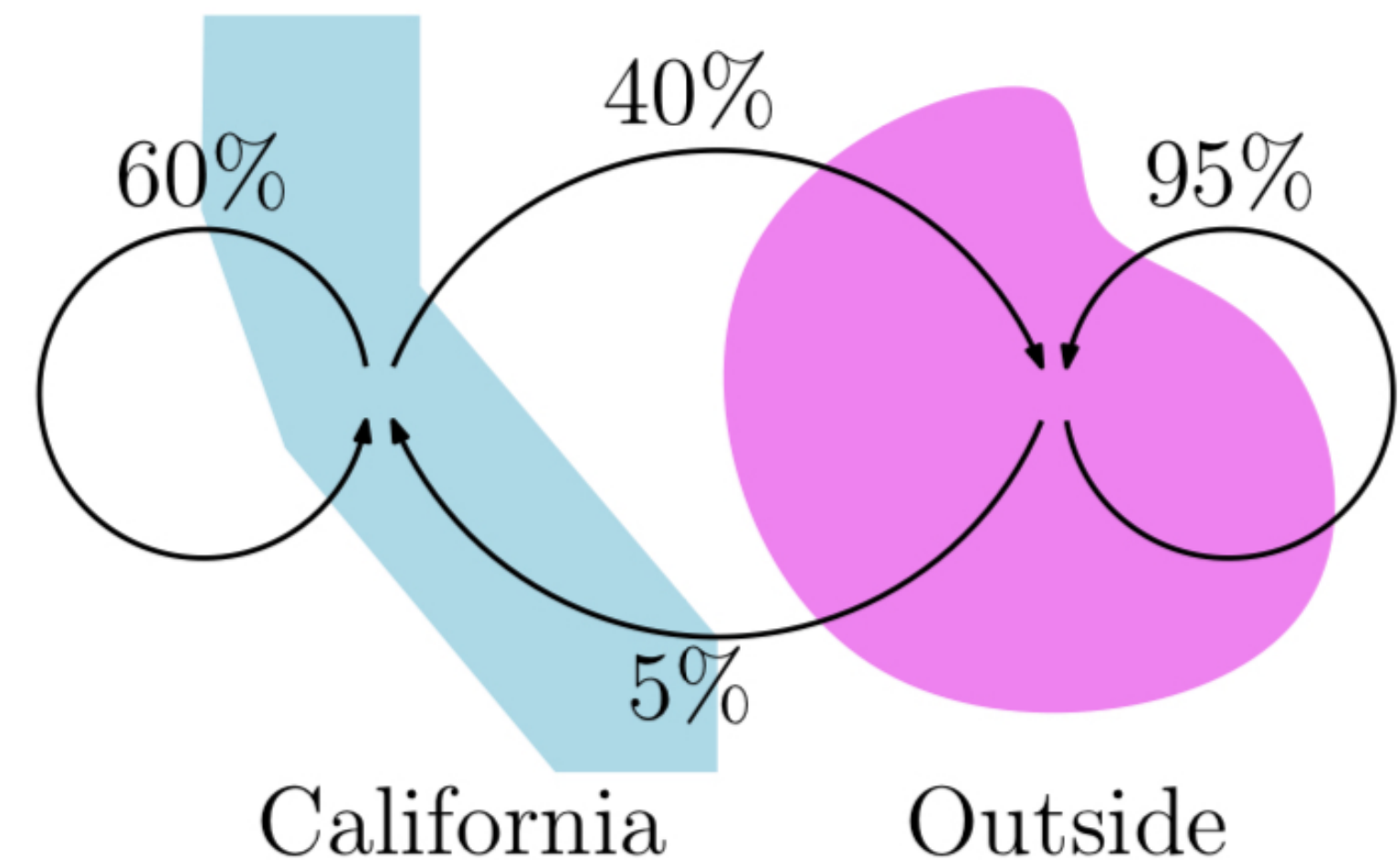
with a corresponding system of linear equations:

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

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?

$$\begin{aligned} \mathbf{x}^{(1)} &= \mathbf{Ax}^{(0)} \\ \mathbf{x}^{(2)} &= \mathbf{Ax}^{(1)} = \mathbf{AAx}^{(0)} = \mathbf{A}^2\mathbf{x}^{(0)} \\ &\vdots \\ \mathbf{x}^{(t)} &= \underbrace{\mathbf{AA}\dots\mathbf{A}}_{t \text{ products}} \mathbf{x}^{(0)} = \mathbf{A}^t\mathbf{x}^{(0)} \end{aligned}$$



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

Population Change

Modeling with a transition matrix

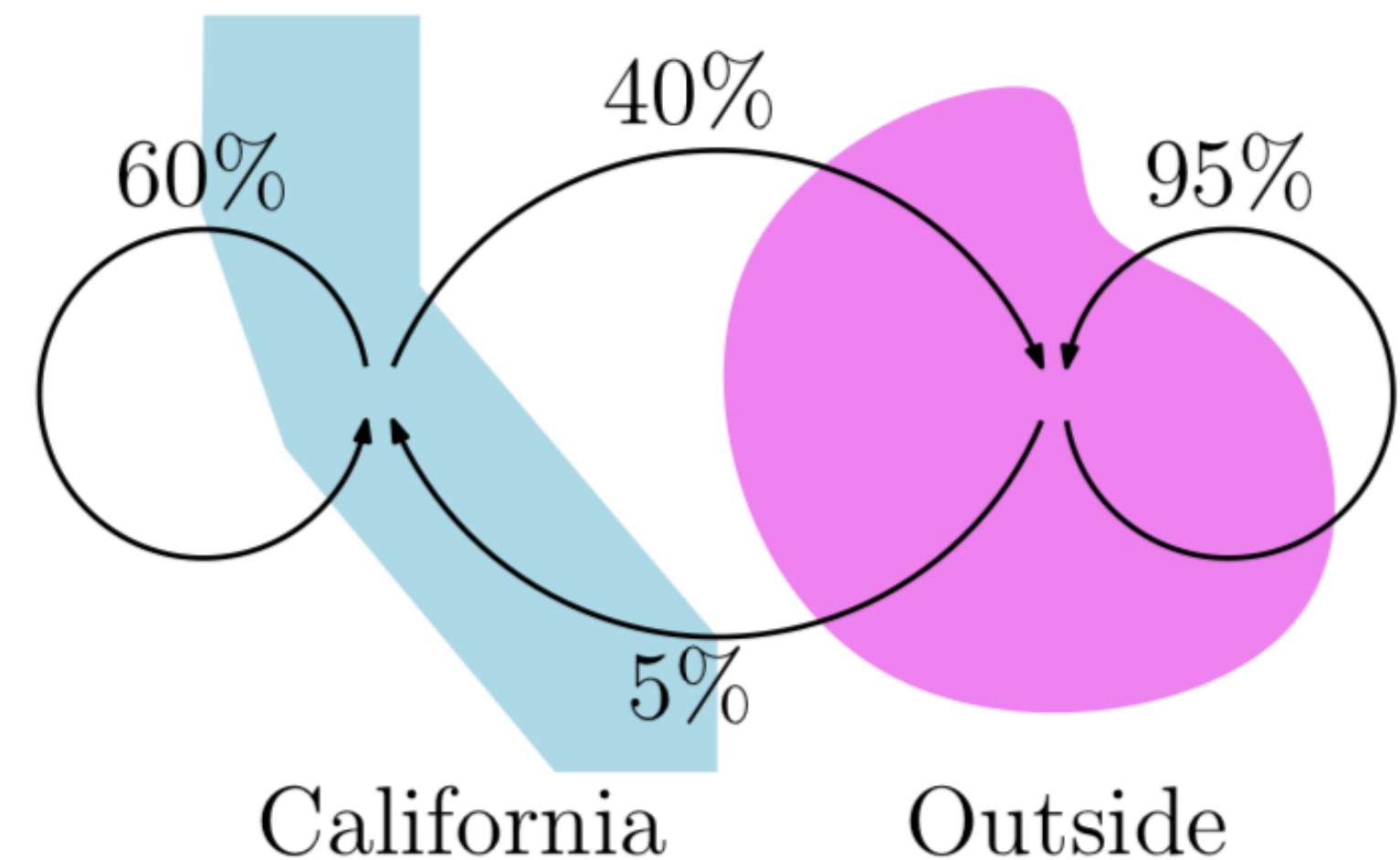
$$\mathbf{Ax} = \begin{bmatrix} \text{in} \rightarrow \text{in} & \text{out} \rightarrow \text{in} \\ \text{in} \rightarrow \text{out} & \text{out} \rightarrow \text{out} \end{bmatrix} \begin{bmatrix} x_{in} \\ x_{out} \end{bmatrix} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix} \begin{bmatrix}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

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...

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

t times.

Population Change

Easy computation 😊

Assume I gave you a couple of vectors, *Magically*, $\mathbf{u} = (1, 8)$ 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:

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

$$\mathbf{A}^t\mathbf{v} = \begin{bmatrix} 0.6 & 0.05 \\ 0.4 & 0.95 \end{bmatrix}^t \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \left(\frac{11}{20}\right)^t \begin{bmatrix} -1 \\ 1 \end{bmatrix} \implies \mathbf{A}^t\mathbf{v} = \left(\frac{11}{20}\right)^t \mathbf{v} \quad \textcircled{2}$$

Population Change

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

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

$$\begin{aligned} \textcircled{1} \mathbf{A}^t \mathbf{u} &= \mathbf{u} \\ \textcircled{2} \mathbf{A}^t \mathbf{v} &= \left(\frac{11}{20} \right)^t \mathbf{v} \end{aligned}$$

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.

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

scalars

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

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

and starting

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:

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

In matrix form:

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

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:

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

$$\mathbf{V}^{-1} = \begin{bmatrix} \uparrow & \uparrow \\ \mathbf{u} & \mathbf{v} \\ \downarrow & \downarrow \end{bmatrix}^{-1}$$

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

$$A \begin{bmatrix} 1 \\ 8 \end{bmatrix} = \begin{bmatrix} 1 \\ 8 \end{bmatrix}$$

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:

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



Diagonal Matrices

Why we like diagonal matrices

Multiplying diagonal matrices with themselves many times is easy:

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

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:

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

*depended on A
(transition matrix)*

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

Eigendecomposition

Intuition and Definition

Eigenvectors and eigenvalues

Intuition

$$\rightarrow T_A: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

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

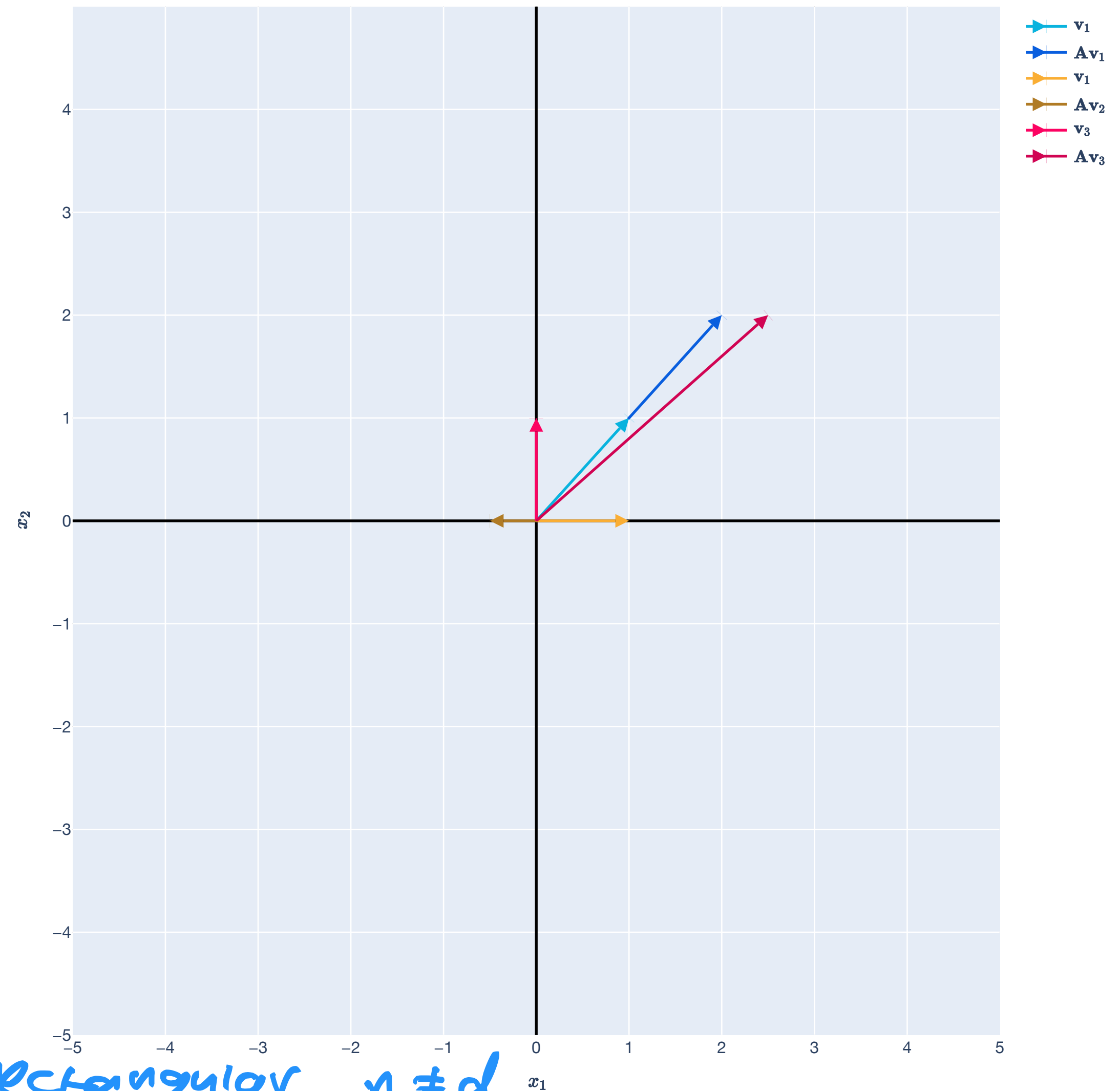
This represents a linear transformation from \mathbb{R}^d to \mathbb{R}^d .

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!

doesn't make sense to ask for rectangular $n \neq d$.



Eigenvectors and eigenvalues

Definition

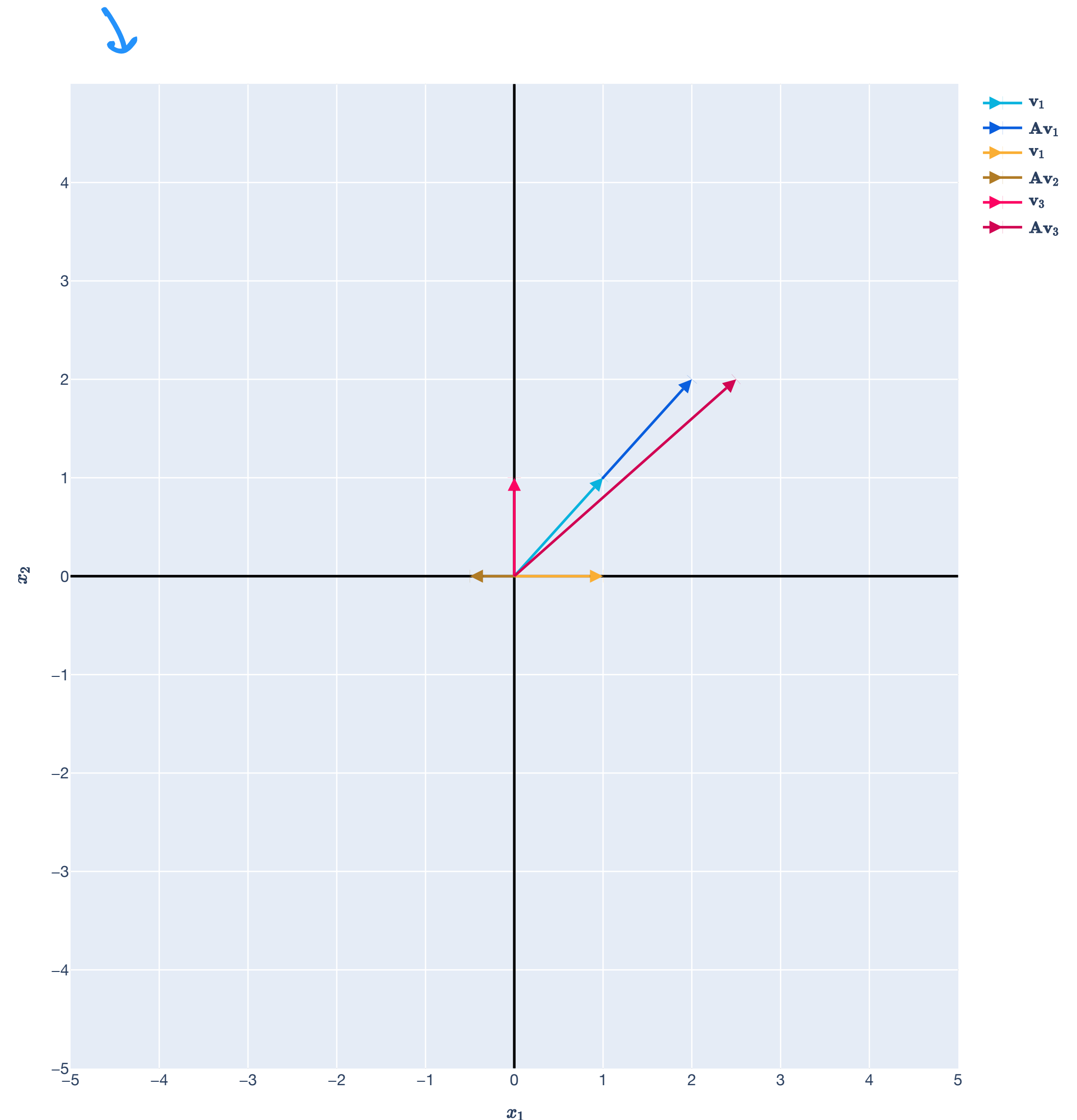
Let $\mathbf{A} \in \mathbb{R}^{d \times d}$ be a *square* matrix.

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

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v}.$$

The scalar λ is the eigenvalue associated with the eigenvector \mathbf{v} .

Eigenvectors/eigenvalues are properties of square matrices!



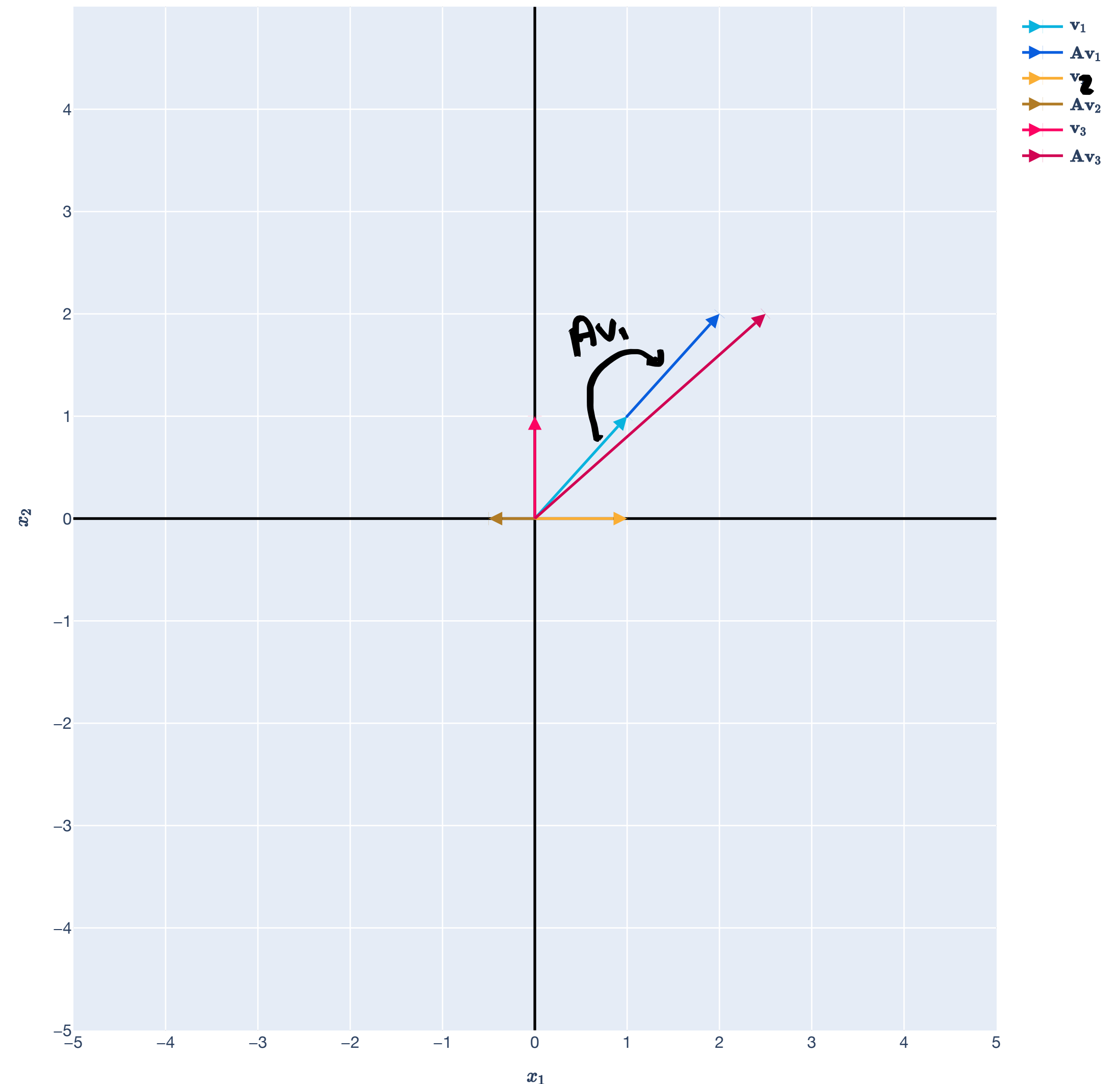
Eigenvectors and eigenvalues

Example

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

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

What happens to the vector $\mathbf{v}_1 = (1,1)$?



Eigenvectors and eigenvalues

Example

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

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

What happens to the vector $\mathbf{v}_2 = (1, 0)$?

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



Eigenvectors and eigenvalues

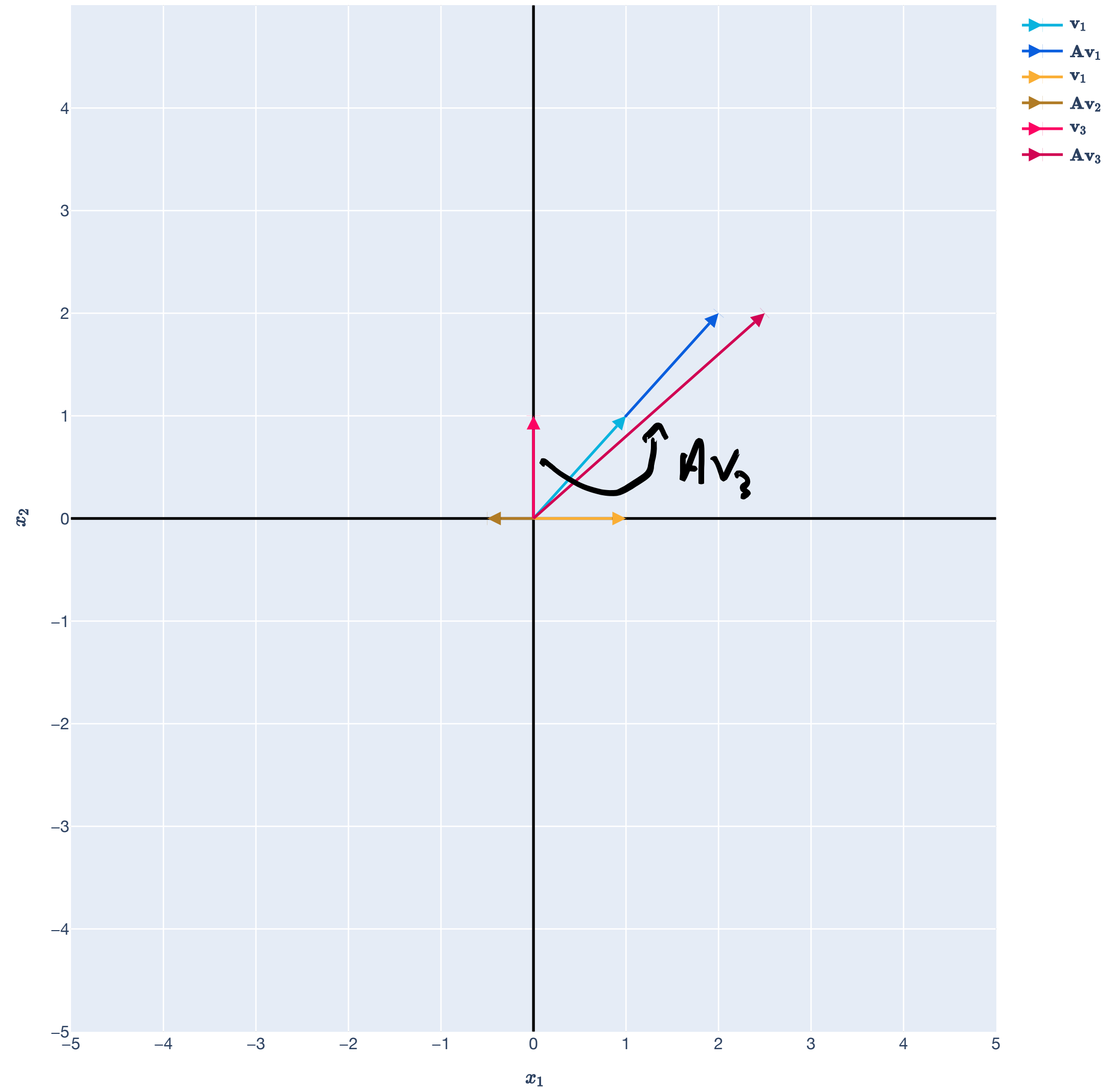
Example

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

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

What happens to the vector $\mathbf{v}_3 = (0, 1)$?

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



Eigenvectors and eigenvalues

Example

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

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

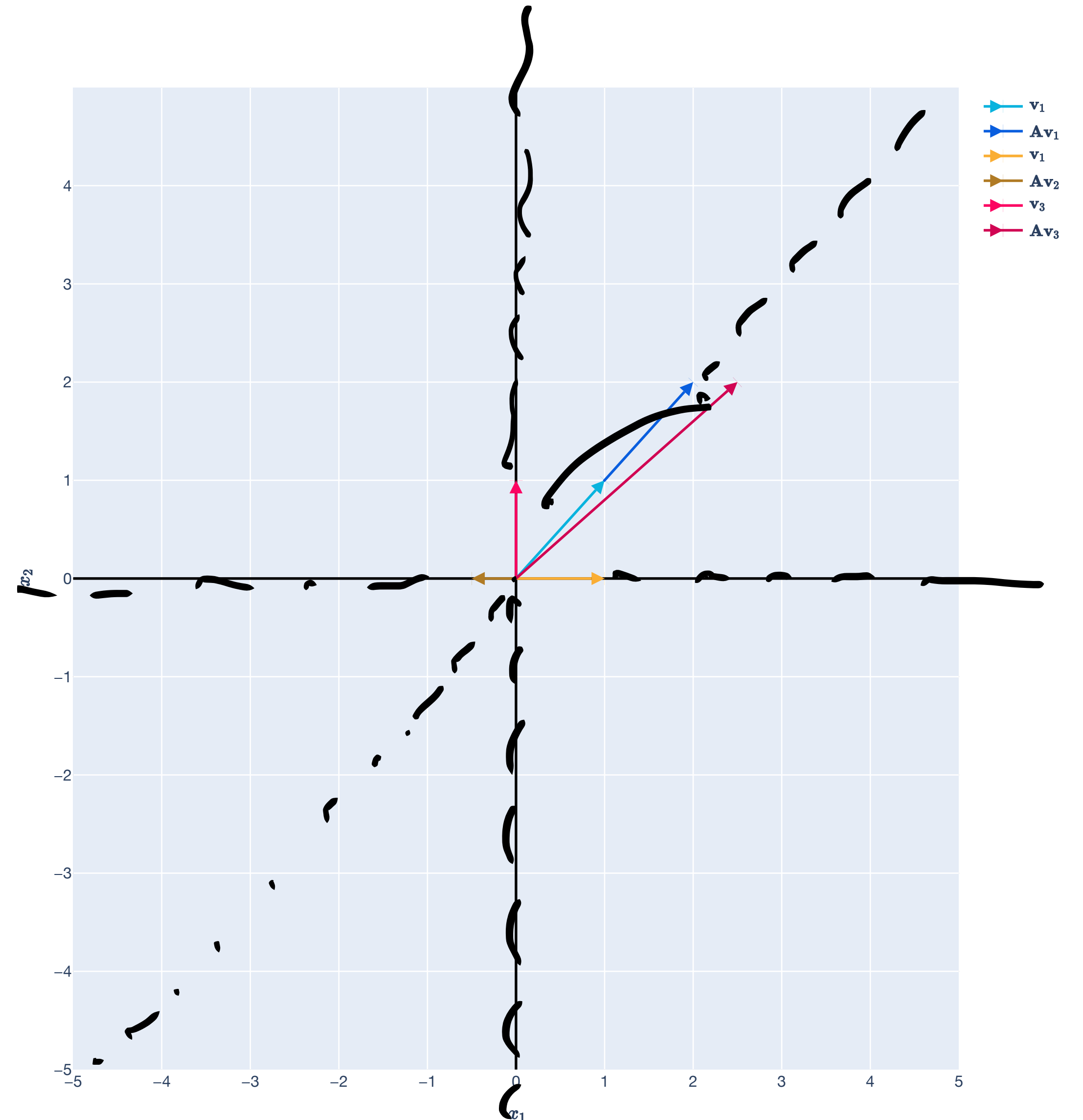
Eigenvectors (with eigenvalues $\lambda_1 = 2$ and $\lambda_2 = -1/2$):

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

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

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 ^A 3 Blue | Braun

Example

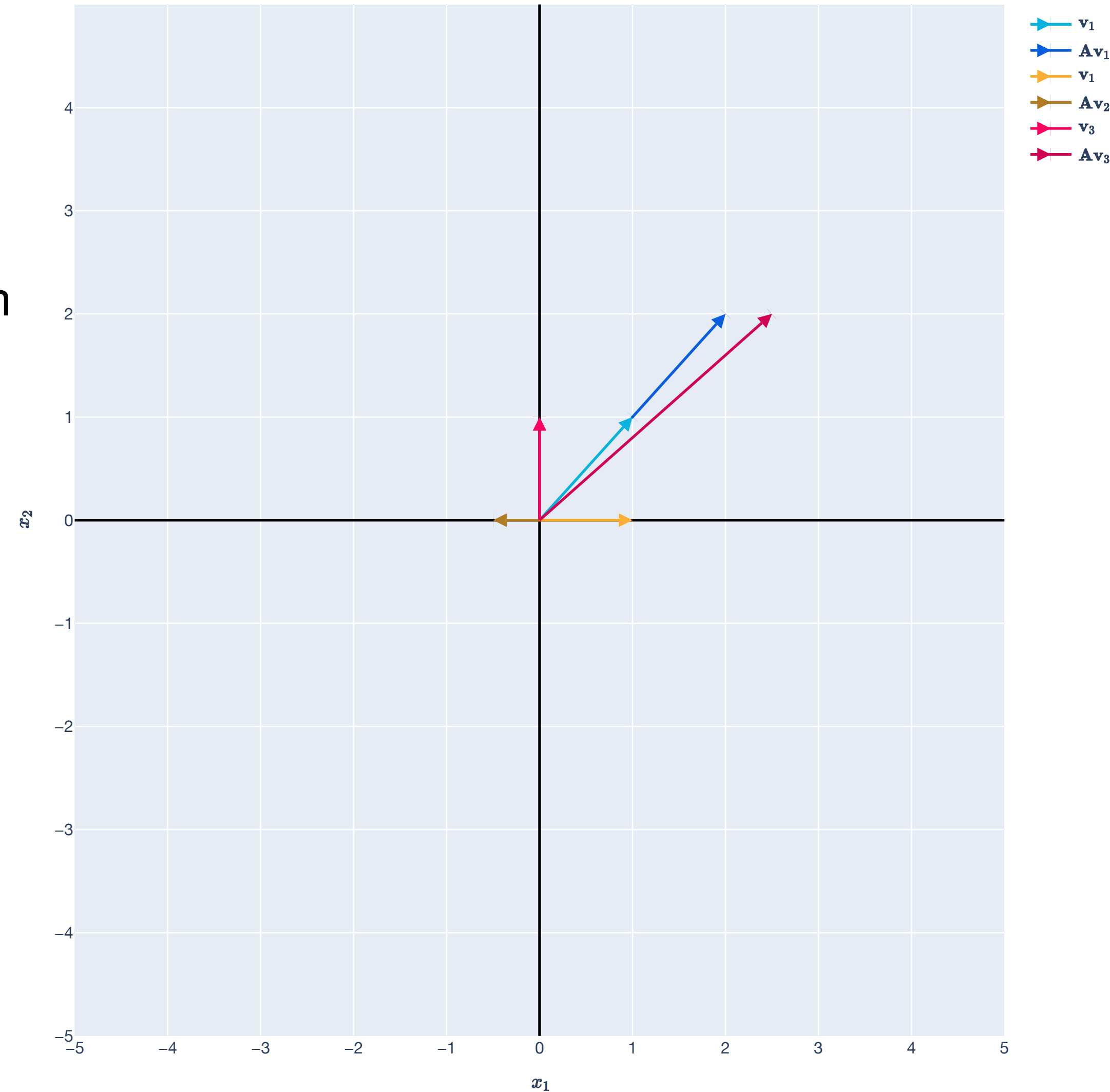
$$\mathbf{A} = \begin{bmatrix} -1/2 & 5/2 \\ 0 & 2 \end{bmatrix}$$

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

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

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

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



Eigenvectors and eigenvalues

Example

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

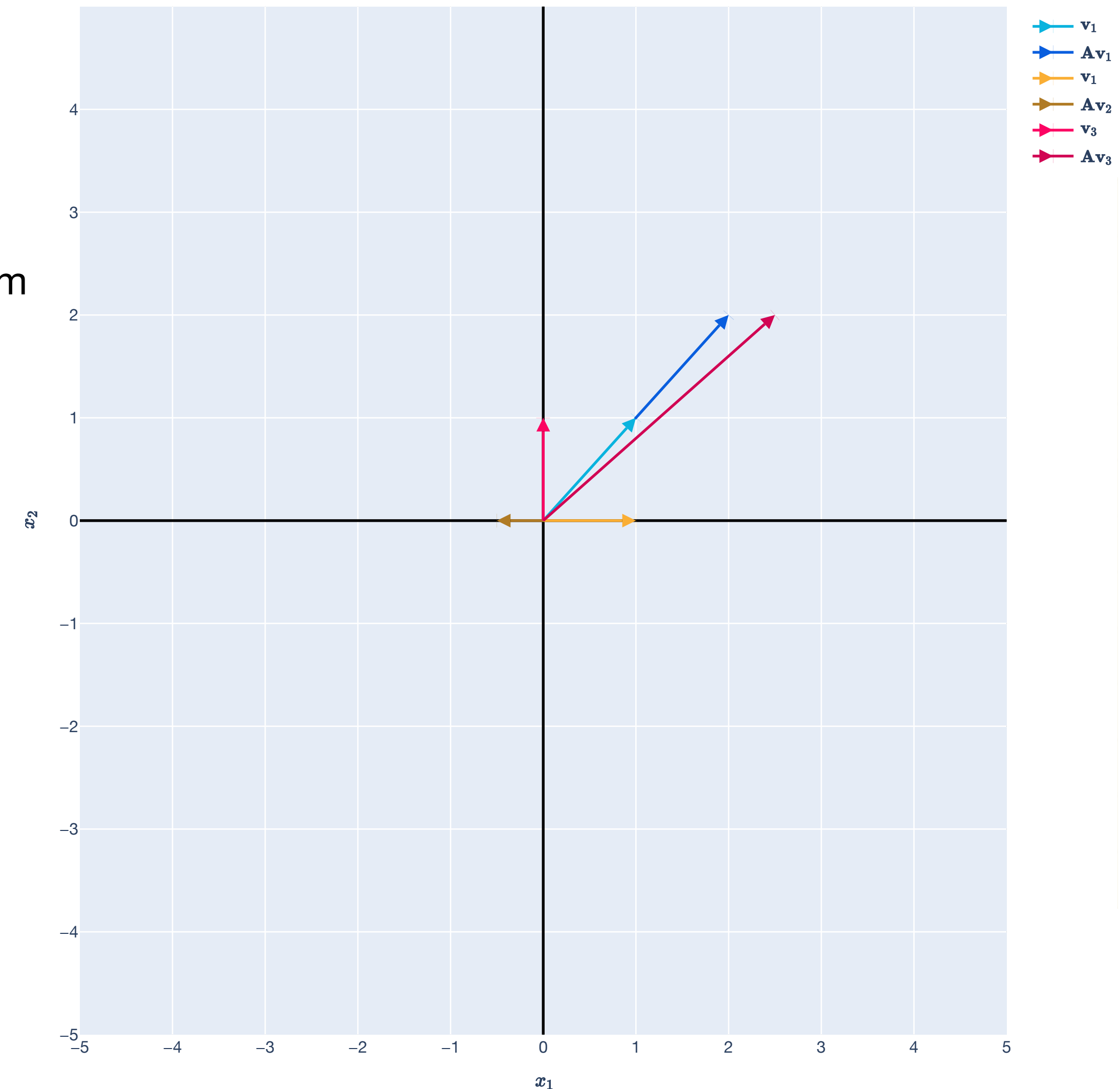
$\mathbf{v}_1 = (1,1)$ and $\mathbf{v}_2 = (1,0)$ are linearly independent *eigenvectors* — they form a basis for \mathbb{R}^2 . Their *eigenvalues* are $\lambda_1 = 2$ and $\lambda_2 = -1/2$.

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:

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



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 :

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

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

coordinates of \mathbf{x} in the eigen vector basis.

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:

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

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:

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

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

Eigendecomposition

Definition

$$\boxed{\begin{bmatrix} a \\ b \end{bmatrix} = V^{-1}x}$$

$$V \begin{bmatrix} a \\ b \end{bmatrix} = x$$

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

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

$$\begin{aligned} Av_1 &= \lambda_1 v_1 \\ &\vdots \\ Av_d &= \lambda_d v_d \end{aligned}$$

$$V^{-1}x \rightarrow \lambda$$

scales.

Then, A has the eigendecomposition:

$$A = V \Lambda V^{-1} = \begin{bmatrix} \uparrow & \dots & \uparrow \\ v_1 & \dots & v_d \\ \downarrow & \dots & \downarrow \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \dots & \lambda_d \end{bmatrix} \begin{bmatrix} \uparrow & \dots & \uparrow \\ v_1 & \dots & v_d \\ \downarrow & \dots & \downarrow \end{bmatrix}^{-1}$$

Such a matrix is said to be diagonalizable.

→ Eigendecomposition exists.

Eigendecomposition

Example

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

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

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

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

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

Eigendecomposition

Example

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

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

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

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

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

\mathbf{Q} $\mathbf{\Lambda}$ \mathbf{Q}^{-1}

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

$$\det(\mathbf{A} - \lambda\mathbf{I}) = 0 \quad \times$$

numpy. np.linalg.eig

Eigendecomposition

Connection with SVD

Connection with SVD

Eigendecomposition from SVD

$$\mathbb{R}^d \rightarrow \mathbb{R}^d$$

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

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

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

$$\approx \boxed{\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T}$$

Connection with SVD

Eigendecomposition from SVD

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

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

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

$$\begin{aligned} \mathbf{A} &= \underline{\mathbf{X}}^T \underline{\mathbf{X}} \rightarrow \mathbf{I} \\ &= \underline{\mathbf{V}} \underline{\mathbf{\Sigma}}^T \underline{\mathbf{U}}^T \underline{\mathbf{U}} \underline{\mathbf{\Sigma}} \underline{\mathbf{V}}^T \\ &= \underline{\mathbf{V}} \underline{\mathbf{\Sigma}}^T \underline{\mathbf{\Sigma}} \underline{\mathbf{V}}^T \end{aligned}$$

Connection with SVD

Eigendecomposition from SVD

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

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

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

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

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

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

Connection with SVD

Eigendecomposition from SVD

$$A = \underbrace{X^T X}$$

PS (2): for proof. $n \times n$ $n \times d$, $d \times d$

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

$$\underbrace{\left\{ \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0, \right.}_{\text{right singular vecs.}} \quad V = \begin{bmatrix} | & & | \\ v_1 & \dots & v_d \\ | & & | \end{bmatrix}$$

and let v_1, \dots, v_d be the columns of $V \in \mathbb{R}^{d \times d}$. Then, each v_i is an eigenvector for A with corresponding eigenvalue $\lambda_i = \sigma_i^2$, and the eigendecomposition of A is:

$$A = V \Lambda V^T$$

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:

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

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

$$\begin{bmatrix} \lambda_1 & & & \\ & \dots & & \\ & & \lambda_r & \\ & & & 0 & \dots & 0 \end{bmatrix}$$

A handwritten diagram illustrating the dimensions of matrices in the equation $A = X^T X$. The matrix A is enclosed in a box and labeled as $d \times d$. The matrix X is labeled as $d \times n$. The product $X^T X$ is labeled as $n \times d$.

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:

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

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

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

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

Positive Semidefinite (PSD) Matrices

Example

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

Positive Semidefinite (PSD) Matrices

Example

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

Its “square root” is the matrix

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

To verify:

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

PSD Matrices and Eigendecomposition

Connection to eigenvalues

$$X = U \Sigma V^T$$

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

$$A = V \Lambda V^T,$$

with orthonormal eigenvectors v_1, \dots, 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 $\lambda_1, \dots, \lambda_d \geq 0$.

d nonnegative eigenvalues.

Positive Semidefinite (PSD) Matrices

Example

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

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

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

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

$\lambda_i \geq 0$.

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

Positive Semidefinite (PSD) Matrices

Third definition

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

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

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

$$\underbrace{\mathbf{x}}_{\mathbb{R}^d}^T \mathbf{A} \underbrace{\mathbf{x}}_{\mathbb{R}^d} \in \mathbb{R}.$$

Positive Semidefinite (PSD) Matrices

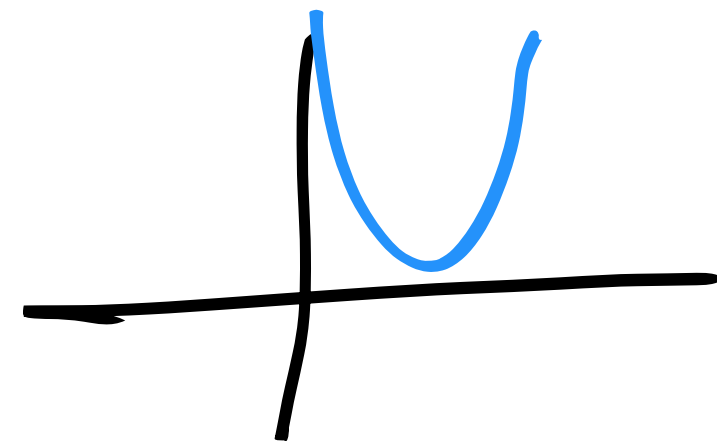
Example

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

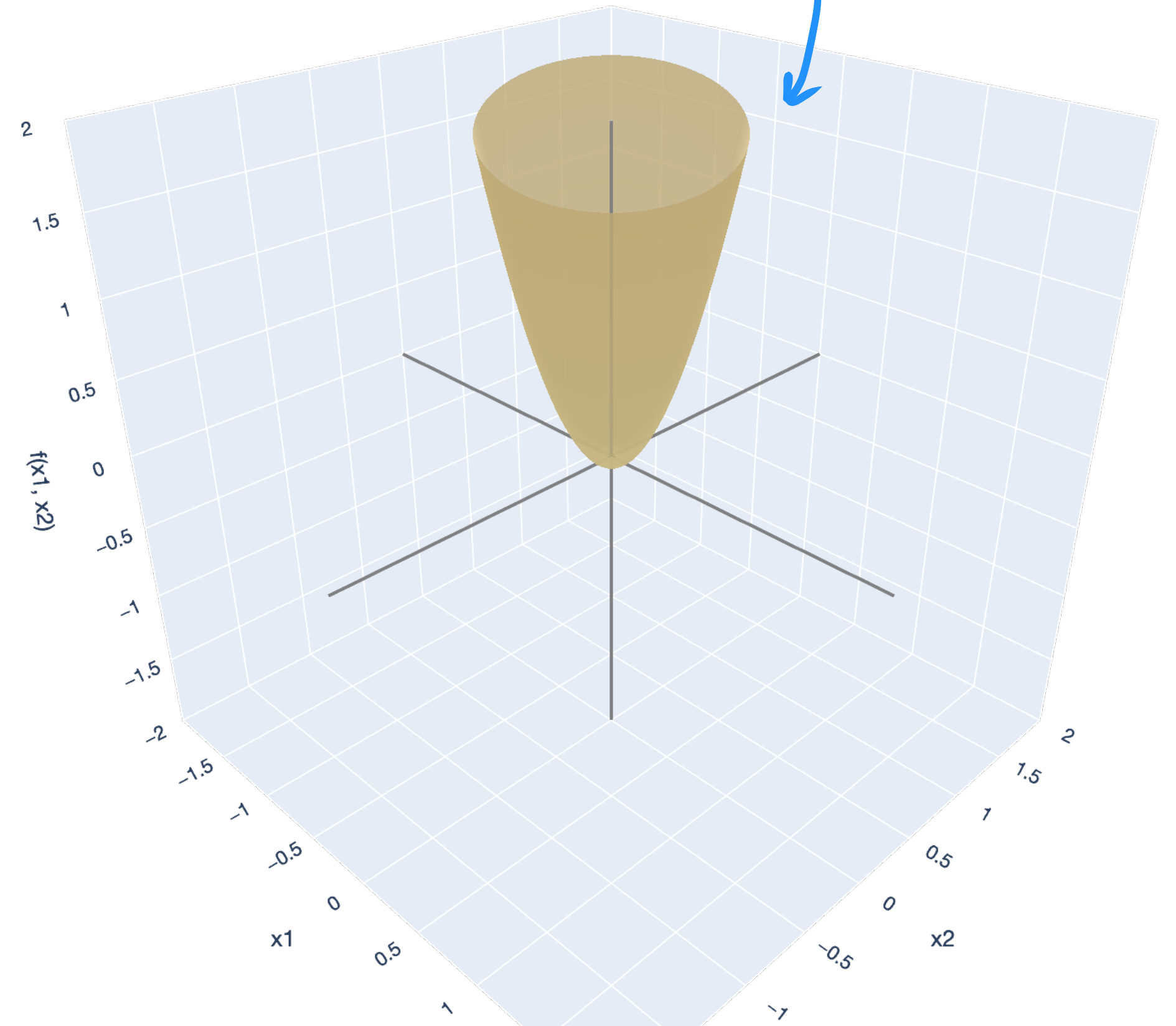
Consider any vector $\mathbf{x} = (x_1, x_2) \in \mathbb{R}^d$.

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

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



$$f(x_1, x_2) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$$



— 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...

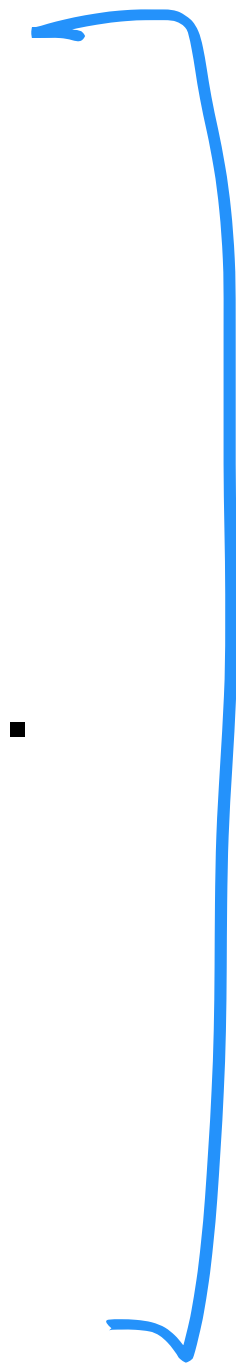
there exists $\mathbf{X} \in \mathbb{R}^{n \times d}$ such that $\mathbf{A} = \mathbf{X}^T \mathbf{X}$.



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



$\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$ 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 $\mathbf{A} = \mathbf{X}^T \mathbf{X}$.

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

$$\begin{aligned} \lambda_1 &= 4 \\ \lambda_2 &= 1. \end{aligned}$$

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

strictly

Spectral Theorem

Statement

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

A: When \mathbf{A} is positive semidefinite! $\rightarrow A^T = A$

$A = X^T X \rightarrow \text{SVD} \rightarrow v_1, \dots, v_d$ is a basis

But even more generally...

Spectral Theorem

Statement

$$A^T = A.$$

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

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

$$\boxed{X = U\Sigma V^T} \leftarrow \text{SVD works for any } \underline{n \times d}.$$

But, in this generality, λ_i can be negative!

$$\mathbb{R}^{100} \rightarrow \mathbb{R}^3 \text{ or } \mathbb{R}^2.$$

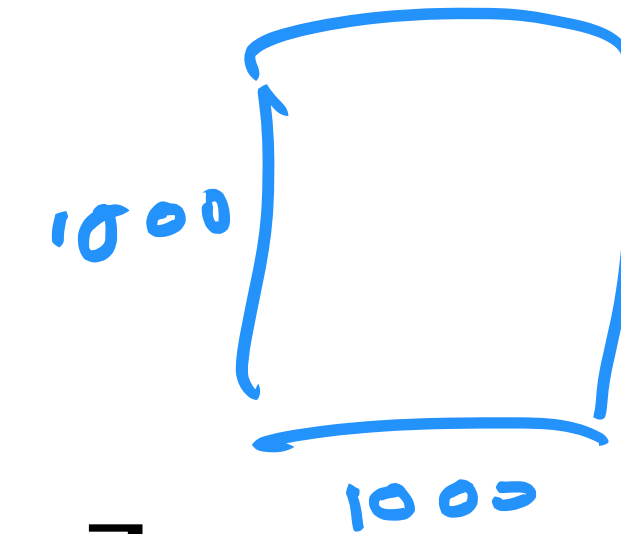
Principal Components Analysis

Application of Eigendecomposition

Principal Components Analysis \equiv Eigen decomposition

Example: “Eigenfaces” and facial recognition

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



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

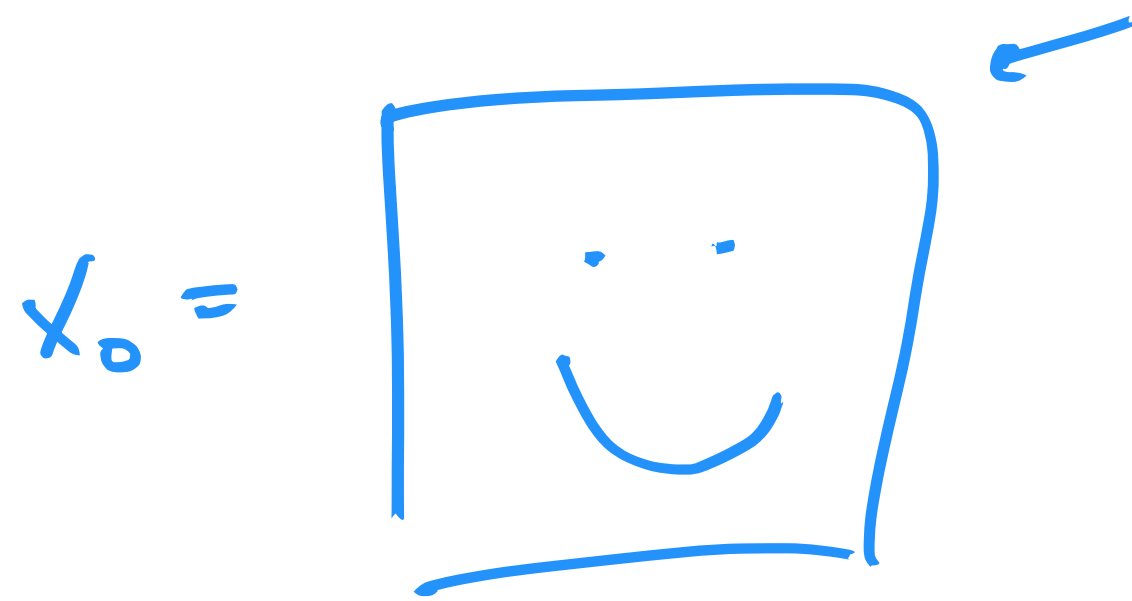
Each row is a “flattened” image vector. Typically, each pixel is in $[0, 255]$ for grayscale images.

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

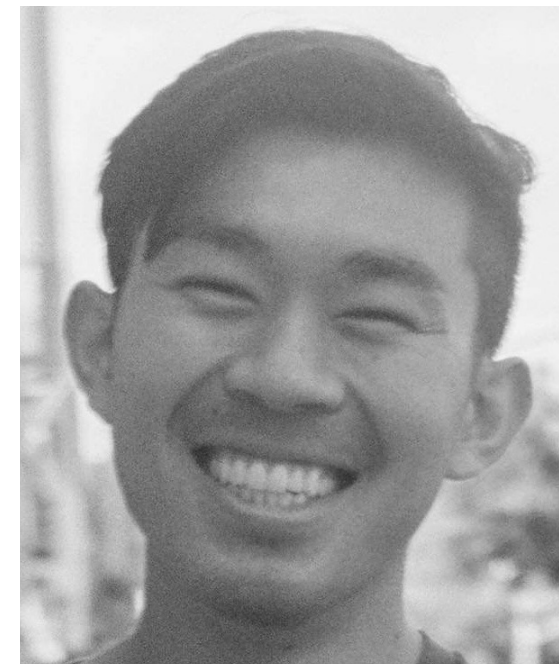
Principal Components Analysis

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 $\|\mathbf{x}_0 - \mathbf{x}_i\|$).

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

Byte.

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!

$$= \underbrace{0.45}_{\hat{w}_1} \mathbf{v}_1 + \underbrace{0.21}_{\hat{w}_2} \mathbf{v}_2 + \underbrace{0.12}_{\hat{w}_3} \mathbf{v}_3 + 0.05 \mathbf{v}_4 + \dots$$

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 $\|\Pi(\mathbf{x}_0) - \Pi(\mathbf{x}_i)\|$.

$$\min_{\hat{\mathbf{w}}} \|\mathbf{x}_0 - \mathbf{V}\hat{\mathbf{w}}\|^2$$

$\hat{\mathbf{w}} \rightarrow \underbrace{(\mathbf{V}\hat{\mathbf{w}})}_{\hat{\mathbf{w}} \in \mathbb{R}^k}$

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

Principal Components Analysis

Example: PCA in 2D

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

Positive Semidefinite

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

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

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

2x2

Principal Components Analysis

Example: PCA in 2D

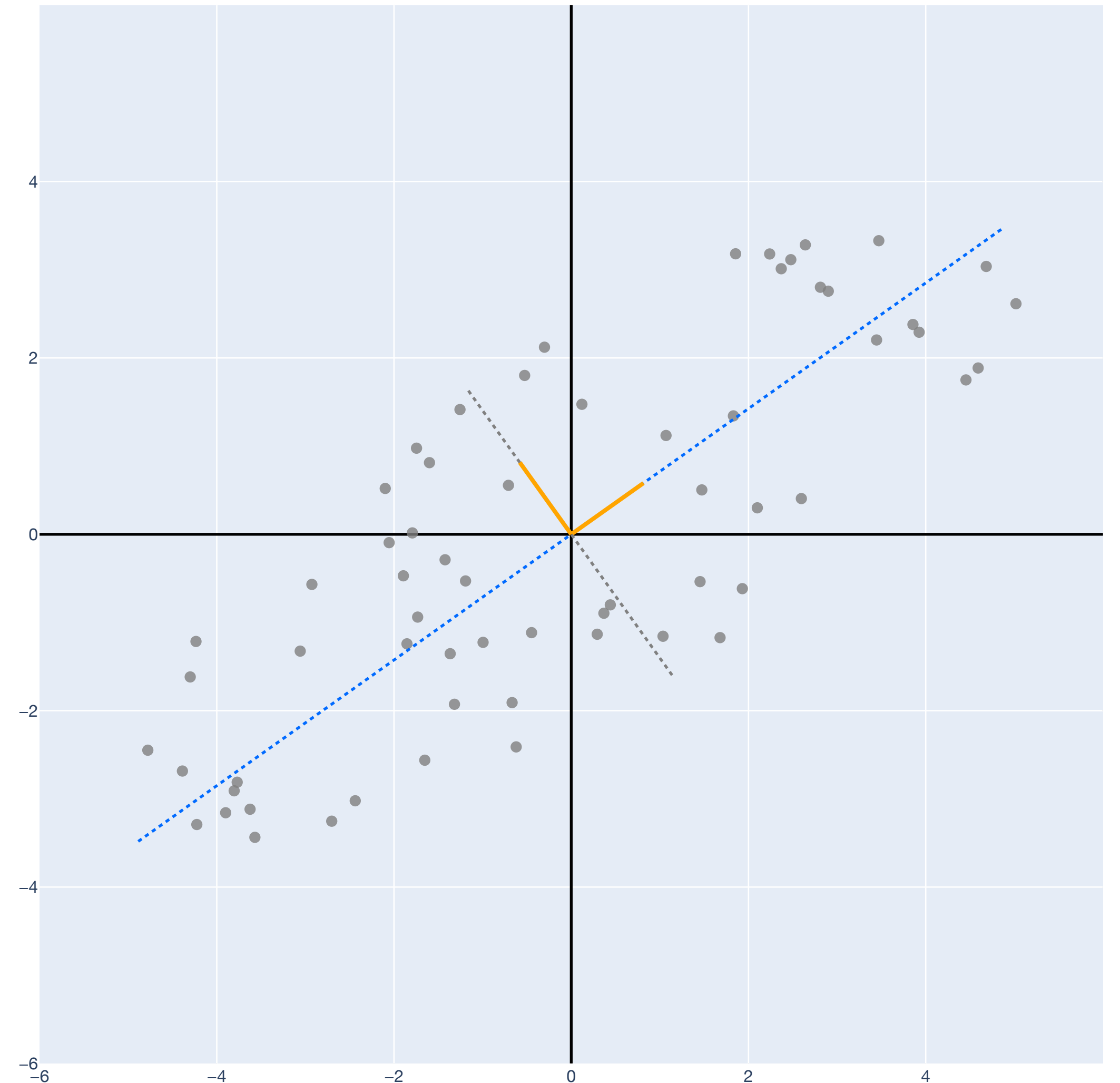
Observed: Matrix of training points $\mathbf{X} \in \mathbb{R}^{n \times 2}$:

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

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

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

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

$C = X^T X$ is symmetric

Eigendecomposition and PCA

\Rightarrow spectral theorem \Rightarrow Eigendecomposition

PCA = Eigendecomposition of the covariance matrix!

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

Therefore, it is diagonalizable with eigendecomposition:

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

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!

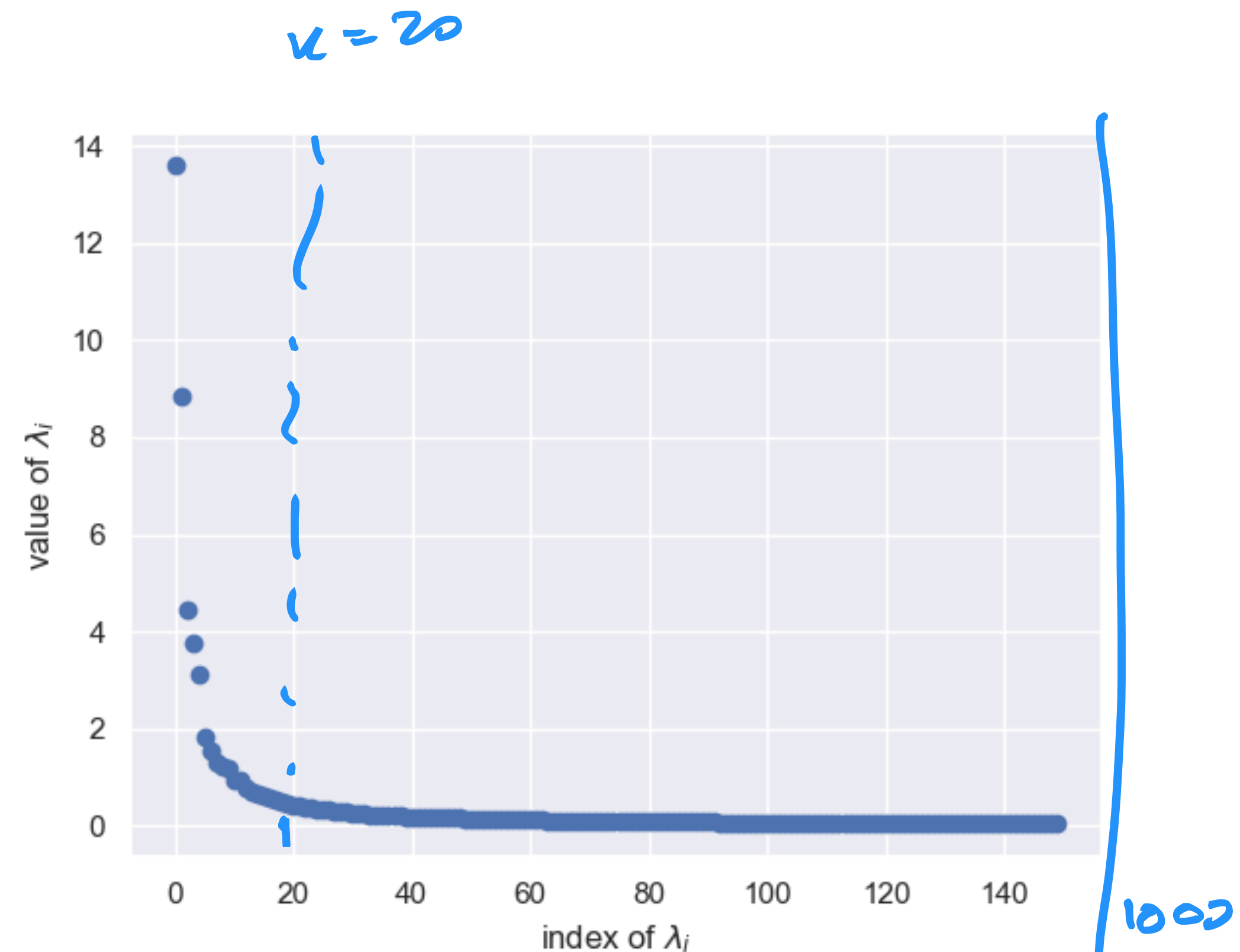
Consider a (column-centered) dataset $\mathbf{X} \in \mathbb{R}^{n \times d}$ and construct its covariance matrix $\mathbf{C} = \mathbf{X}^T \mathbf{X} \in \mathbb{R}^{d \times d}$. By definition, \mathbf{C} is positive semidefinite.

Therefore, it is diagonalizable with eigendecomposition:

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

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.



150

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:

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

(Could have also just taken the right singular vectors of $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$ 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:

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

$$X\mathbf{w}^* = \mathbf{y}.$$

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* $\epsilon = (\epsilon_1, \dots, \epsilon_n) \in \mathbb{R}^n$.

$$y_i = \underbrace{\mathbf{x}_i^\top \mathbf{w}^*}_{\text{true fit}} + \underbrace{\epsilon_i}_{\text{error}} \text{ for all } i \in [n] \quad \epsilon_i \sim \mathcal{D}.$$

$$\mathbf{y} = \mathbf{X}\mathbf{w}^* + \vec{\epsilon} \in \mathbb{R}^n$$

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 $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$?

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\mathbf{w}^* + \epsilon) \\ &= \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}}_{\mathbf{I}} \mathbf{w}^* + \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T}_{\text{OLS}} \epsilon \\ &= \mathbf{w}^* + \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T}_{\text{OLS}} \epsilon\end{aligned}$$

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

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

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 $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$?

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\mathbf{w}^* + \epsilon) \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}\mathbf{w}^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \\ &= \mathbf{w}^* + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon \quad \rightarrow \text{Spectral Thm.}\end{aligned}$$

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

$$\hat{\mathbf{w}} - \mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \epsilon.$$

Error in Regression

Eigendecomposition perspective

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

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

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

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

$$\boldsymbol{\Lambda}^{-1} = \begin{bmatrix} 1/\lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 1/\lambda_d \end{bmatrix},$$

so if λ_i is small, the entries of $\hat{\mathbf{w}}$ blow up!

λ_i is small \rightarrow $1/\lambda_i$ is big.

Gradient Descent

Positive Semidefinite Matrices and Convexity



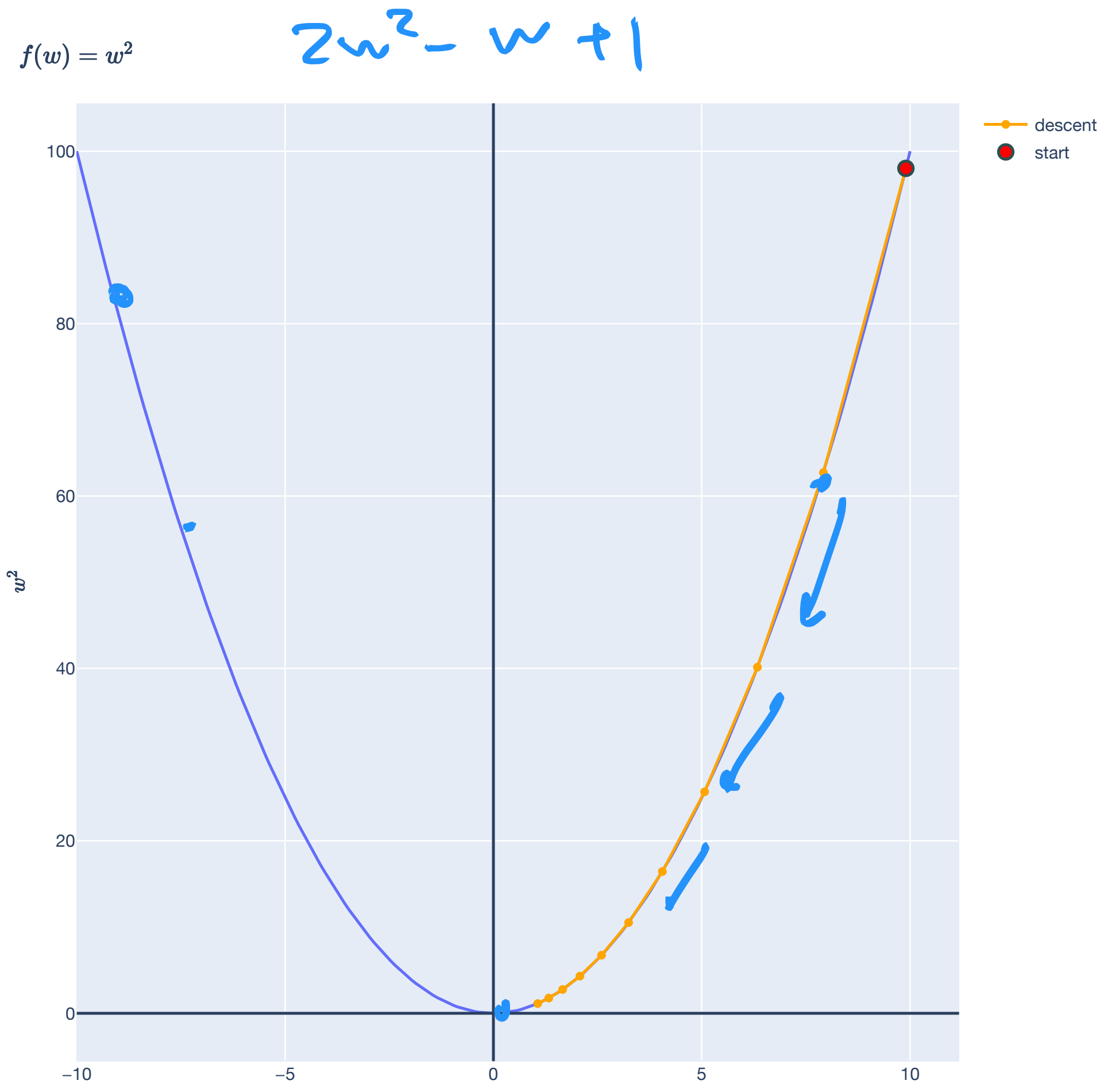
$$x^T A x$$

Lesson Overview

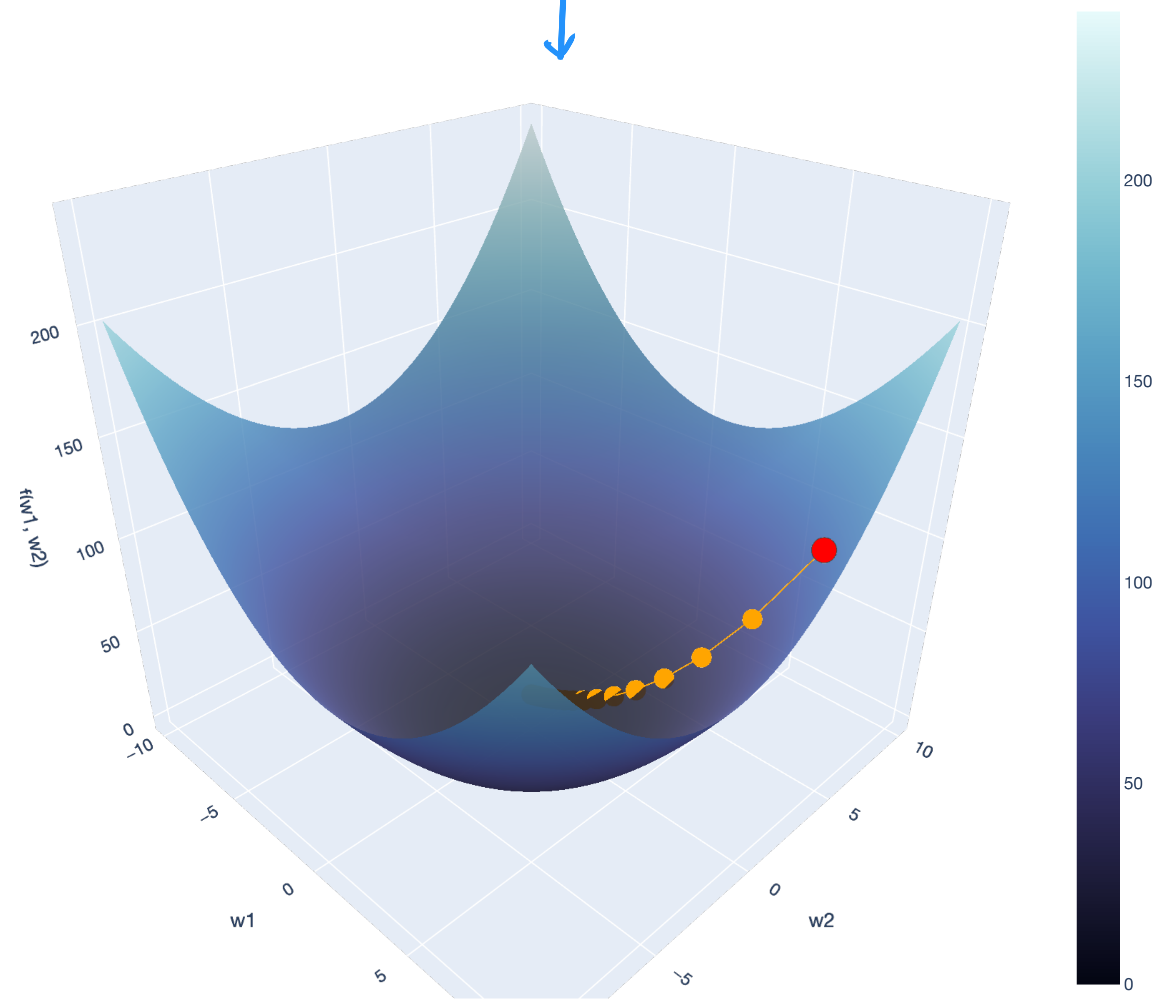
Big Picture: Gradient Descent

$$f(w_1, w_2) = \|Xw - \gamma\|^2$$

fixed X
fixed γ .



$$f(w) = w^2$$

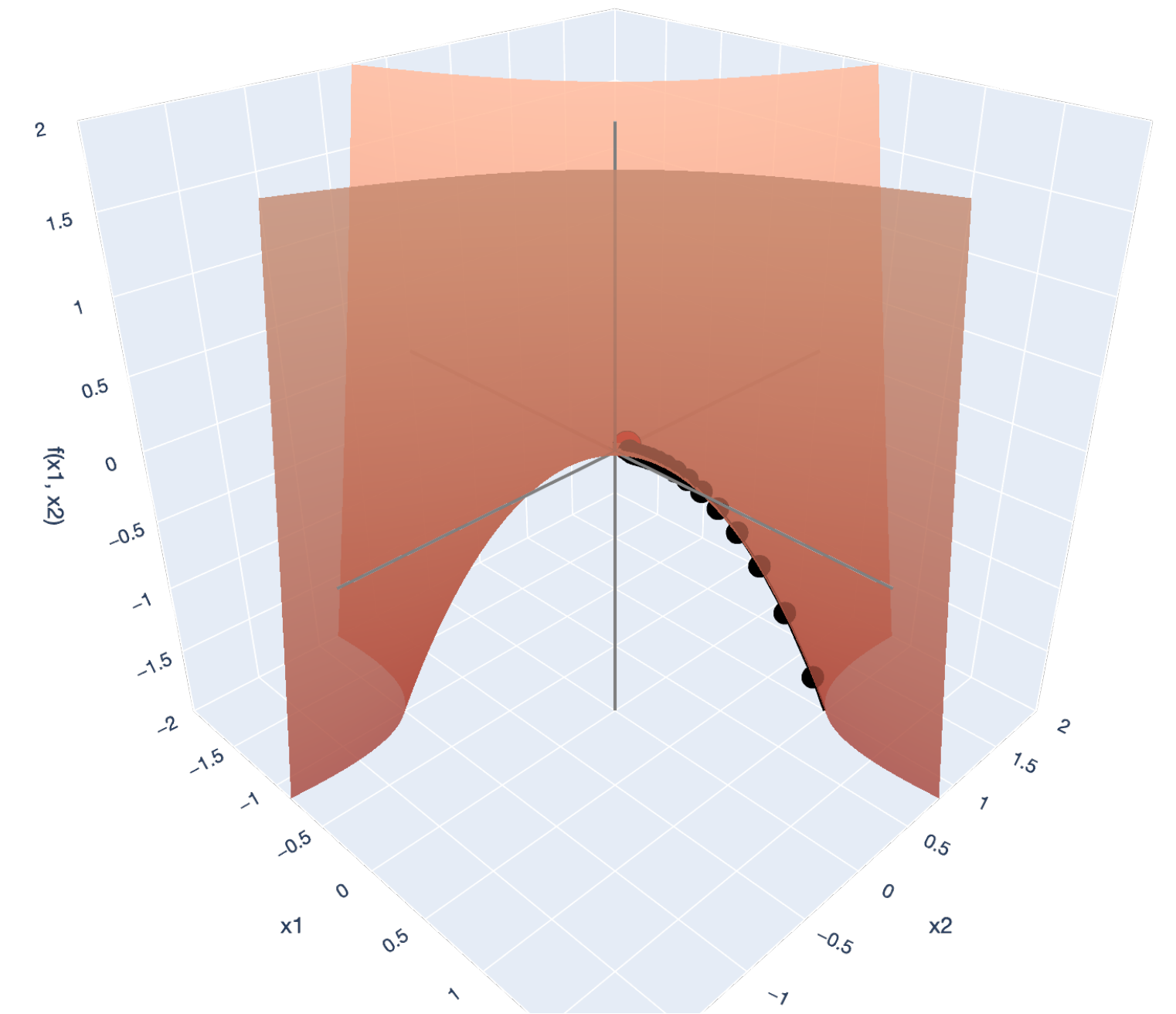
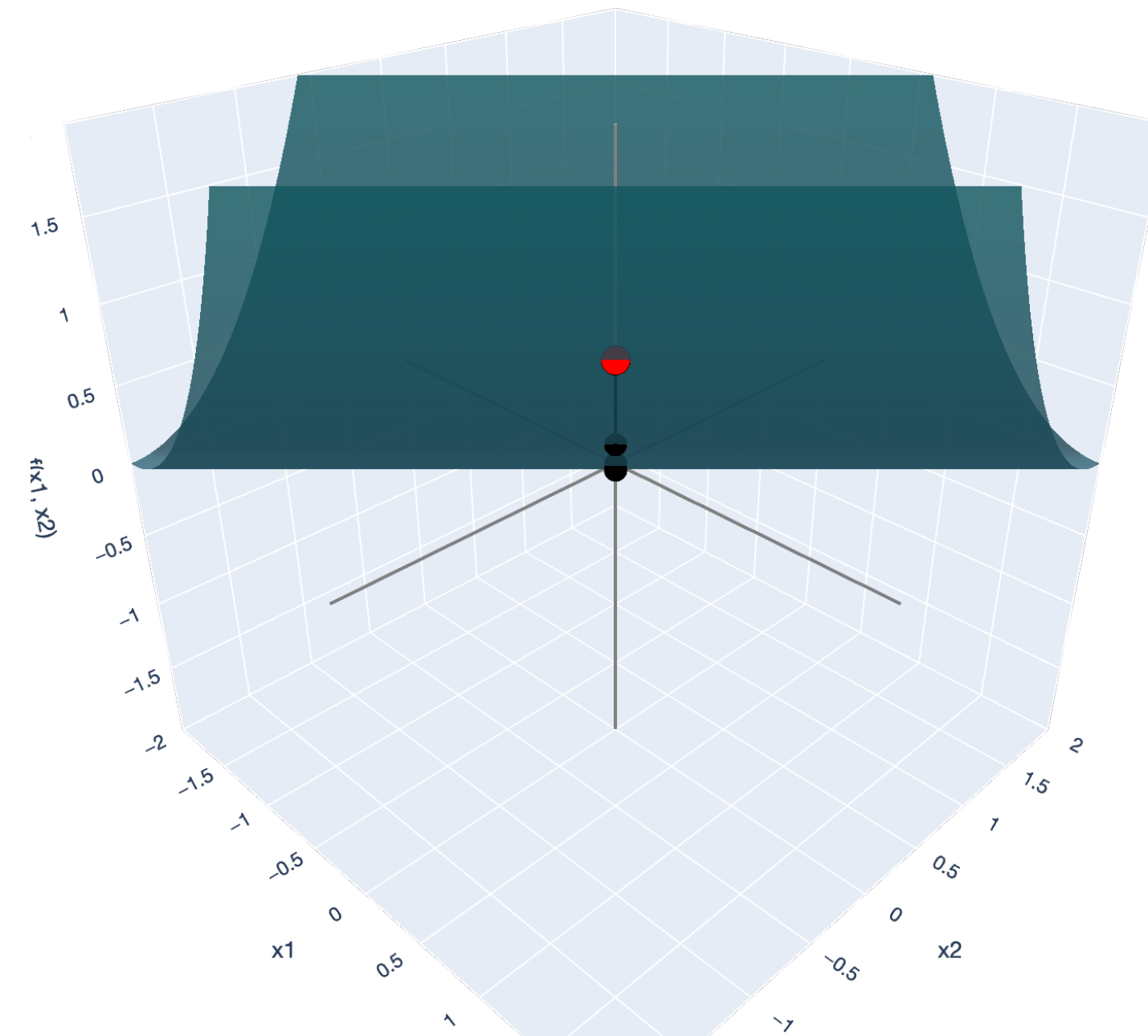
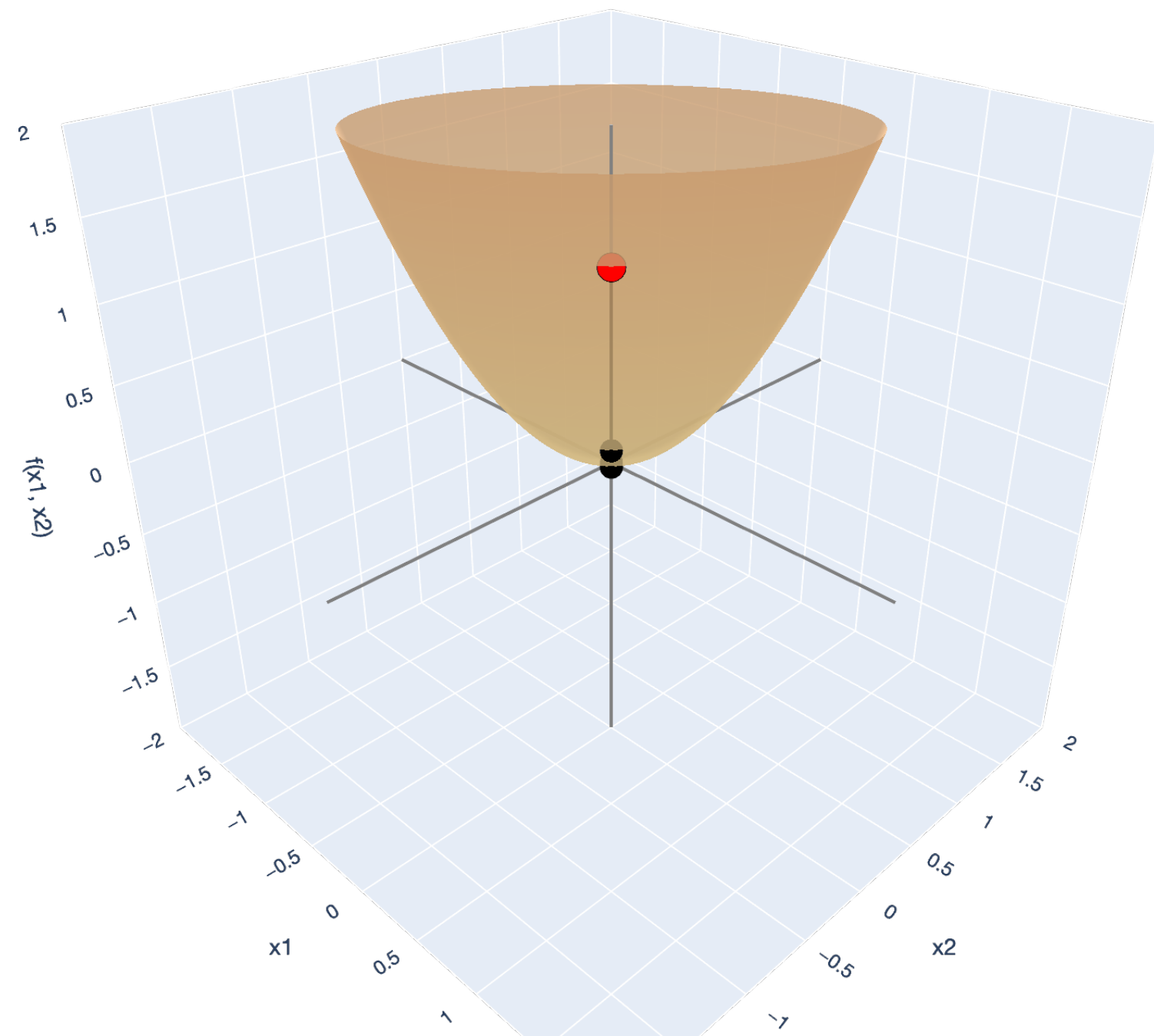


—●— descent ● start

[Click to interact](#)

Lesson Overview

Big Picture: Gradient Descent



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

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

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

Quadratic Forms

2D Example

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

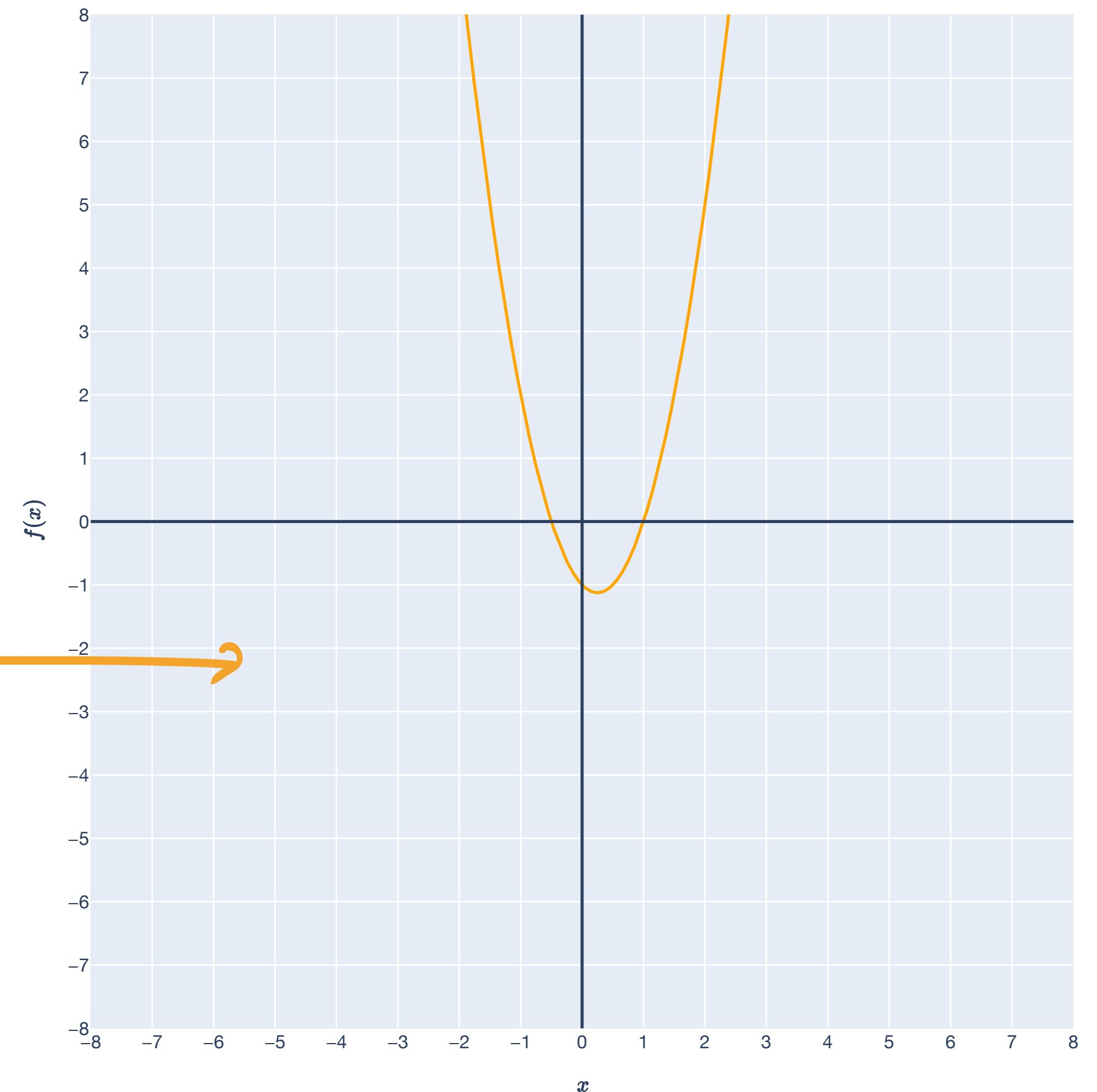
$$f(x) = ax^2 + bx + c,$$

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

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

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

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



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

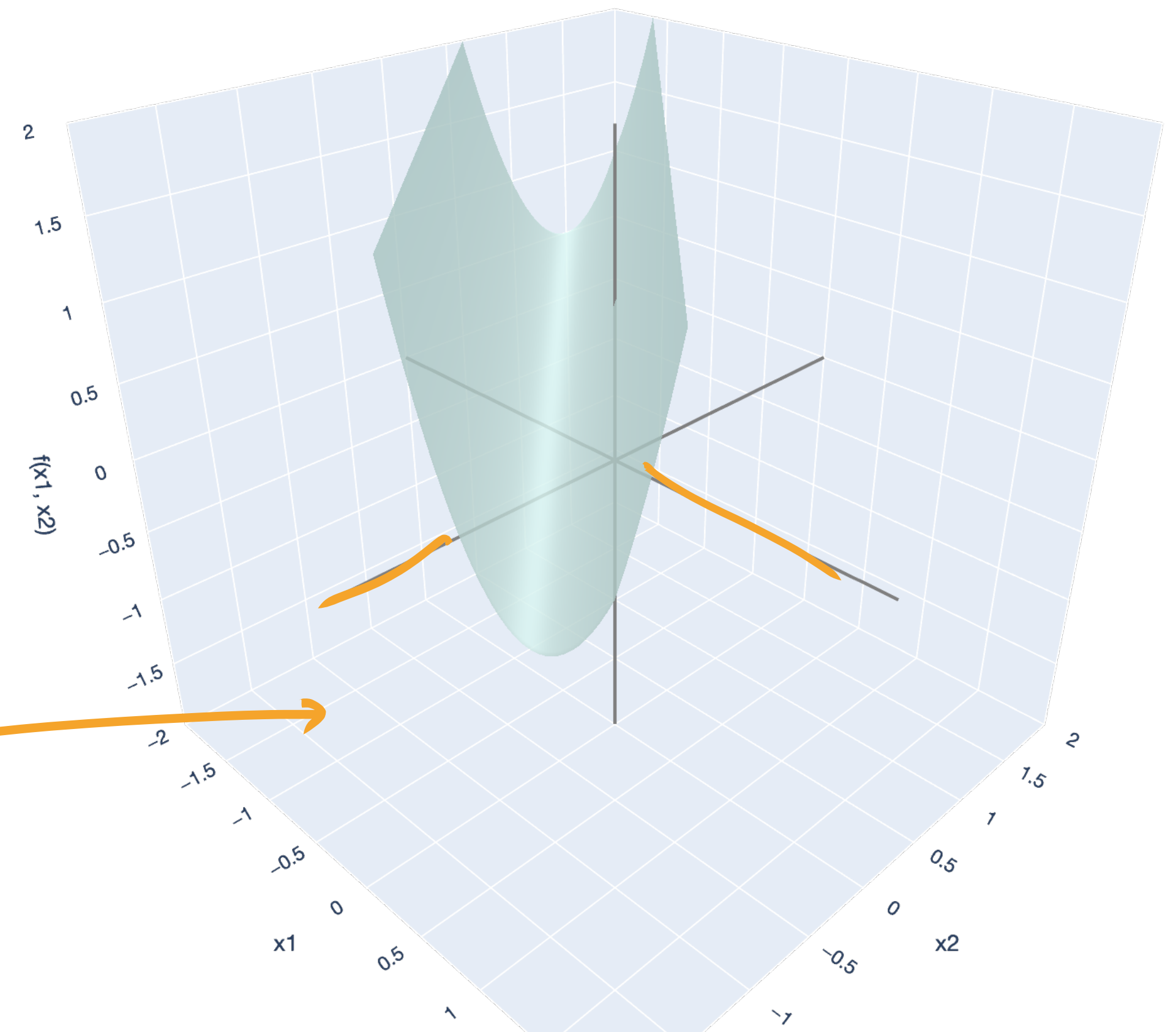
In 3D, a *quadratic function* $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ has the form

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

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

Example:

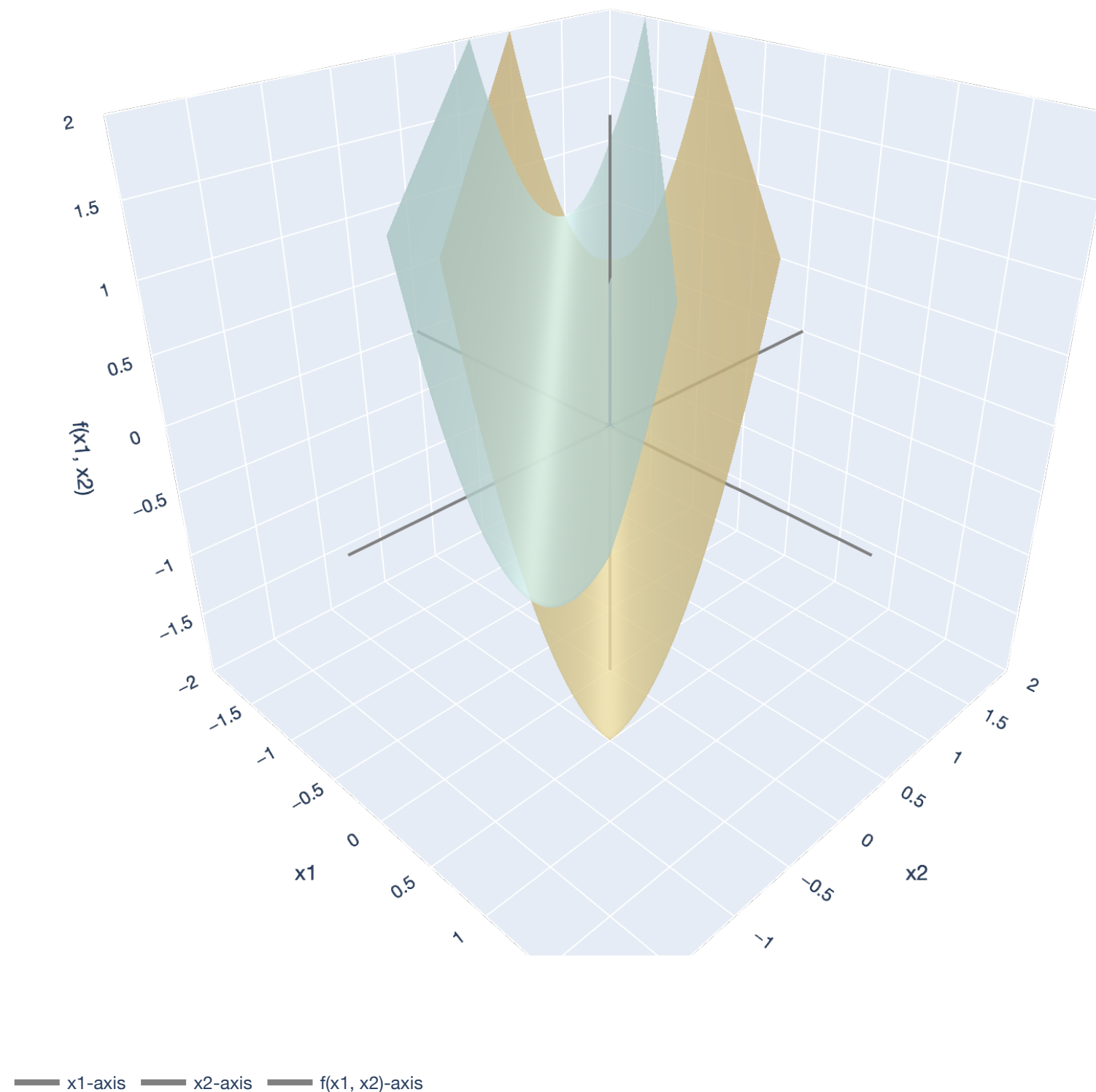
$$f(x) = 2x^2 + 4xy + 2y^2 + 2x + 2y + 1$$



Quadratic Forms

3D Example

$$f(x) = \underbrace{2x^2 + 4xy + 2y^2}_{\text{quadratic form}} + \underbrace{2x + 2y + 1}_{\text{linear form}} \text{ vs. } \underbrace{f(x) = 2x^2 + 4xy + 2y^2}_{\text{quadratic form}}$$



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!

$$\underbrace{f(x) = ax^2 + 2bxy + cy^2}_{\text{Quadratic Form}} .$$

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:

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

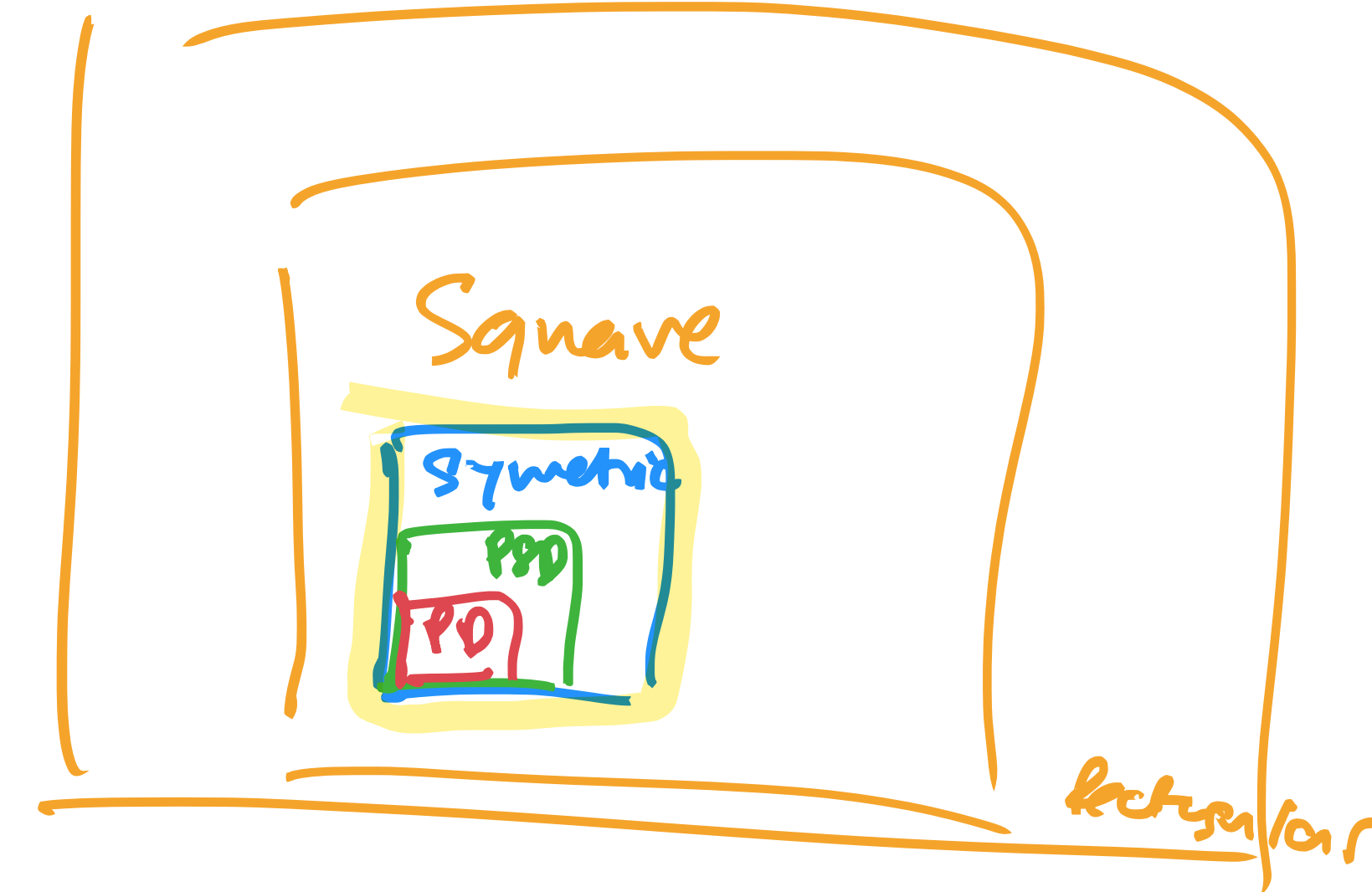
$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$$

Symmetric

Spectral
Thm
 \implies

Diagonalizable

Eigenvalue & Eigen vectors



Quadratic Forms

Relationship with matrices and eigenvalues

Consider a quadratic form:

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

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

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

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

Quadratic Forms

Relationship with matrices and eigenvalues

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

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

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

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

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

$$\boxed{\bar{\mathbf{x}}^T \mathbf{\Lambda} \bar{\mathbf{x}}}$$

$$\begin{aligned} \mathbf{x} &= M_1 \mathbf{v}_1 + M_2 \mathbf{v}_2 \\ \mathbf{x} &= \begin{bmatrix} v_1 & v_2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} M_1 \\ M_2 \end{bmatrix} \\ &\downarrow \\ \mathbf{Q}^T &= \mathbf{Q}^{-1} \\ &\downarrow \\ \bar{\mathbf{x}} &= \mathbf{Q}^T \mathbf{x} \end{aligned}$$

Quadratic Forms

Relationship with matrices and eigenvalues

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

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

There are three possibilities:

1. λ_1 and λ_2 are both positive (positive definite).
2. λ_1 or λ_2 is zero, and the other is positive (positive semidefinite). $\lambda_1, \lambda_2 \geq 0$.
3. λ_1 or λ_2 is negative (indefinite).

Quadratic Forms

Example: positive definite

$$\begin{bmatrix} x & y \end{bmatrix} \begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 2x - y \\ -x + 2y \end{bmatrix}$$

Example:

$$\downarrow = 2x^2 - xy - xy + 2y^2 = 2x^2 - 2xy + 2y^2$$

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

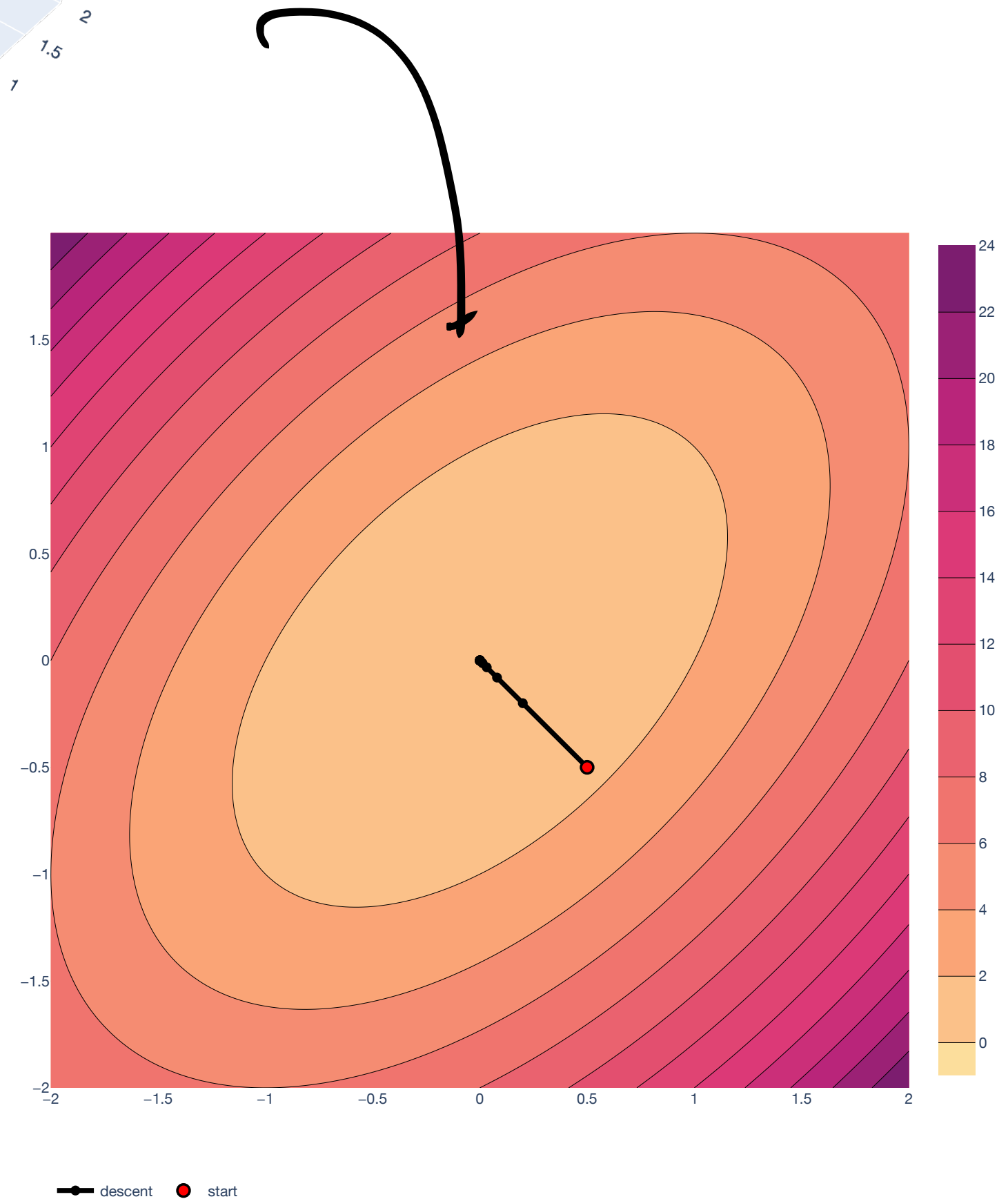
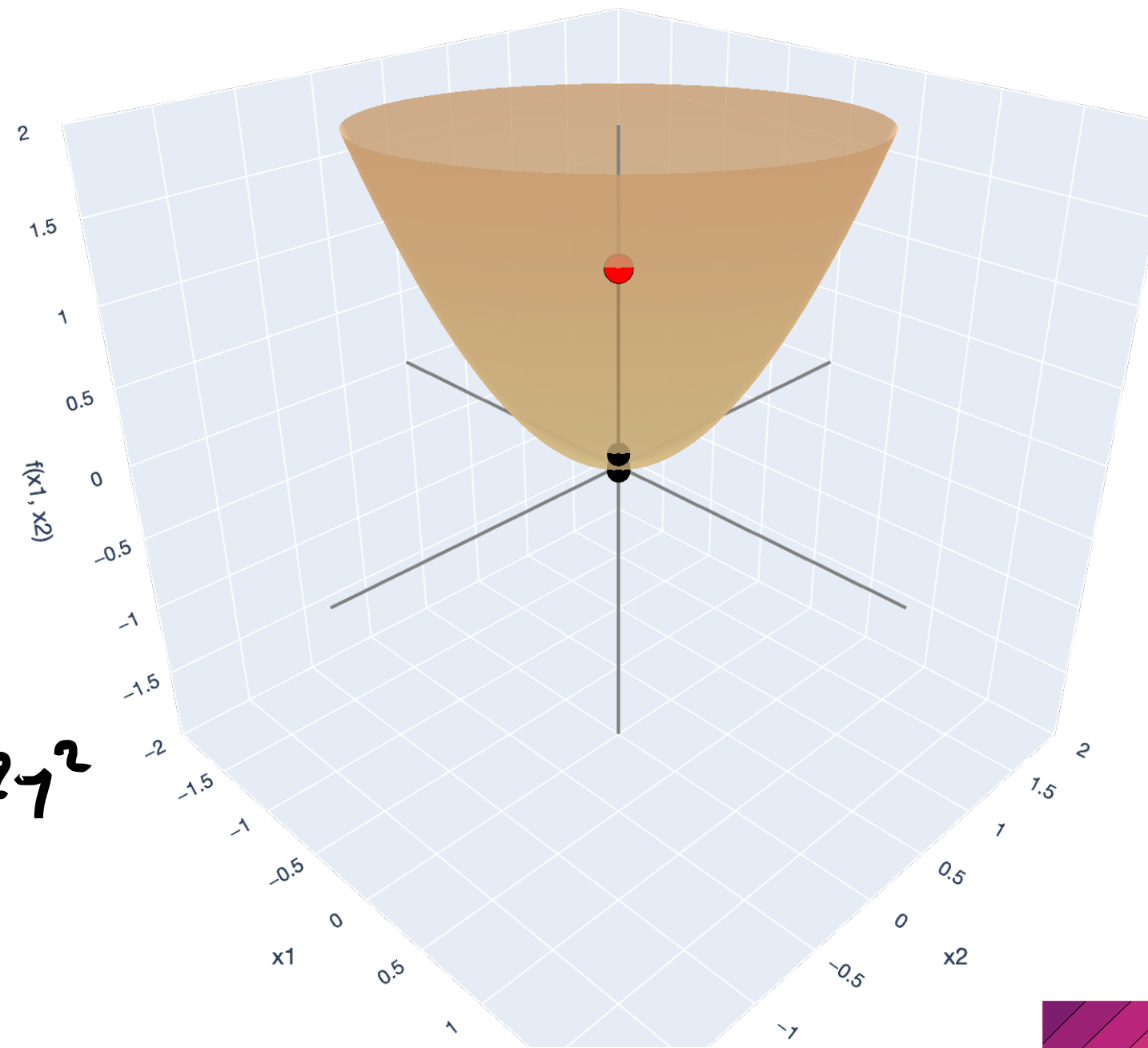
Eigendecomposition:

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

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

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

$$Q^T$$



Quadratic Forms

Example: positive semidefinite

Example:

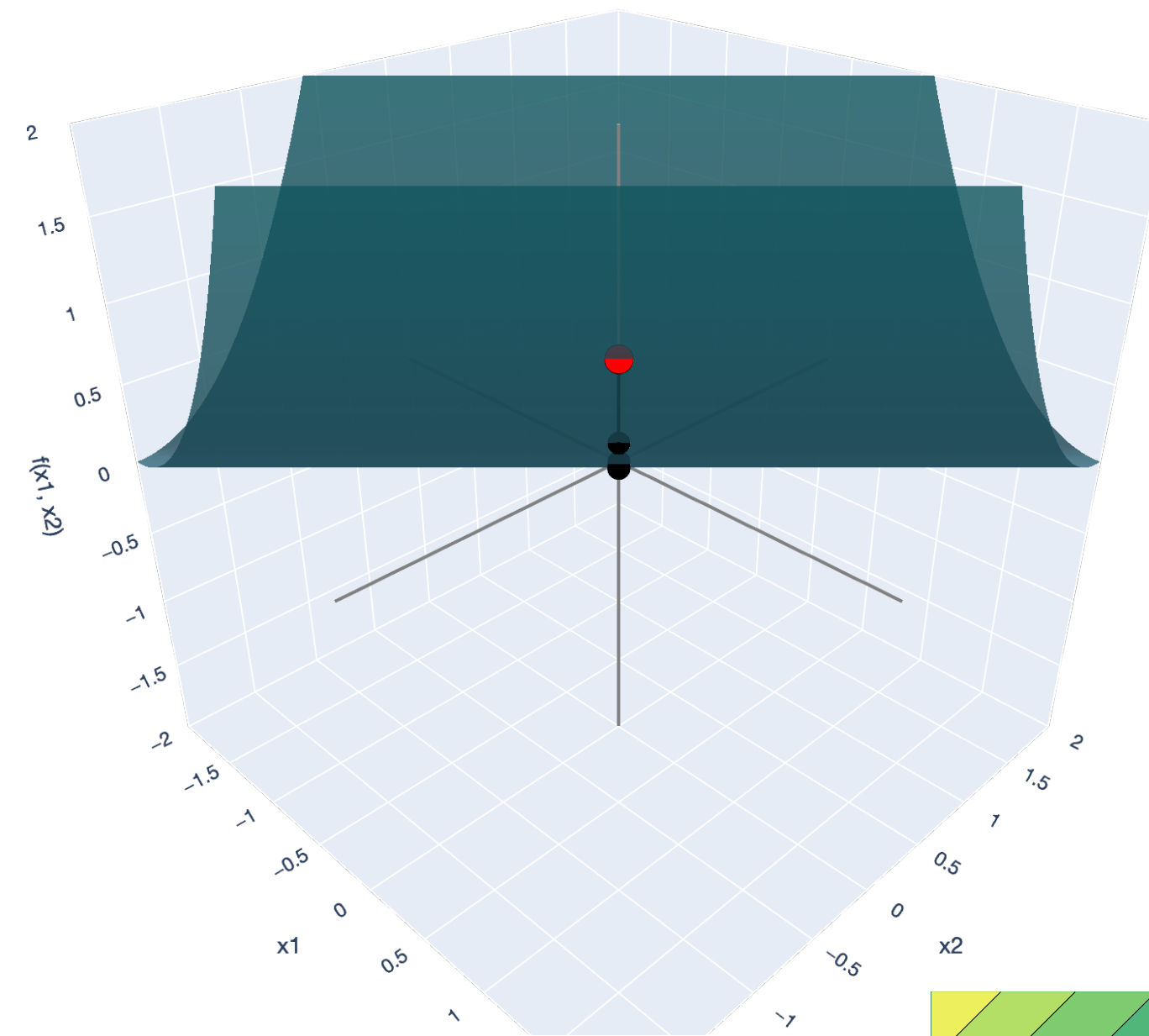
$$f(x, y) = [x \quad y] \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

Eigendecomposition:

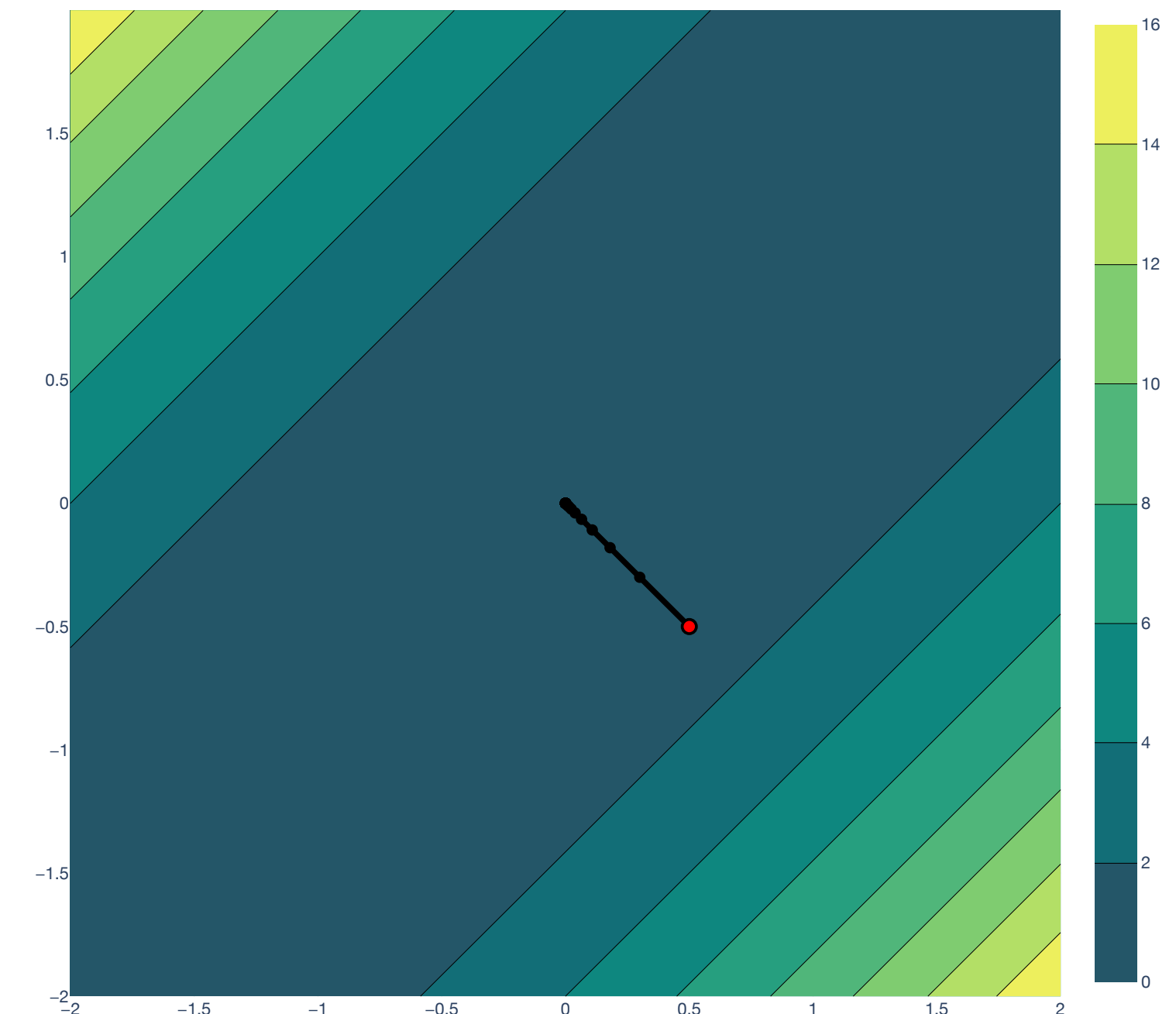
$$\begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} = \begin{bmatrix} -1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} -1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

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

$\lambda_1 = 2$
 $\lambda_2 = 0$



— x1-axis — x2-axis — f(x1, x2)-axis — descent — start



— descent — start

Quadratic Forms

Example: indefinite

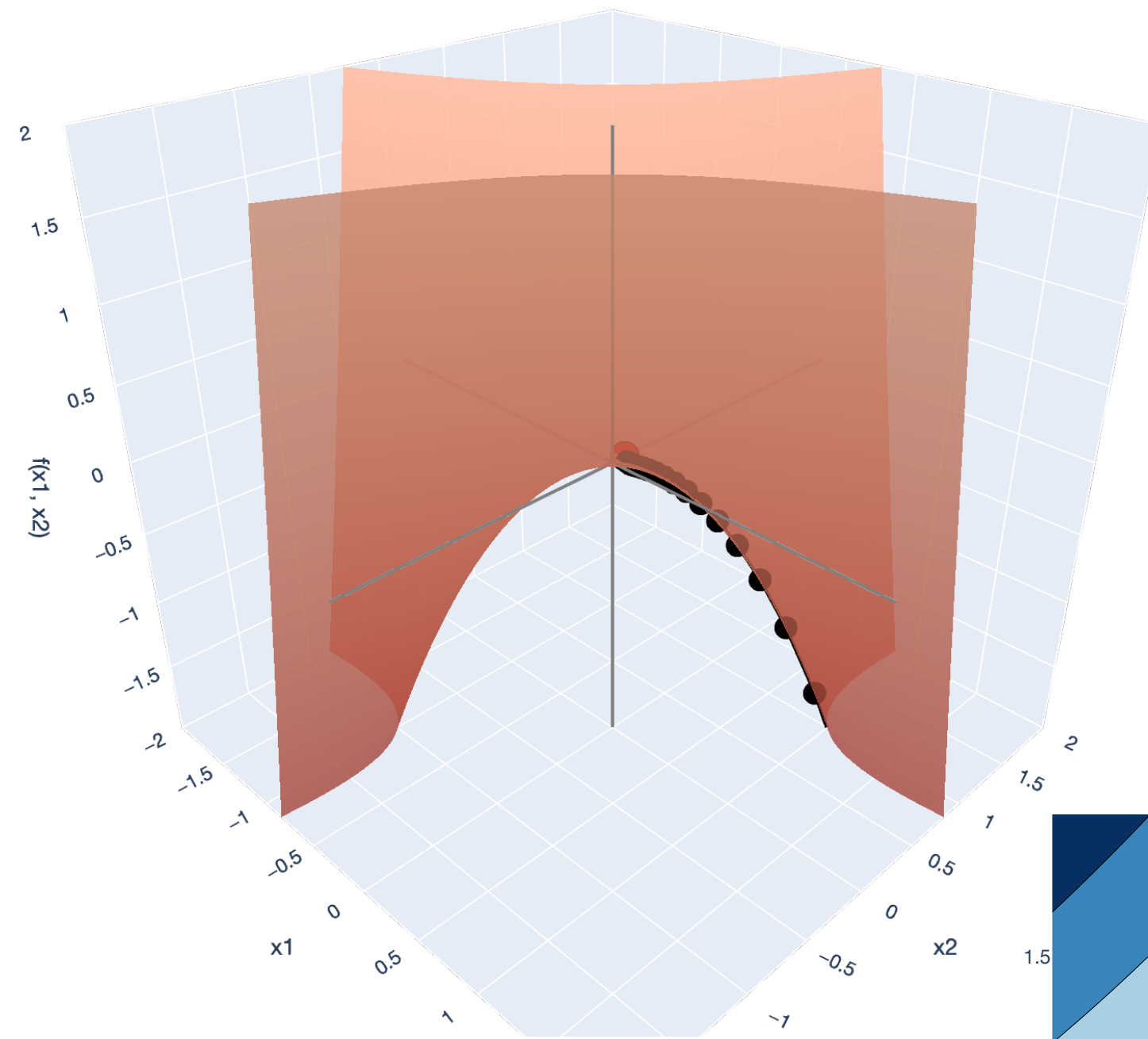
Example:

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

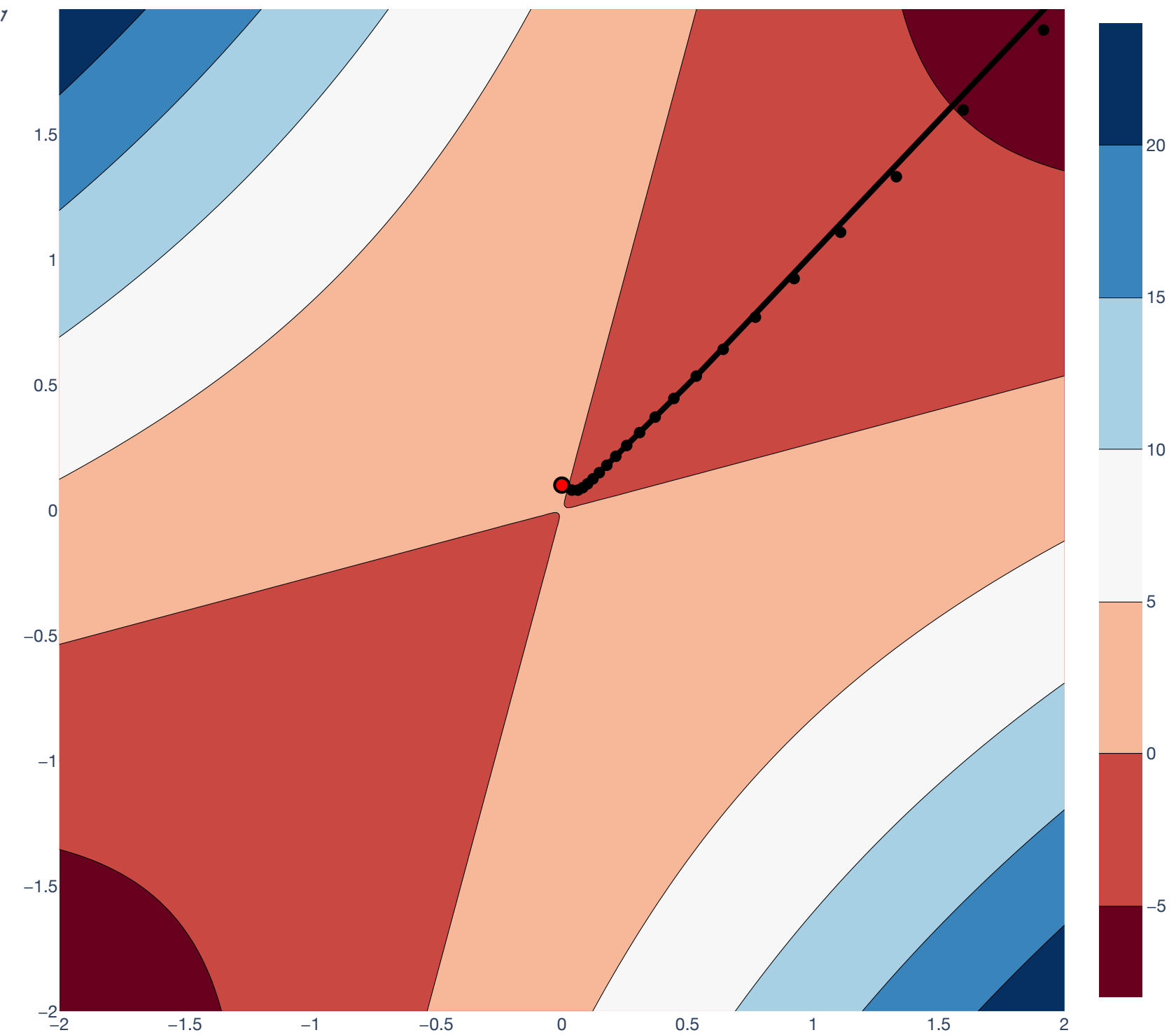
Eigendecomposition:

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

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



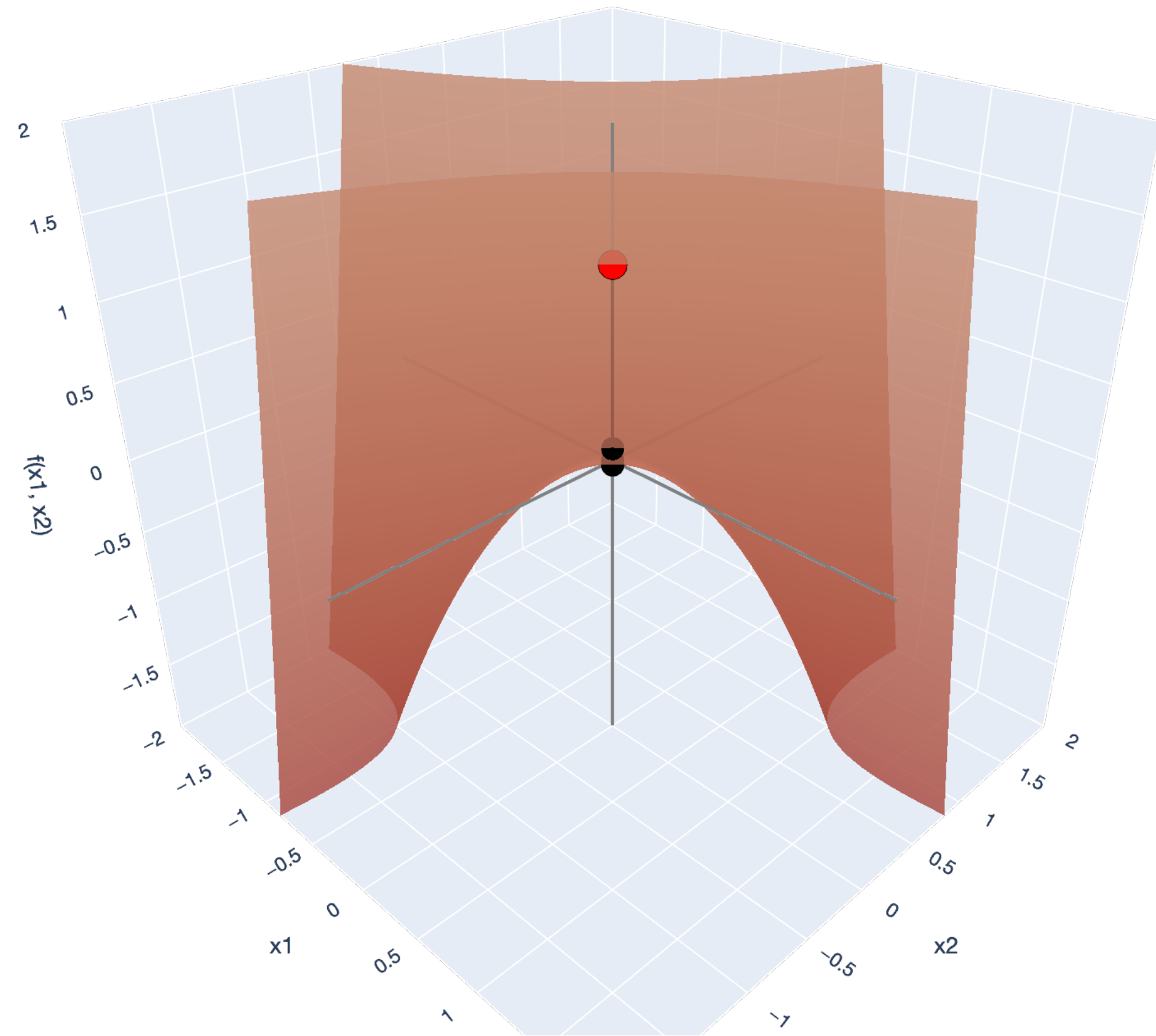
— x1-axis — x2-axis — f(x1, x2)-axis — descent ● start



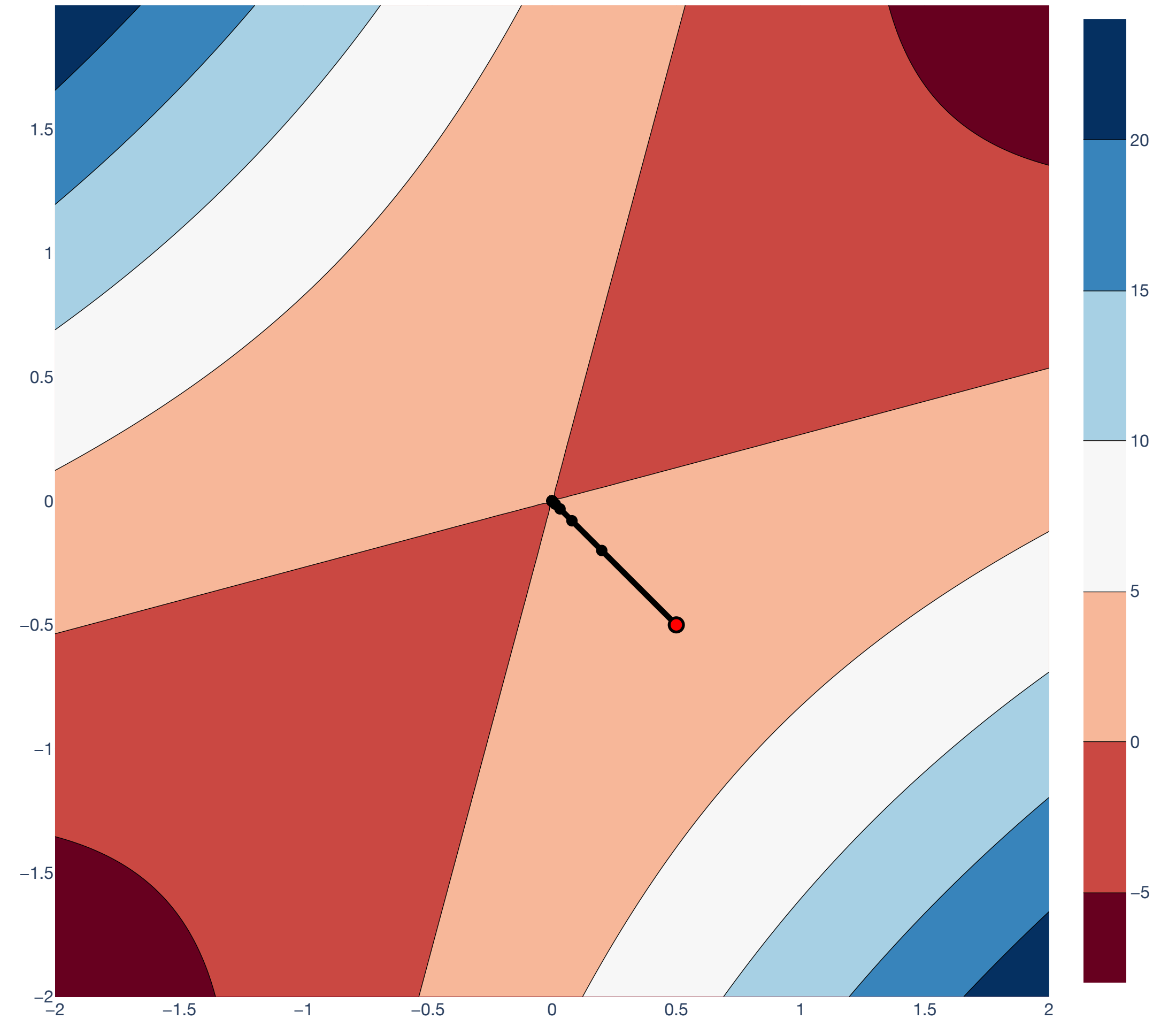
— descent ● start

Quadratic Forms

Example: indefinite



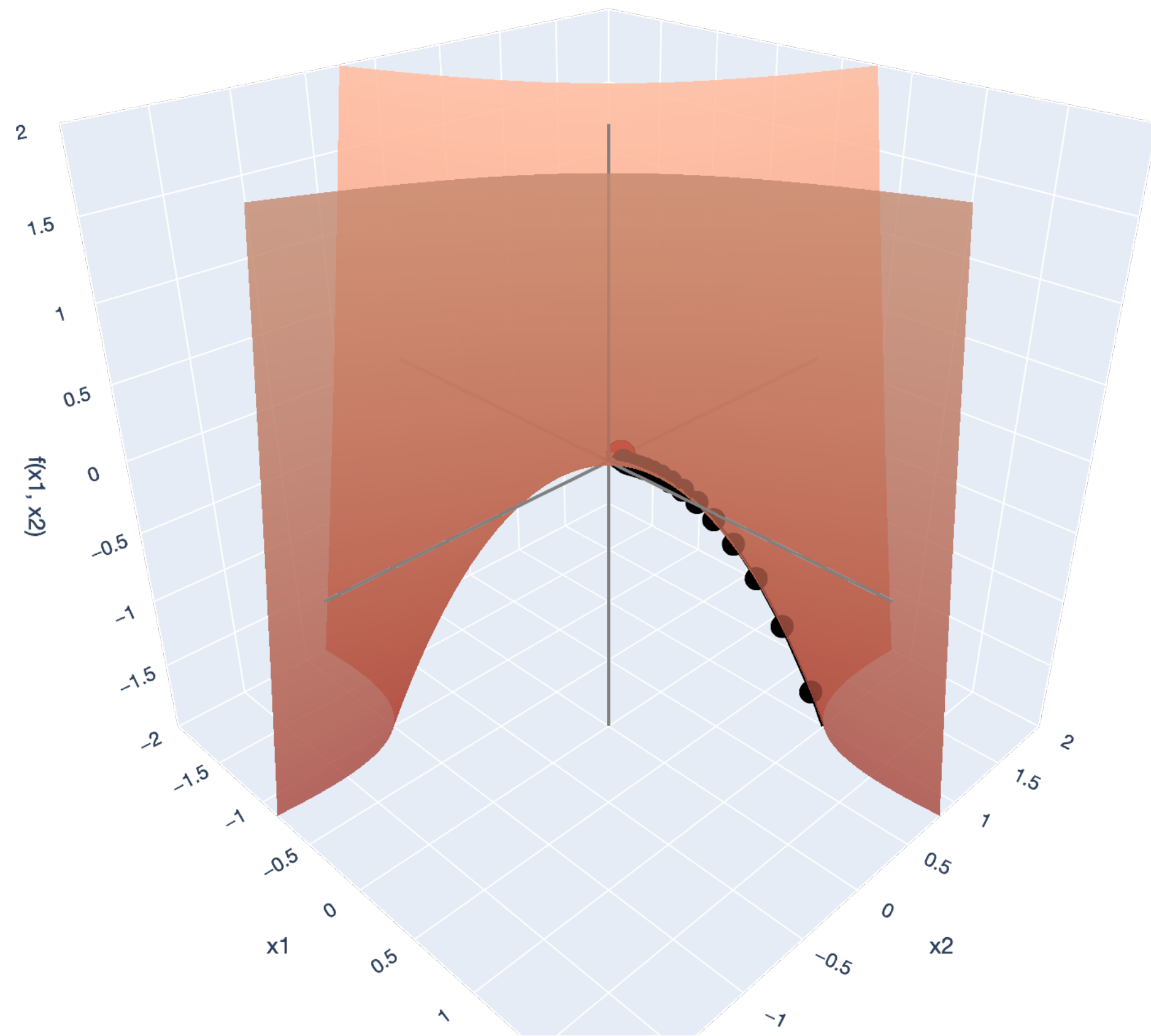
— x1-axis — x2-axis — f(x1, x2)-axis —●— descent ● start



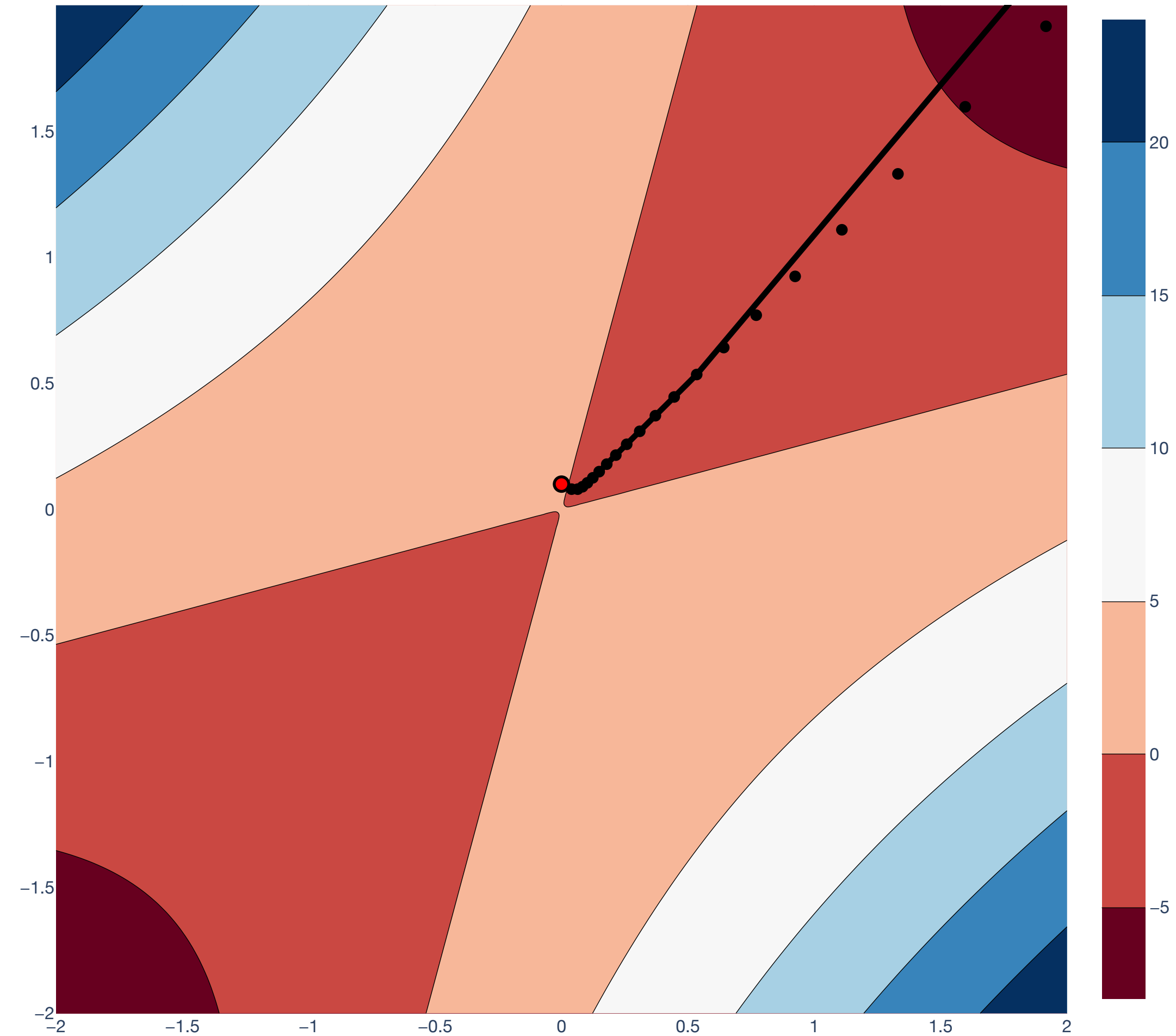
—●— descent ● start

Quadratic Forms

Example: indefinite



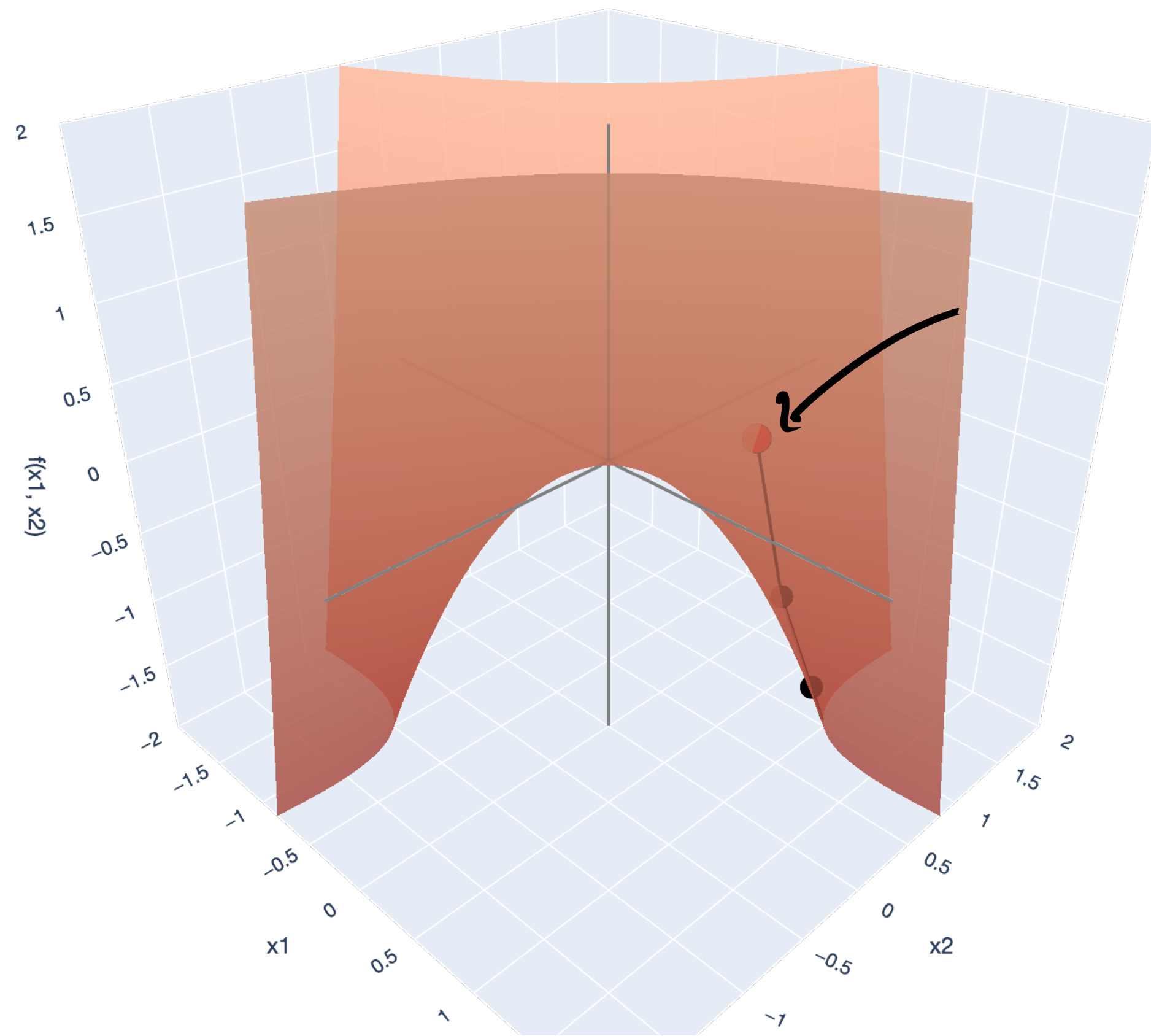
— x1-axis — x2-axis — f(x1, x2)-axis —●— descent ● start



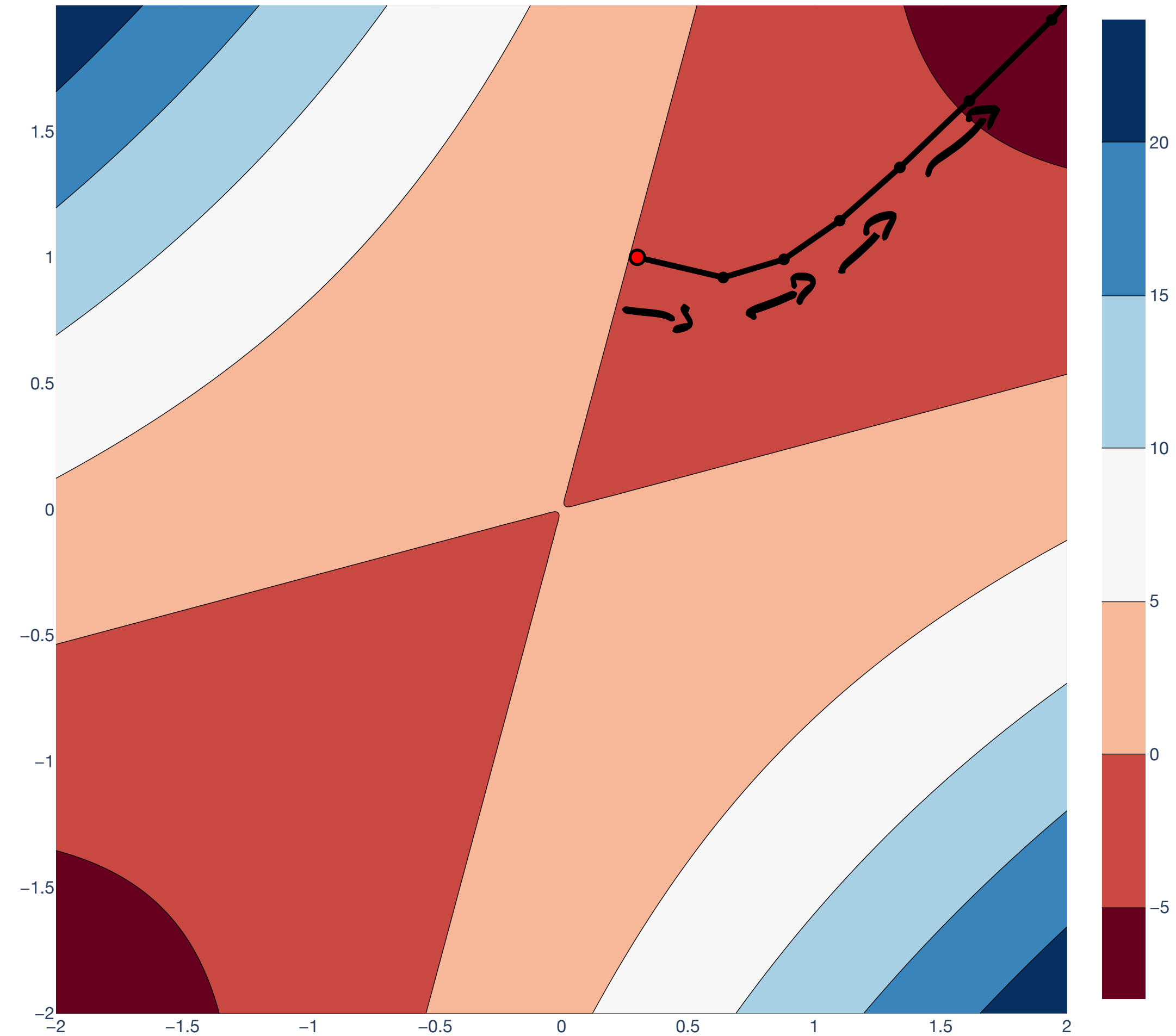
—●— descent ● start

Quadratic Forms

Example: indefinite



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



—●— descent ● start

Least Squares

Example of quadratic form

Consider the familiar function we've been thinking about:

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

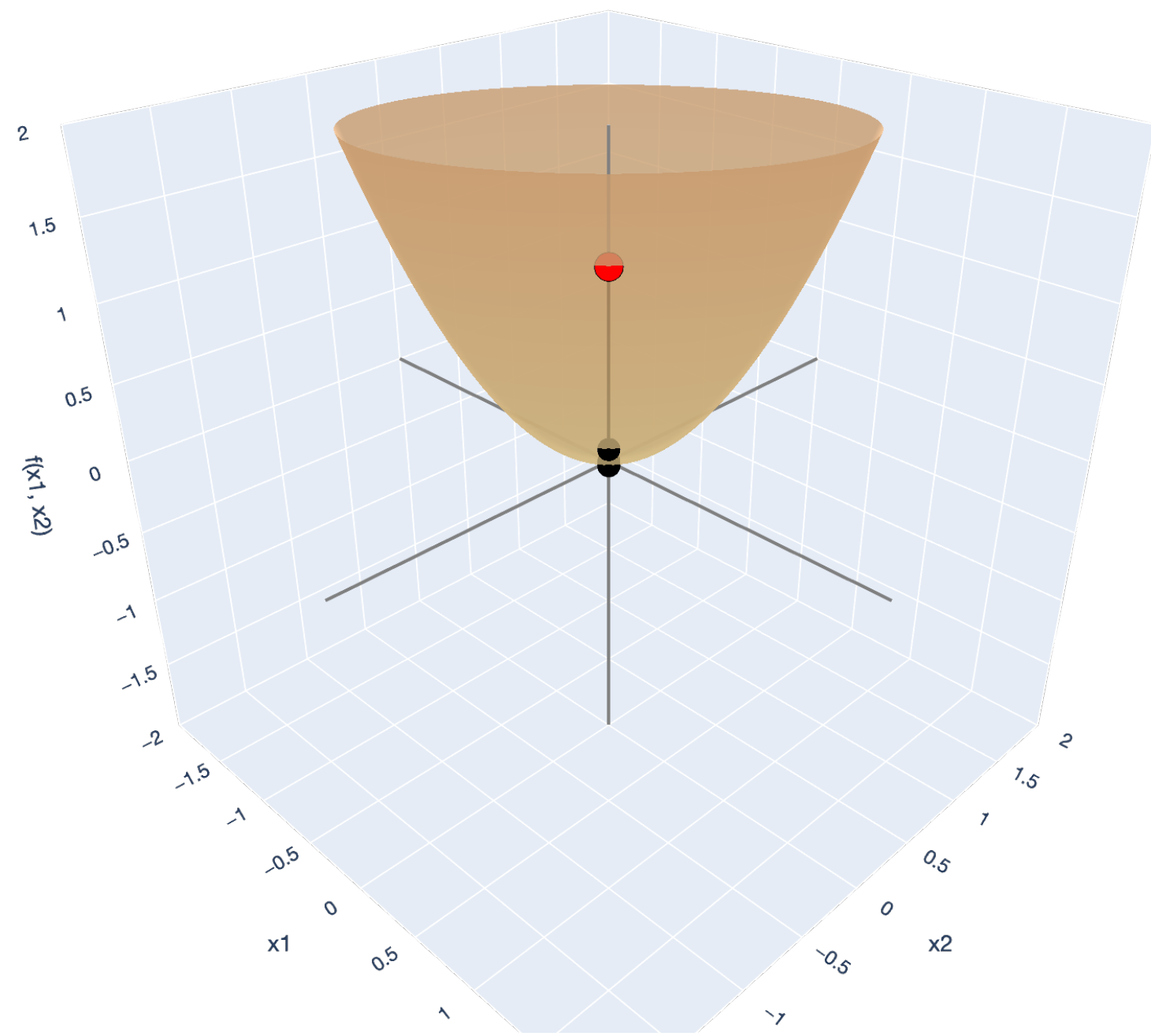
$$(\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) = \mathbf{w}^\top (\mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\mathbf{w}^\top (\mathbf{X}^\top \mathbf{y}) + \mathbf{y}^\top \mathbf{y}.$$

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

$$A = \mathbf{X}^\top \mathbf{X} \text{ is } \underline{\underline{\text{PSD}}},$$

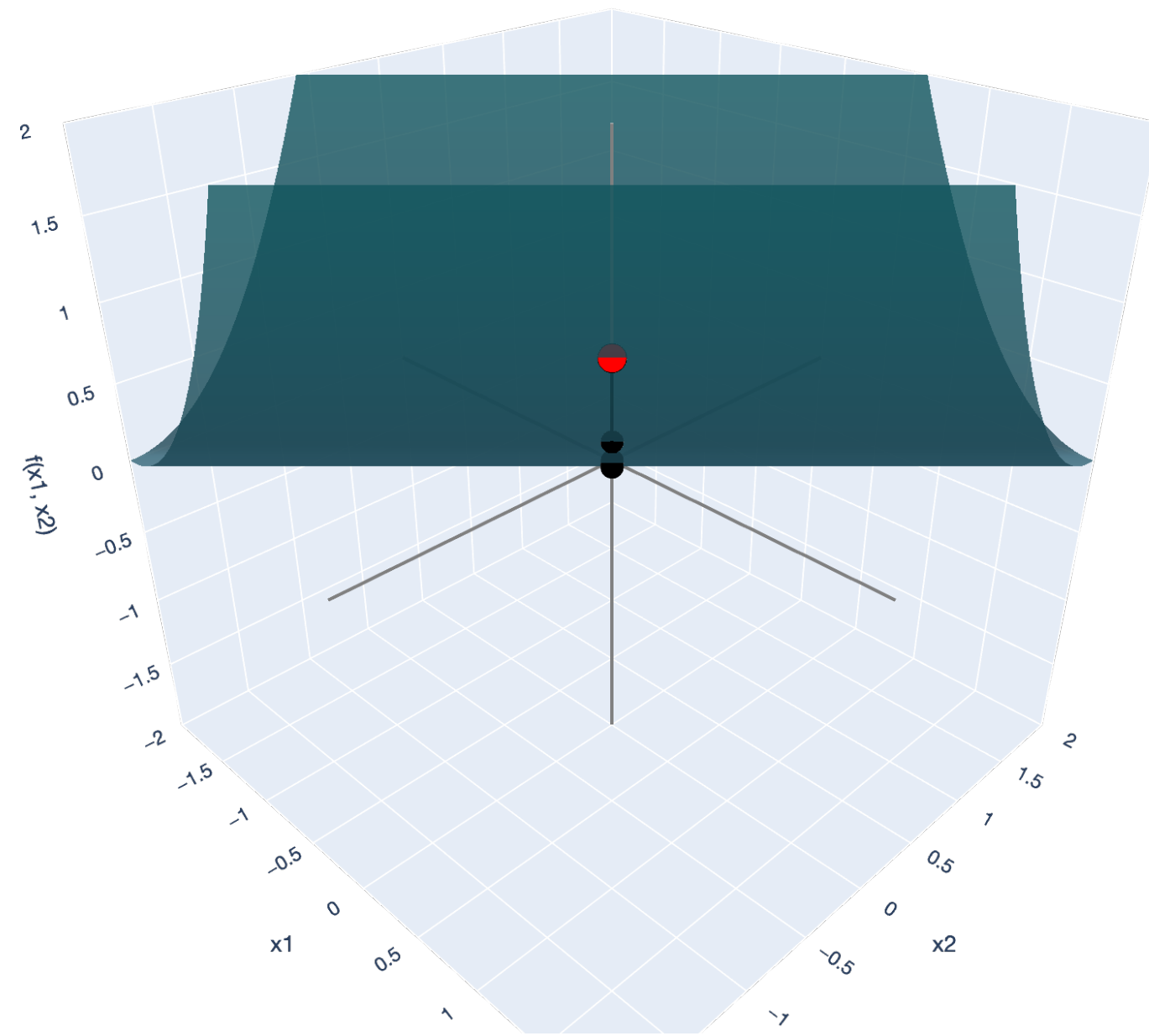
Gradient Descent

Preview



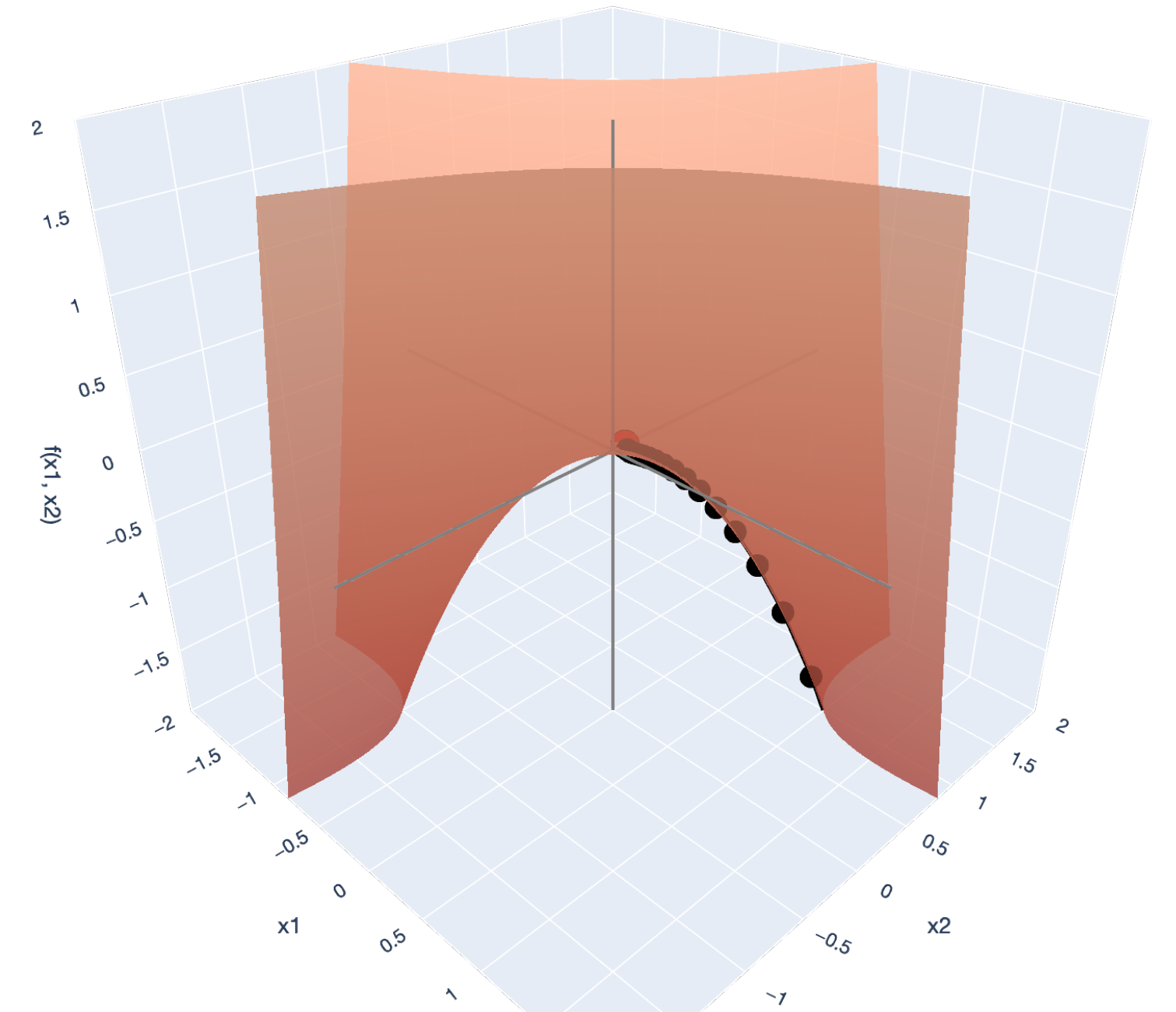
— x1-axis — x2-axis — f(x1, x2)-axis —● descent ● start

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



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

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



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

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

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.

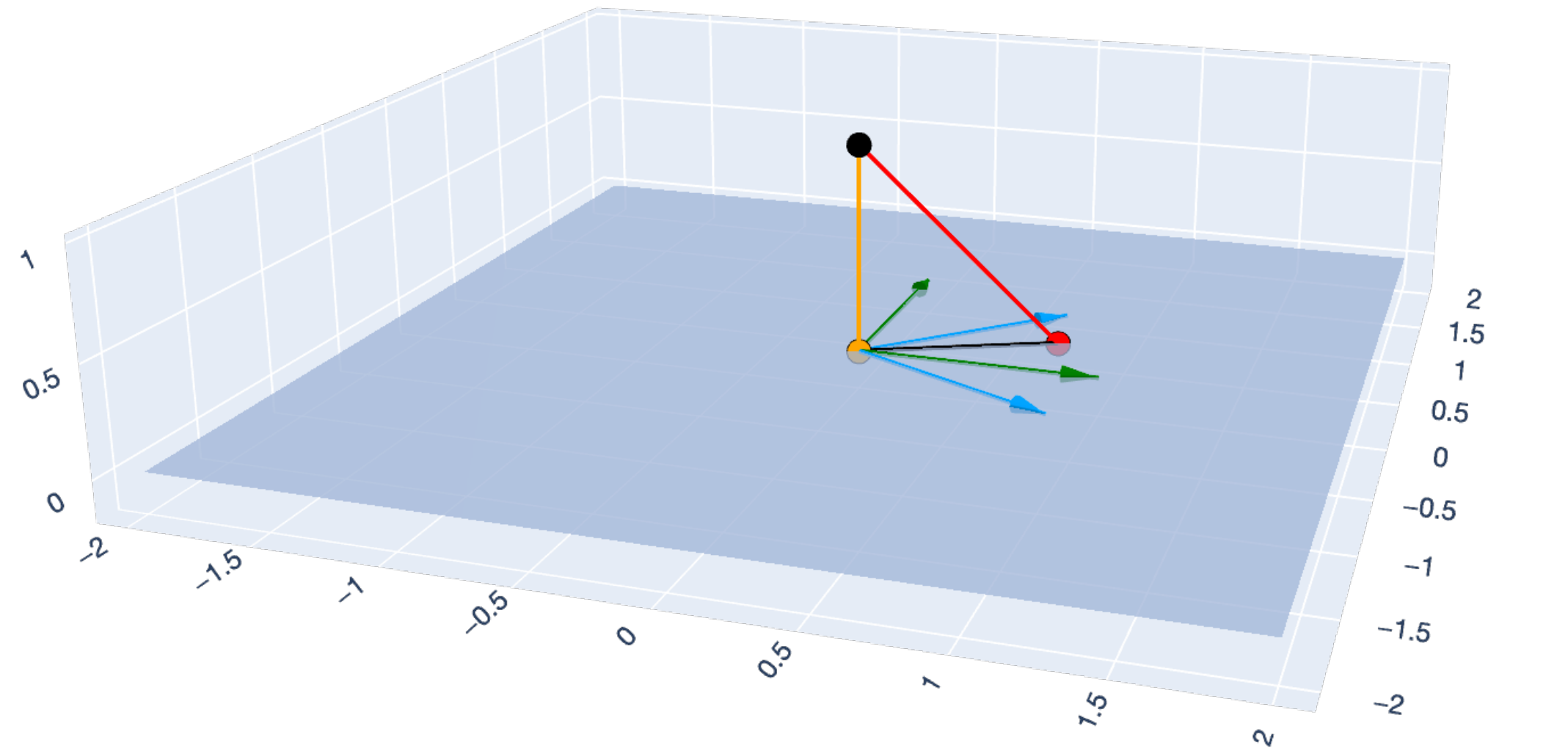
Spectral Theorem. Symmetric matrices are always diagonalizable.

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

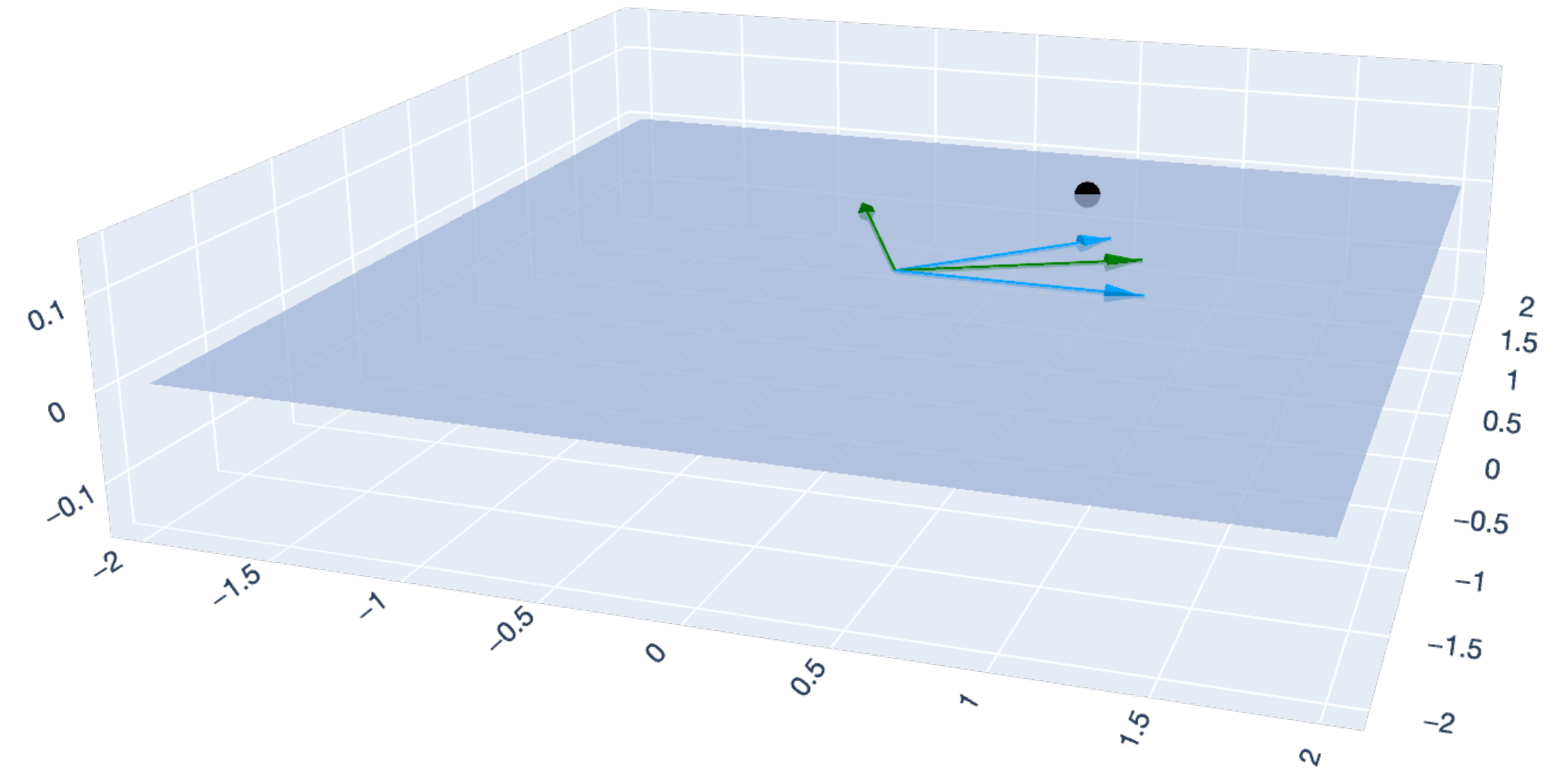
Lesson Overview

Big Picture: Least Squares

$$\hat{w} = X^+ y$$



— x1
 — x2
 — u1
 — u2
 — $y - \hat{y}$
 — $\sim y - \hat{y}$
 — $\sim y - y$
 ● y
 ● \hat{y}
 ● $\sim y$

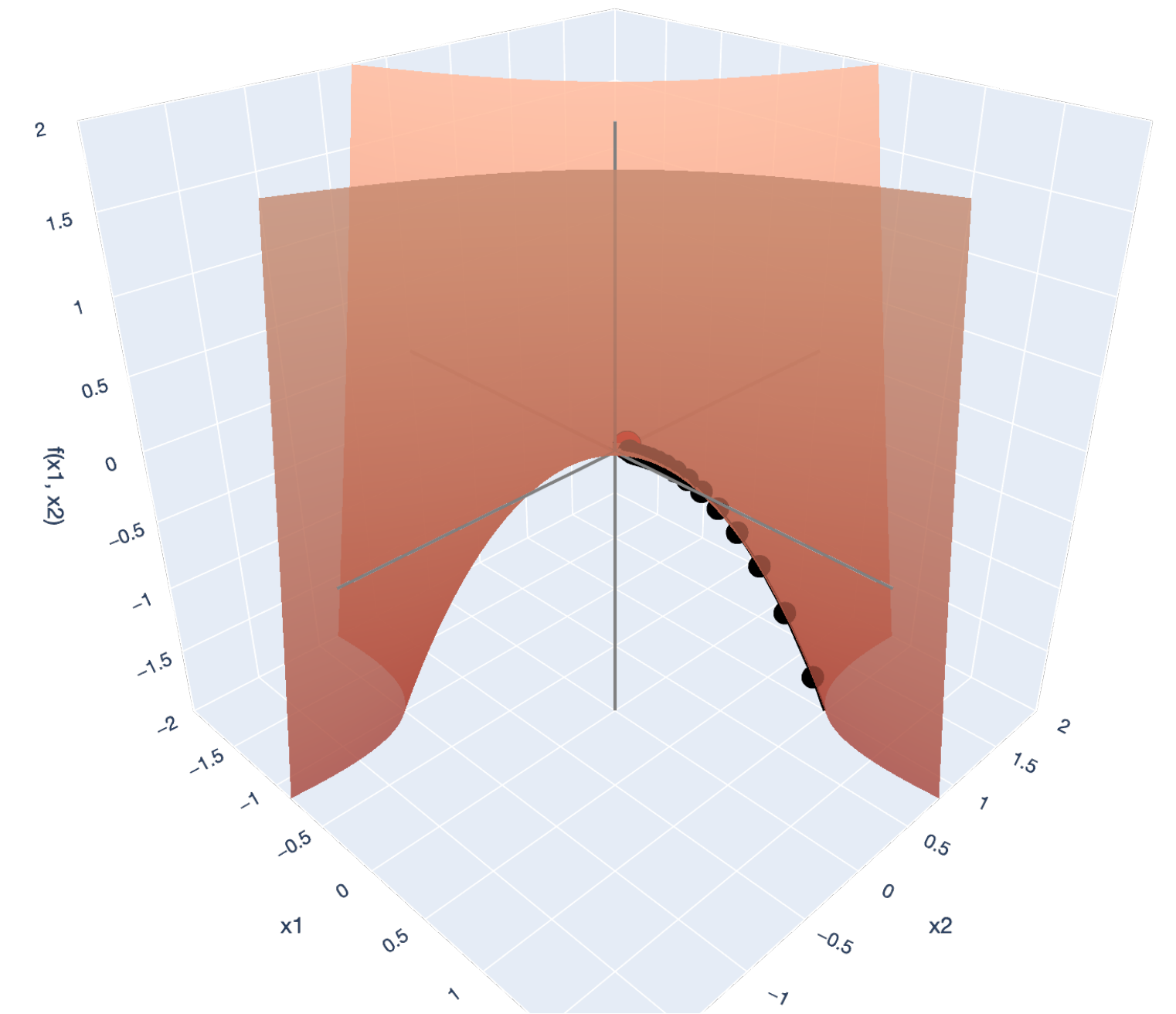
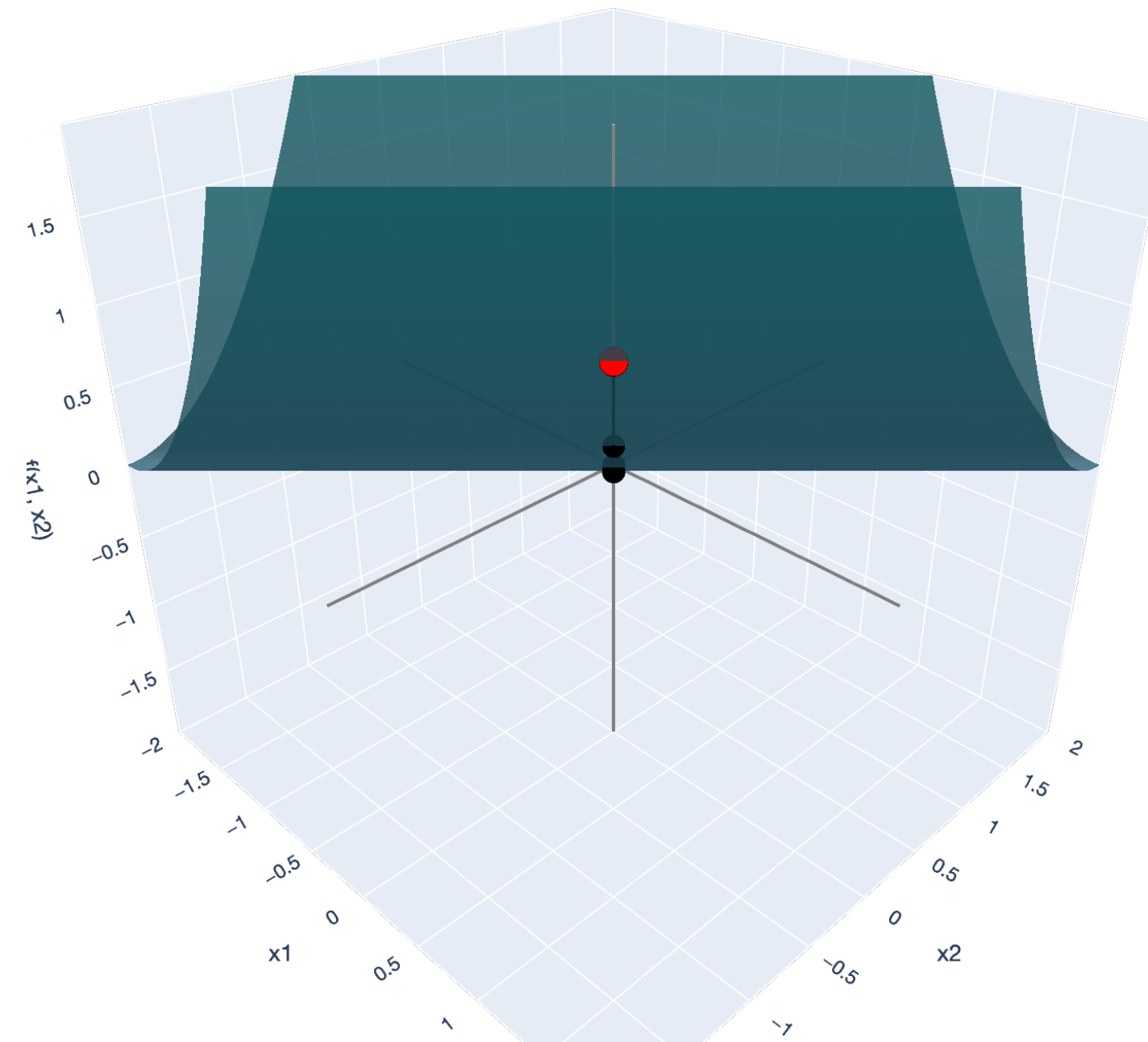
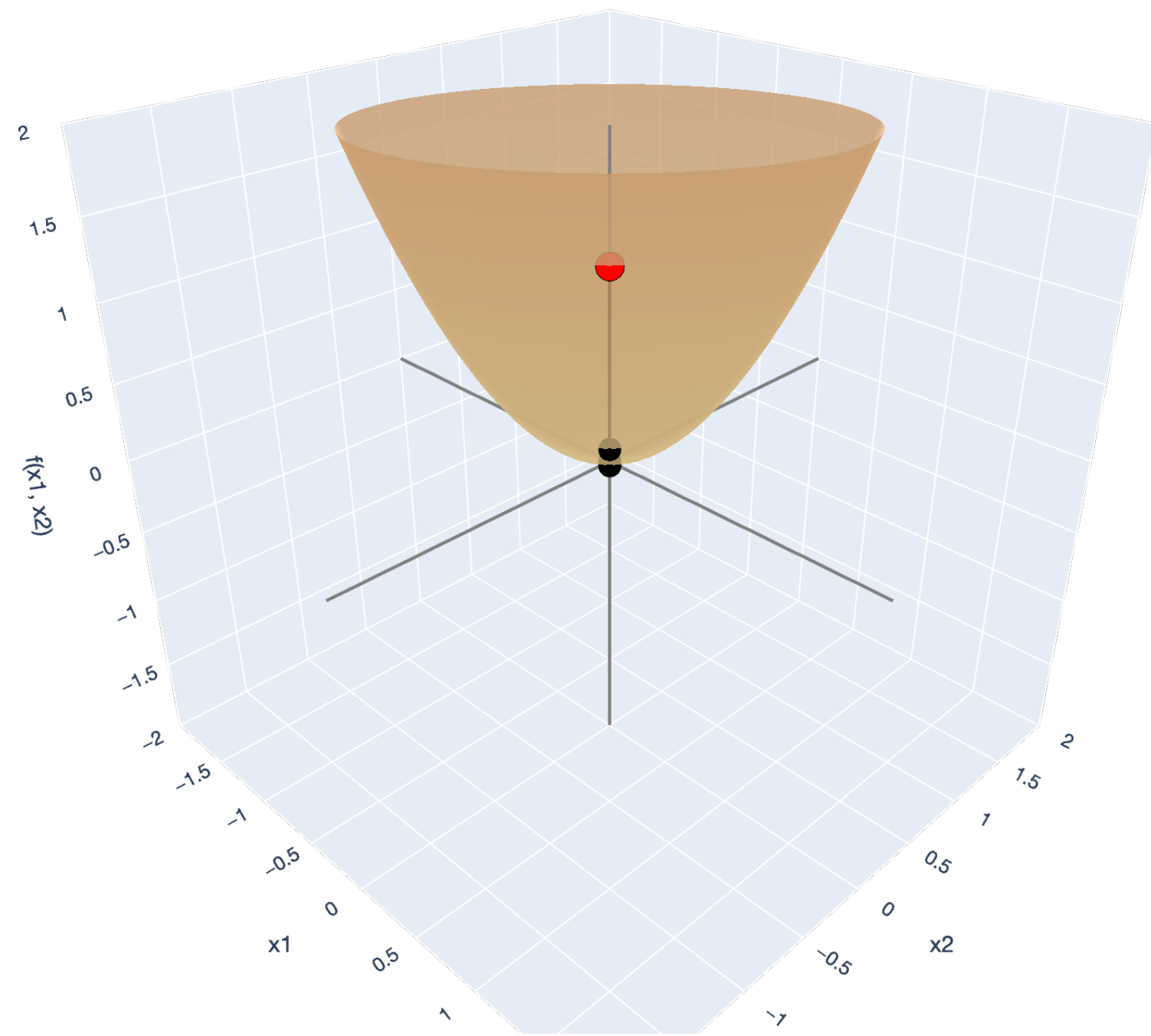


— x1
 — x2
 — u1
 — u2
 ● y

Lesson Overview

Big Picture: Gradient Descent

QUADRATIC FUNCTIONS



x1-axis x2-axis f(x1, x2)-axis descent start

x1-axis x2-axis f(x1, x2)-axis descent start

x1-axis x2-axis f(x1, x2)-axis descent start

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.