

Math for Machine Learning

Week 2.1: Singular Value Decomposition

By: Samuel Deng

Logistics & Announcements

- DVE 11:59 PM ~~ENIGHT~~ (Mon, Jul 8): Project 1st reflection.
→ Pick a paper, follow the instructions under "Project:"
- DVE 11:59 PM (Thurs.): PS1.
- PS2 released Tues. (~ Noon).

⑥
Late Pays.

BREAKS: 10 minutes / 5 minutes each.

Roll Ev. com / Sandenog.

Lesson Overview

⌘ MATRICES AS LINEAR TRANSFORMATIONS

- Orthogonal complement and properties of projection. We go over several useful properties of the projection operation.
- Derivation of the singular value decomposition (SVD). We derive the SVD from the “best-fitting subspace” problem using all the properties of projection.
- SVD Definition. We go over the definition of SVD and the geometric intuition as the factorization of a data matrix.
- Application of SVD: rank- k approximation. We state and give an example of rank- k approximation, a common data compression technique using SVD.
- Pseudoinverse. We unify our OLS solution from the perspective of SVD and the notion of the pseudoinverse, a generalization of inverses to rectangular matrices.

Lesson Overview

Big Picture: Least Squares

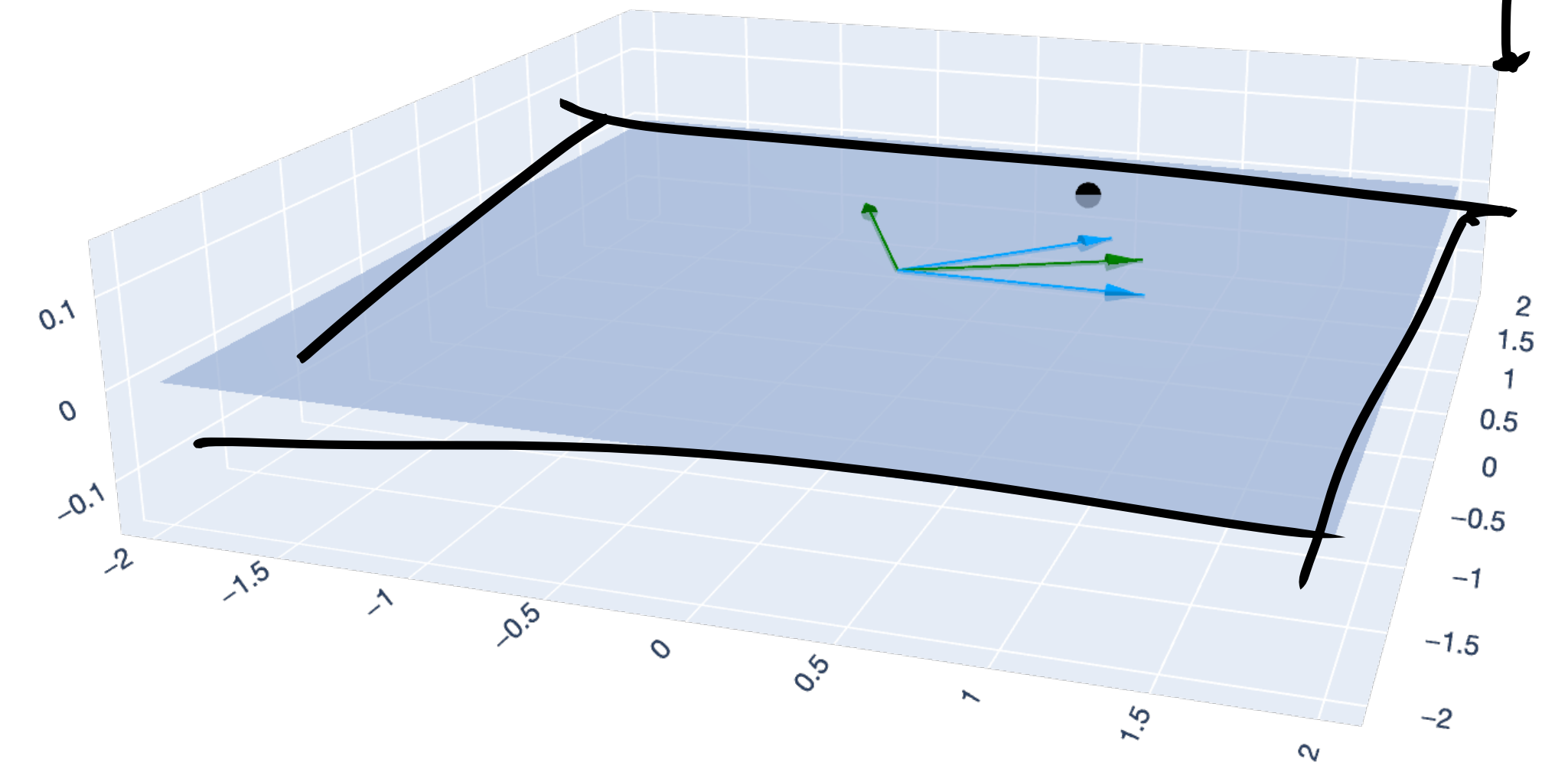
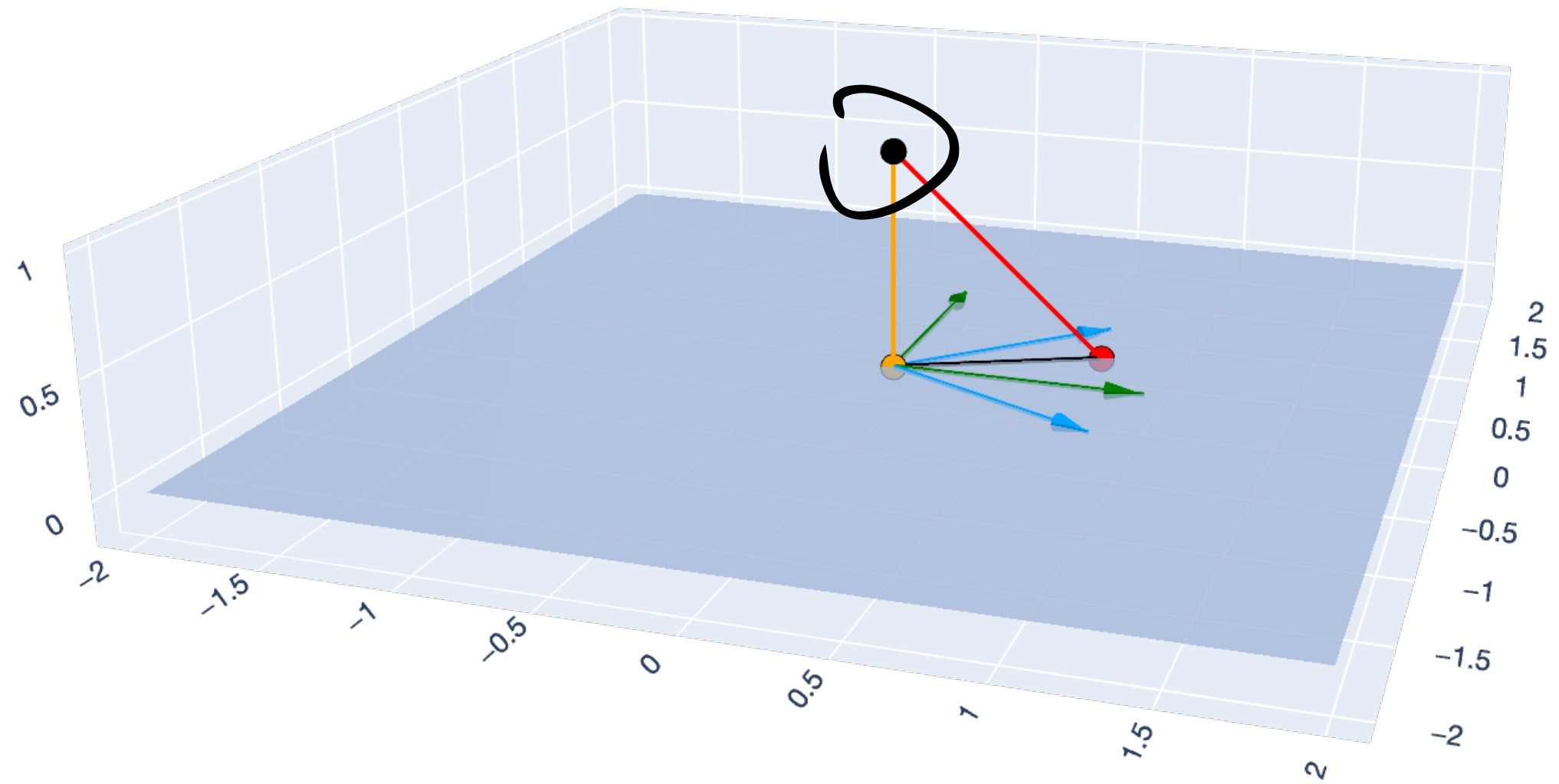
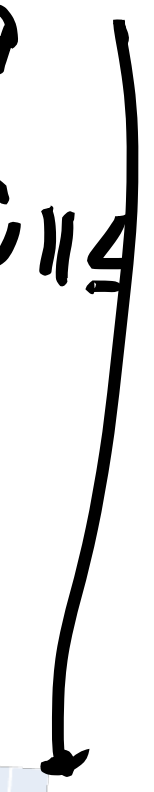
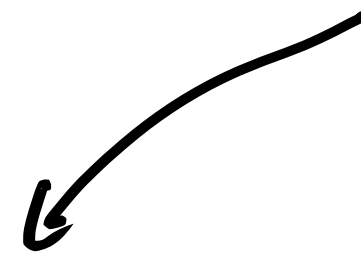
$X \in \mathbb{R}^{n \times d}$
 $n \geq d$
 $\hat{w} = (X^T X)^{-1} X^T y$

$X \hat{w} \approx y$

$d \geq n$

$X w = y$

$\| \hat{w} \|_2$
 $\| w \|_2$
 $\forall w \in \mathbb{R}^d$



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

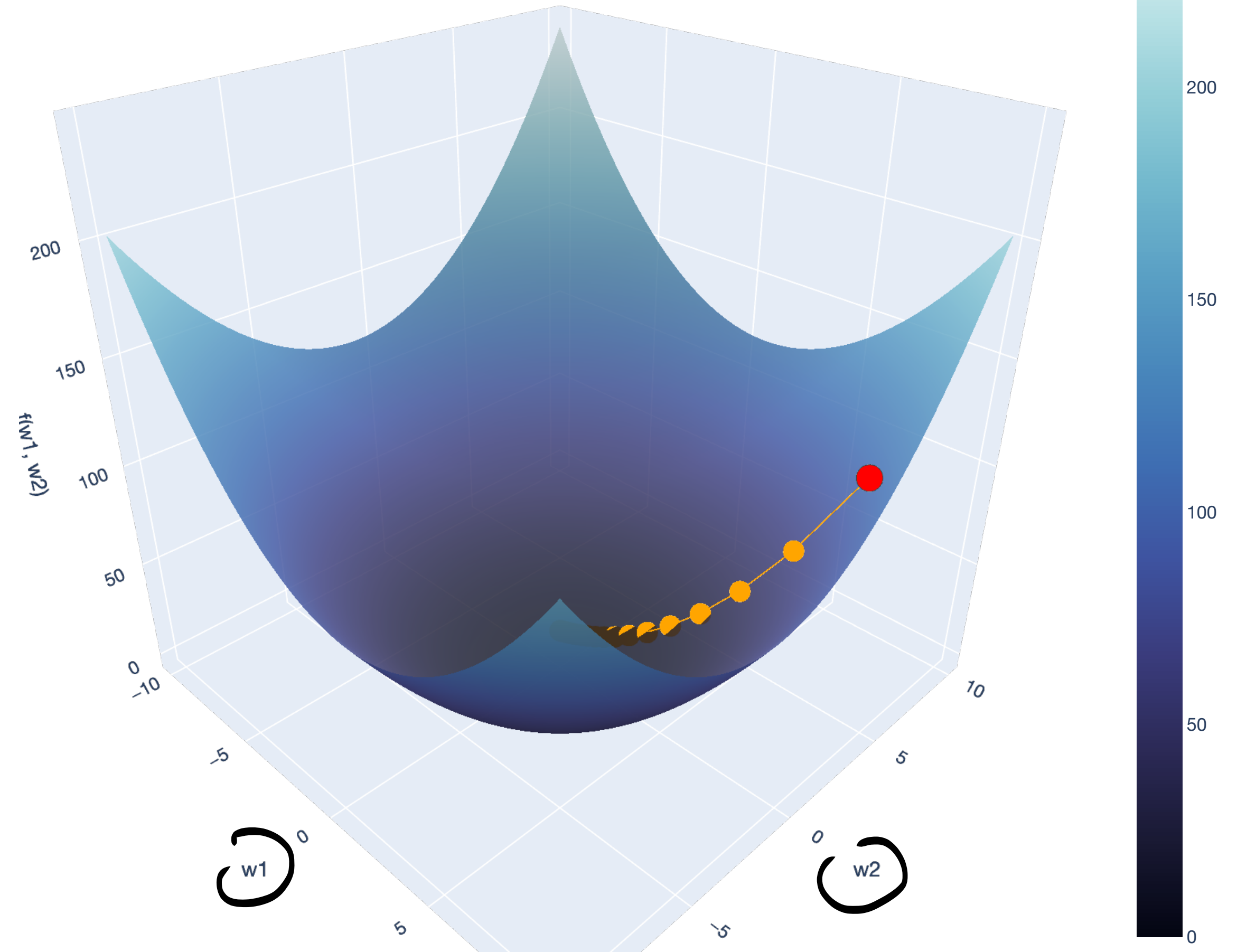
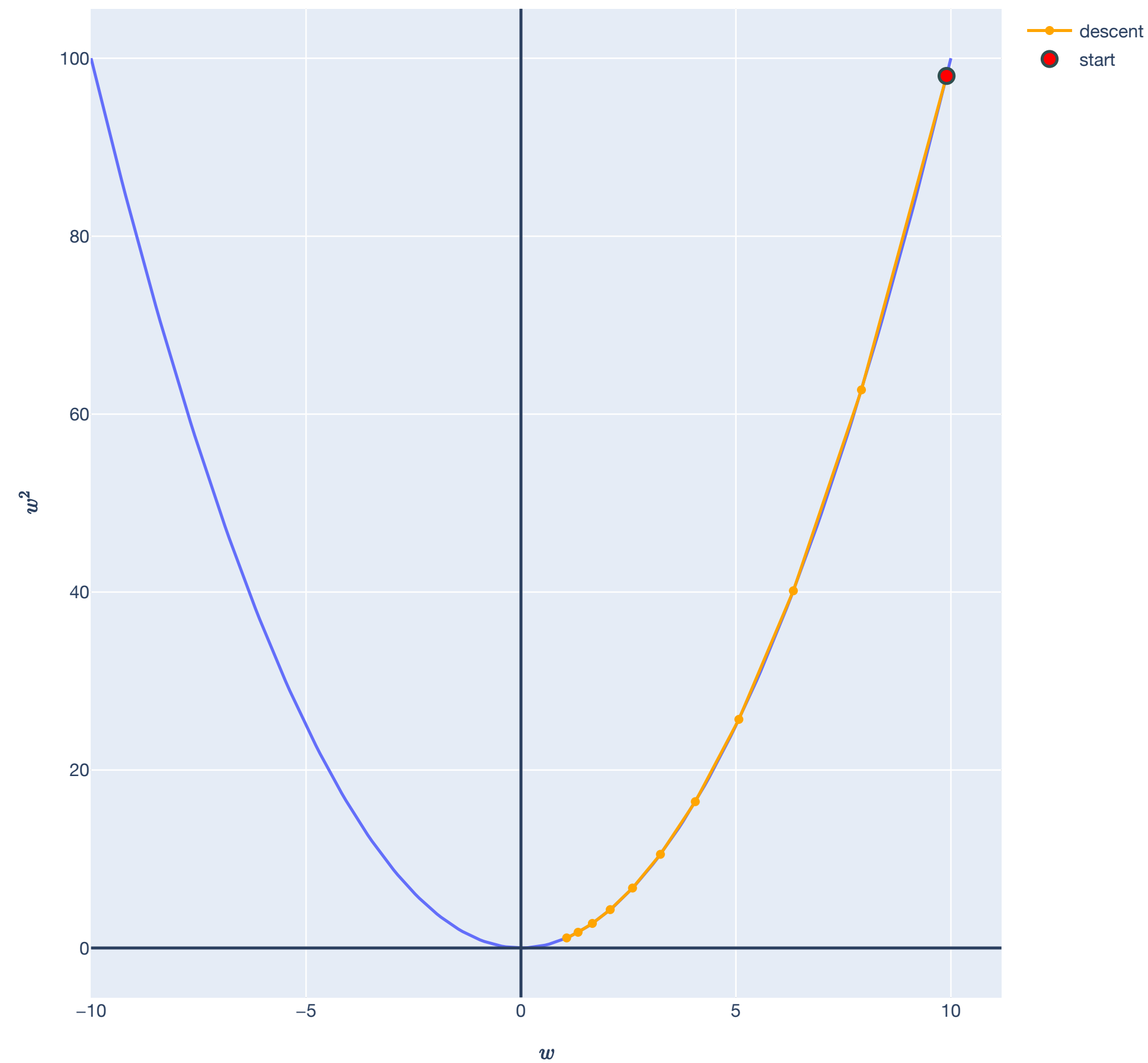
— x1 — x2 — u1 — u2 ● y

Lesson Overview

Big Picture: Gradient Descent

$$w \in \mathbb{R}^2$$
$$f(w) = \|Xw - \gamma\|^2$$

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

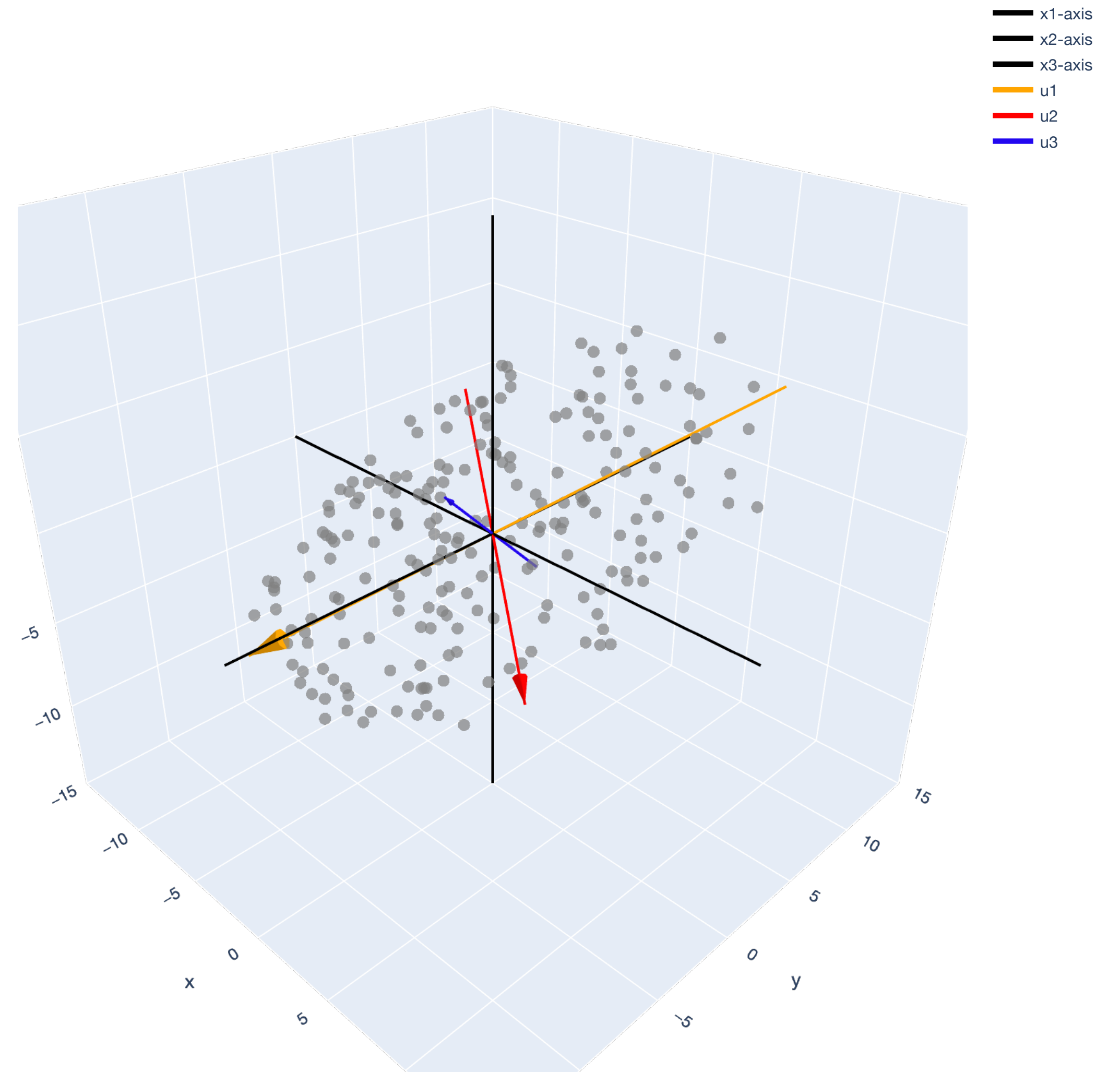
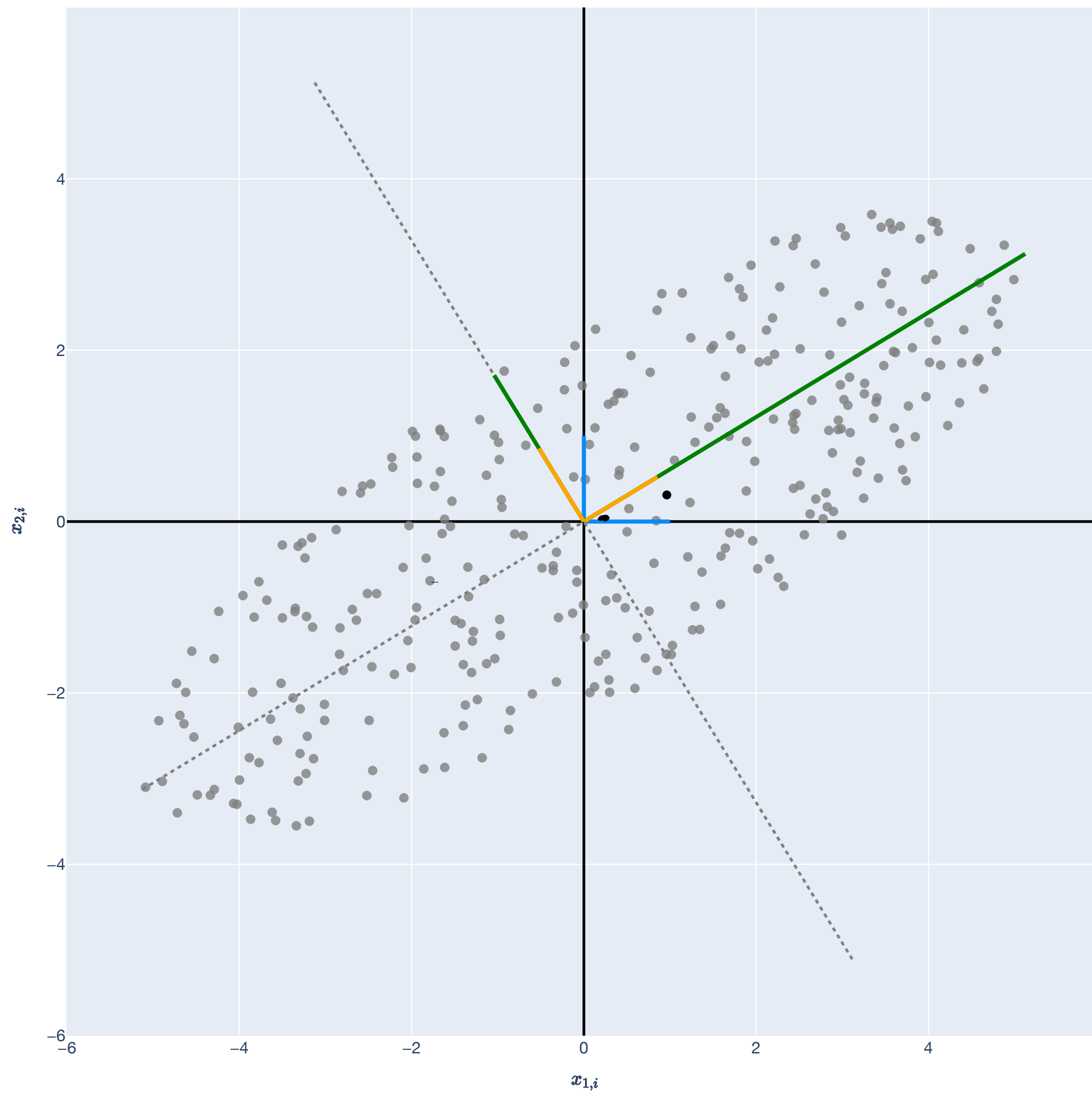


descent start

[Click to interact](#)

Lesson Overview

Big Picture: Singular Value Decomposition (SVD)



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

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

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*.

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

Regression Setup

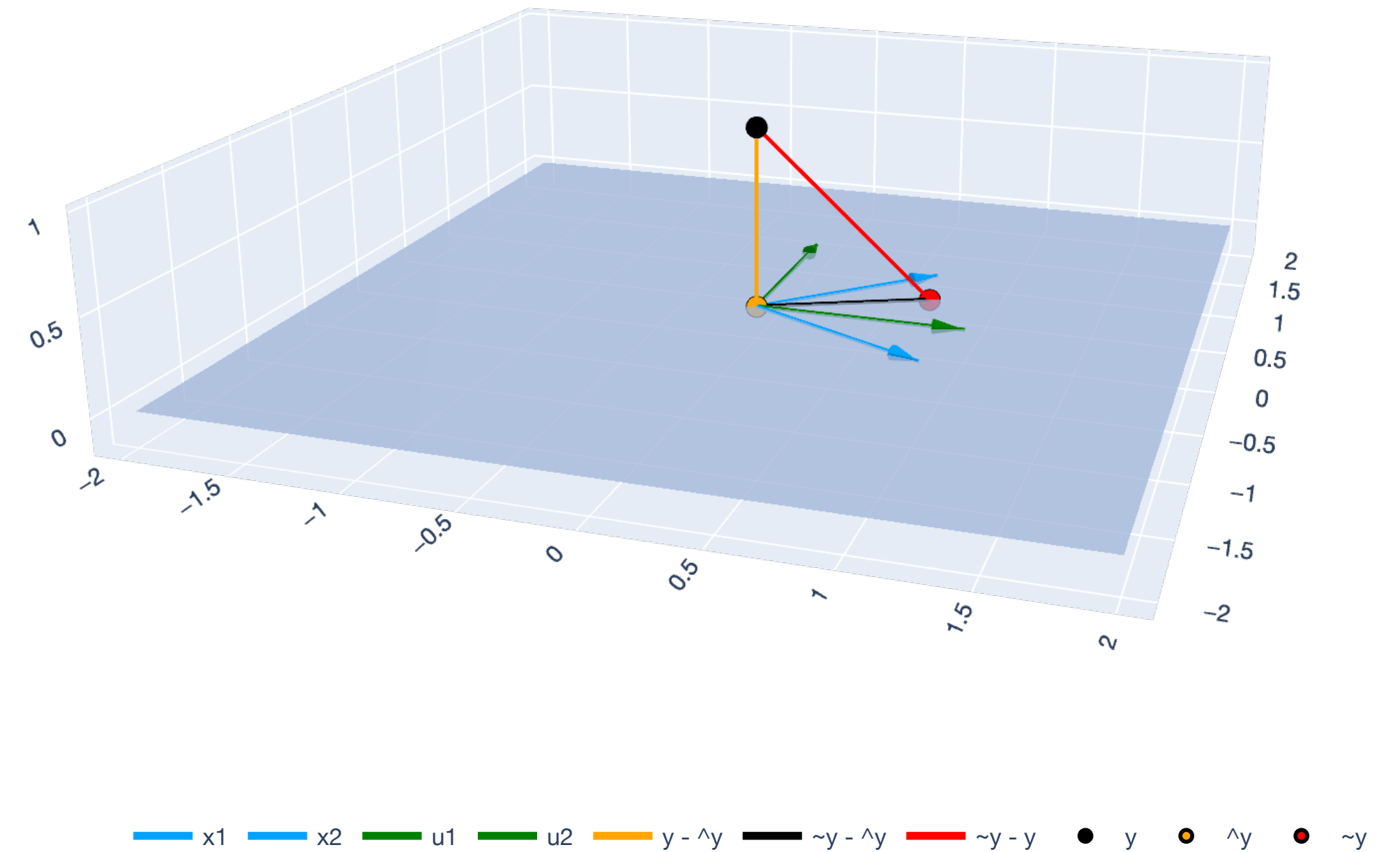
To find $\hat{\mathbf{w}}$, we follow the *principle of least squares*.

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

This gives the predictions $\hat{\mathbf{y}} \in \mathbb{R}^n$ that are close in a least squares sense:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}} \text{ such that } \|\hat{\mathbf{y}} - \mathbf{y}\|^2 \leq \|\tilde{\mathbf{y}} - \mathbf{y}\|^2$$

(for $\tilde{\mathbf{y}} = \mathbf{X}\mathbf{w}$ from any other $\mathbf{w} \in \mathbb{R}^d$).



Least Squares

OLS Theorem

Theorem (Ordinary Least Squares). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^n$. Let $\hat{\mathbf{w}} \in \mathbb{R}^d$ be the least squares minimizer:

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

If $\underline{n} \geq \underline{d}$ and $\underline{\text{rank}(\mathbf{X})} = \underline{d}$, then:

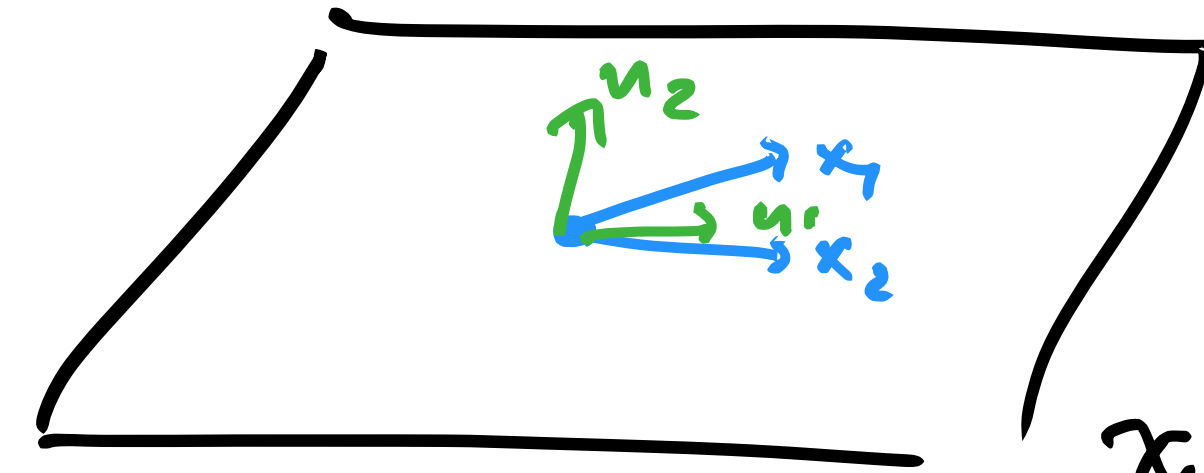
$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

To get predictions $\hat{\mathbf{y}} \in \mathbb{R}^n$:

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

Least Squares

OLS with Orthogonal Basis



$$\mathcal{X} = \text{span}(\text{col}(X))$$

Theorem (OLS with orthogonal basis). Let $\mathcal{X} \subseteq \mathbb{R}^n$ be a subspace and let $\mathbf{u}_1, \dots, \mathbf{u}_d \in \mathbb{R}^n$ be an orthonormal basis for \mathcal{X} , with semi-orthogonal matrix $\mathbf{U} \in \mathbb{R}^{n \times d}$. Let $\mathbf{y} \in \mathbb{R}^n$ and let $\hat{\mathbf{w}} \in \mathbb{R}^d$ be the least squares minimizer:

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

which is solved by:

$$\hat{\mathbf{w}} = \mathbf{U}^T \mathbf{y}.$$

Additionally, the projection $\hat{\mathbf{y}} \in \mathbb{R}^n$ is given by $\Pi_{\mathcal{X}}(\mathbf{y}) = \arg \min_{\hat{\mathbf{y}} \in \mathcal{X}} \|\hat{\mathbf{y}} - \mathbf{y}\|^2$:

$$\hat{\mathbf{y}} = \Pi_{\mathcal{X}}(\mathbf{y}) = \mathbf{U}\mathbf{U}^T \mathbf{y}.$$

Least Squares

OLS with Orthogonal Basis

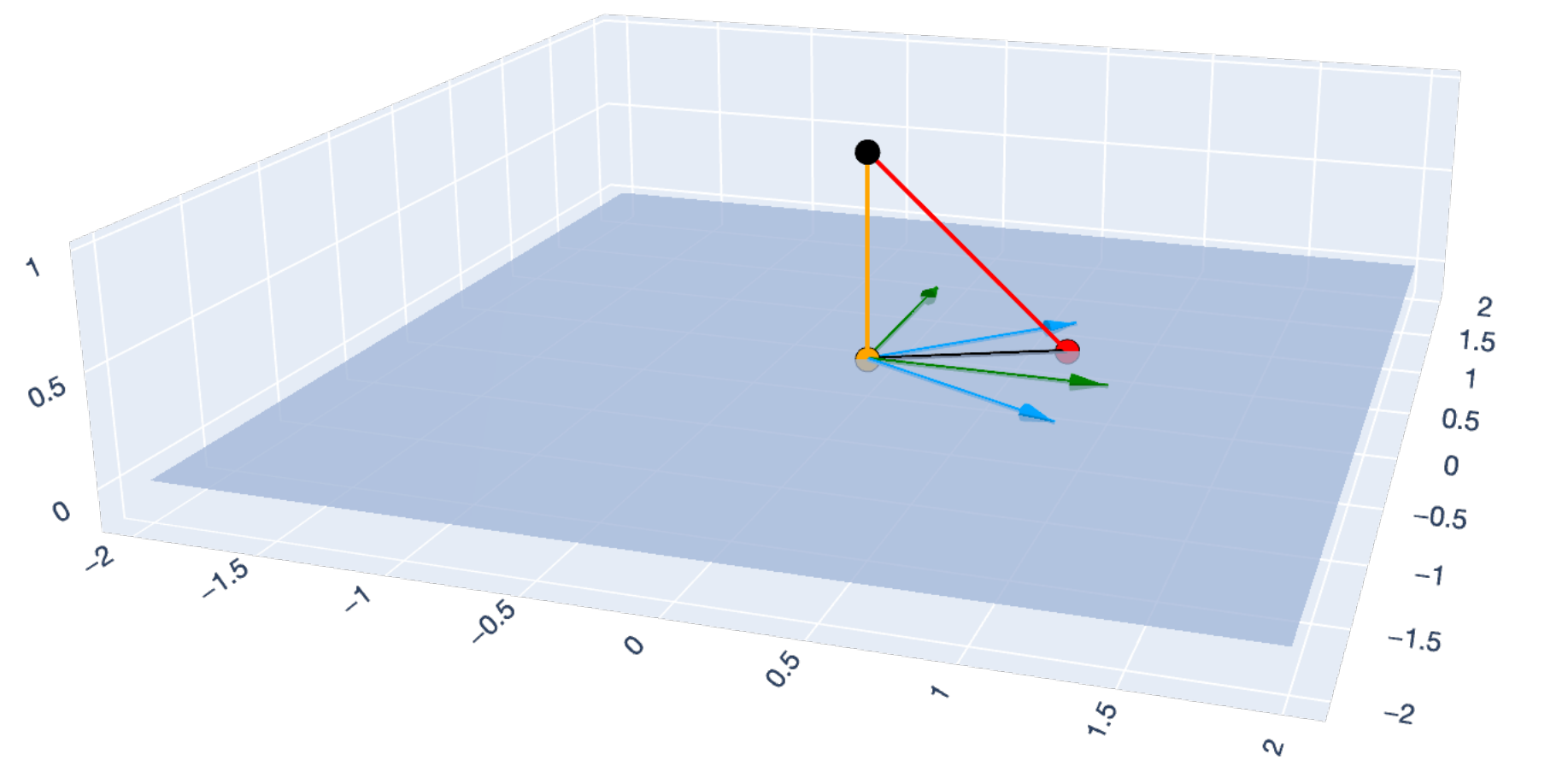
$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\hat{\mathbf{y}} = \Pi_{\mathcal{X}}(\mathbf{y}) = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

if we did have
 u_1, \dots, u_d

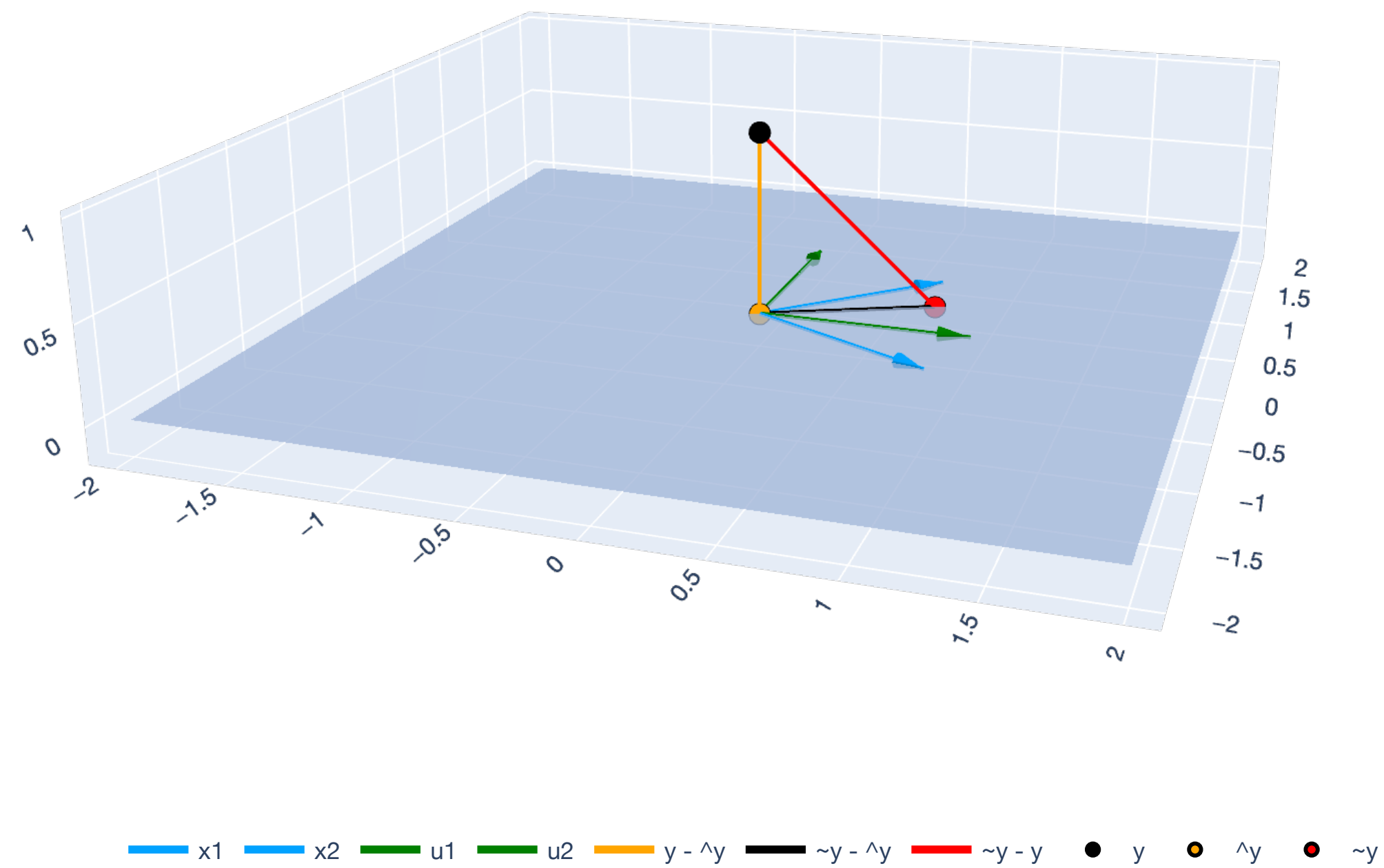
$$\hat{\mathbf{w}}_{onb} = \mathbf{U}^\top \mathbf{y}$$

$$\hat{\mathbf{y}} = \Pi_{\mathcal{X}}(\mathbf{y}) = \mathbf{U}\mathbf{U}^\top \mathbf{y}$$



— x1 — x2 — u1 — u2 — y - \hat{y} — \hat{y} - y — y - \hat{y} • y • \hat{y} • \hat{y}

How to find a good orthogonal basis?

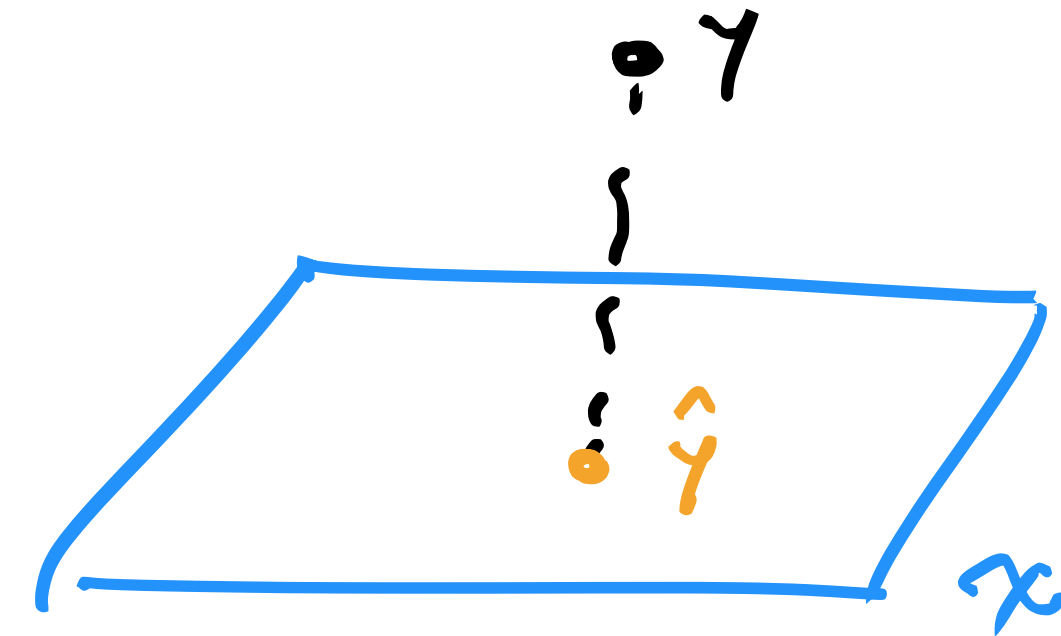


Properties of Projections

Projection Matrices and Orthogonal Complement

Projection

Projection of a vector onto a subspace



For a subspace $\mathcal{X} \subseteq \mathbb{R}^n$, the **projection** of a vector $\mathbf{y} \in \mathbb{R}^n$ onto the set \mathcal{X} is the closest vector $\hat{\mathbf{y}}$ in \mathcal{X} to \mathbf{y} , in a Euclidean distance sense:

$$\hat{\mathbf{y}} = \arg \min_{\hat{\mathbf{y}} \in \mathcal{X}} \|\hat{\mathbf{y}} - \mathbf{y}\| = \|\hat{\mathbf{y}} - \mathbf{y}\|^2.$$

Let $\mathcal{X} = \text{span}(\text{col}(\mathbf{X}))$. Any point $\hat{\mathbf{y}} \in \mathcal{X}$ is a linear combination $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{w}}$, with:

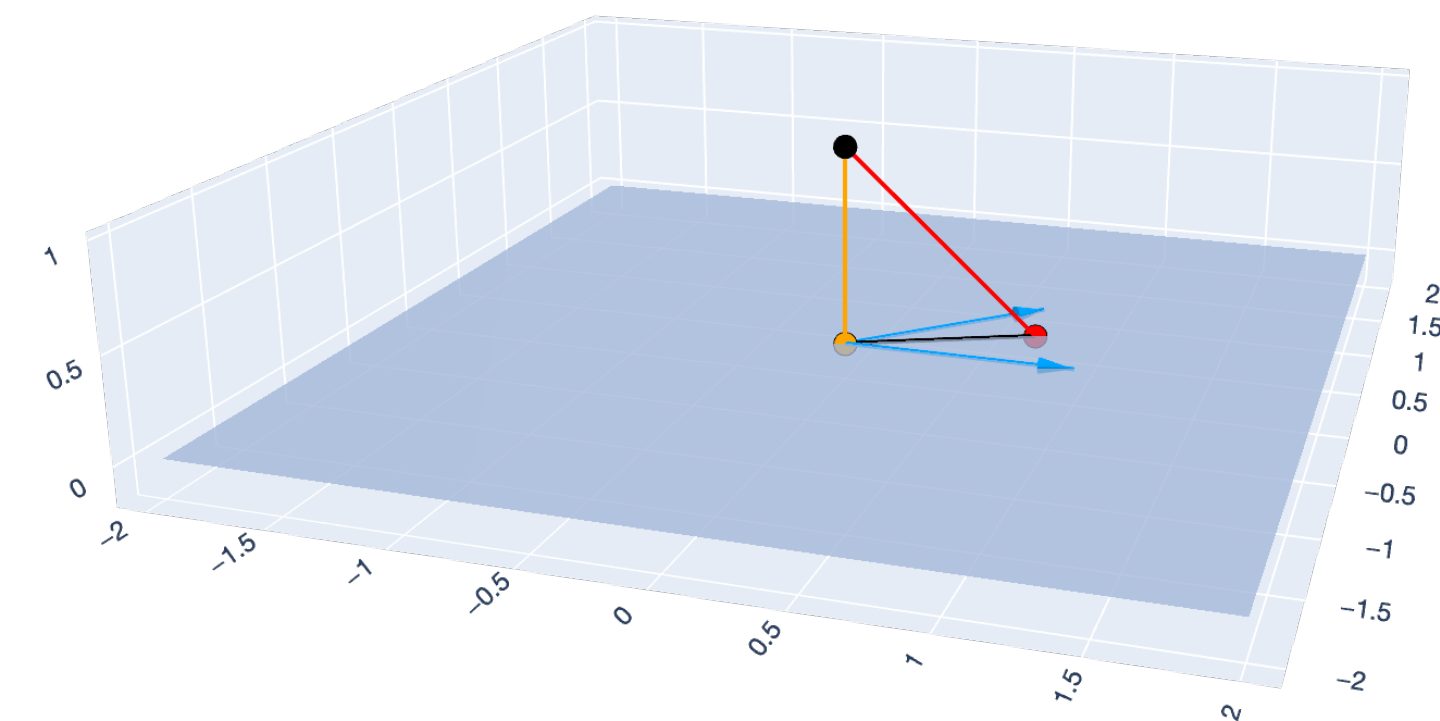
$$\hat{\mathbf{w}} = \arg \min_{\hat{\mathbf{w}} \in \mathbb{R}^d} \|\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}\|^2$$

Least Squares as Projection

Projection Matrix

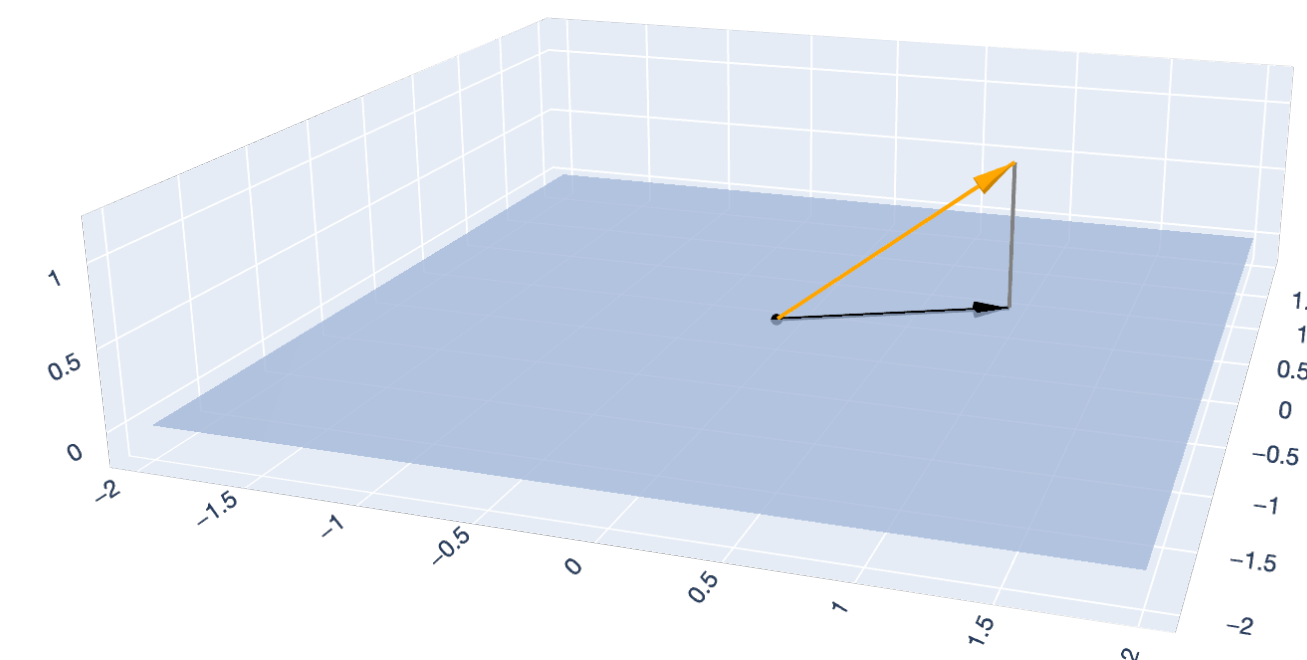
$$\hat{w} = \arg \min_{\hat{w} \in \mathbb{R}^d} \|X\hat{w} - y\|^2$$

This is just least squares! By what we've learned...



— x1 — x2 — y - \hat{y} — $\sim y - \hat{y}$ — $\sim y - y$ • y • \hat{y} • $\sim y$

Click to



— y - proj_y — y — proj_y • origin

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

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

$n \times d$ $d \times d$ $n \times d$ $d \times n$

The **projection matrix** is:

$$P_X = X(X^T X)^{-1} X^T \in \mathbb{R}^{n \times n}$$

$$P_X y = \Pi_X(y)$$

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

α

$n \times n$

Least Squares as Projection

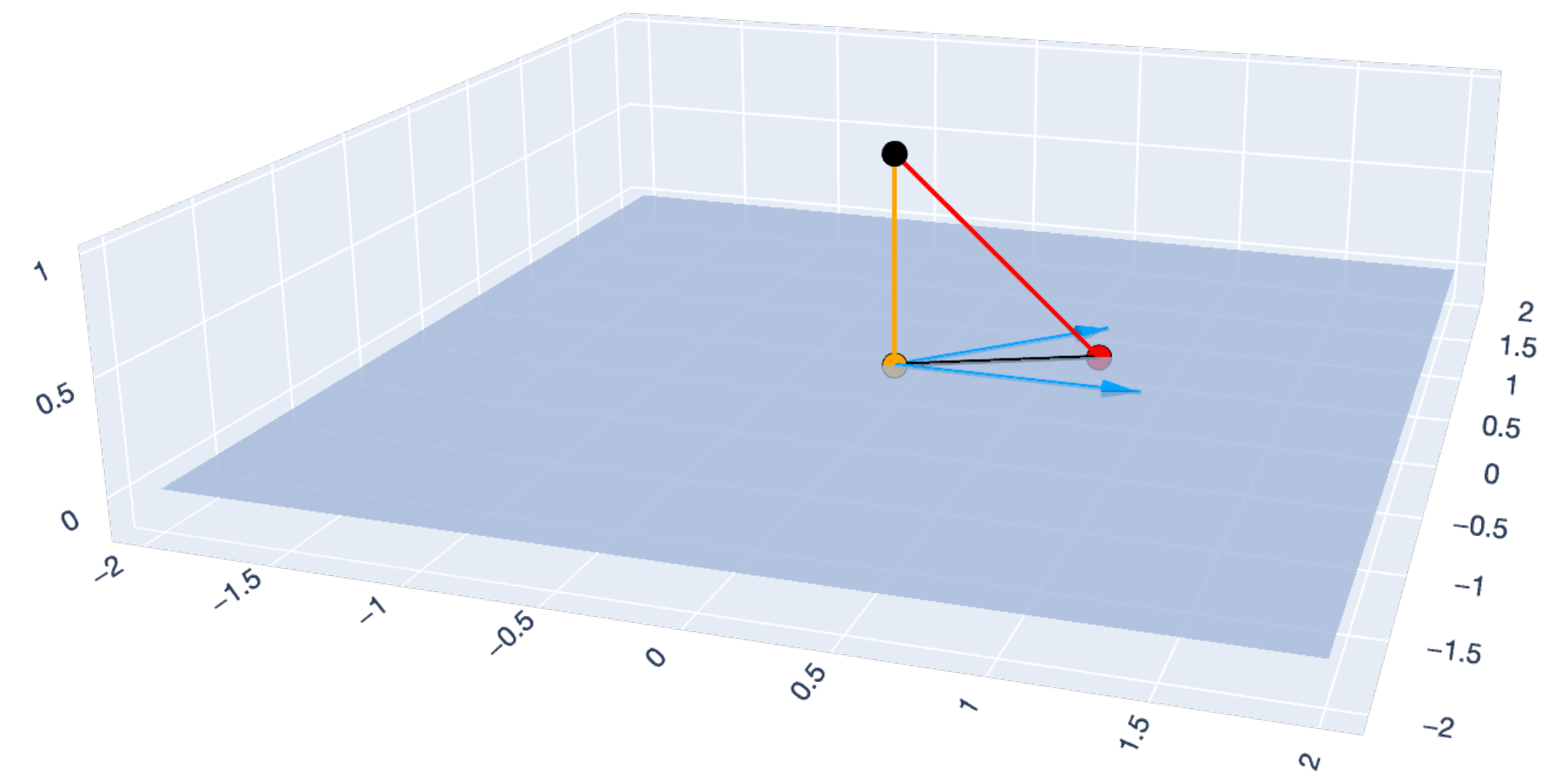
Projection Matrix

Any matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ has a subspace $\mathcal{X} = \text{span}(\text{col}(\mathbf{X}))$.

If the columns $\mathbf{x}_1, \dots, \mathbf{x}_d$ are *linearly independent*, then:

$$\Pi_{\mathcal{X}}^{(\mathcal{Y})} = \underline{\underline{P_{\mathcal{X}}}} \mathbf{y} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y},$$

where $P_{\mathcal{X}} \in \mathbb{R}^{n \times n}$ is a projection matrix.



— x1 — x2 — y - ^y — -y - ^y — -y - y • y • ^y • -y

Click to

Orthogonal Complement

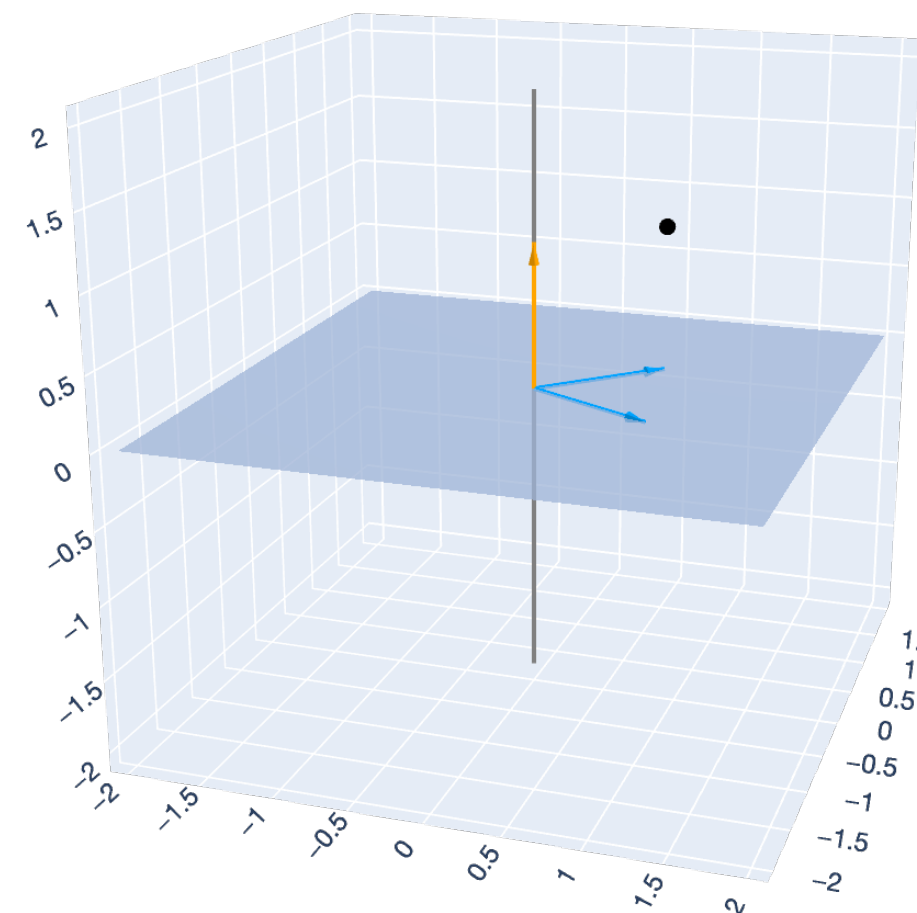
Intuition

Any subspace $A \subseteq \mathbb{R}^n$ has an orthogonal complement A^\perp . All vectors in A are orthogonal to all the vectors in A^\perp , and vice versa.

"A-perp"



Any vector $\mathbf{y} \in \mathbb{R}^n$ can be constructed by adding a vector from A to a vector from A^\perp .



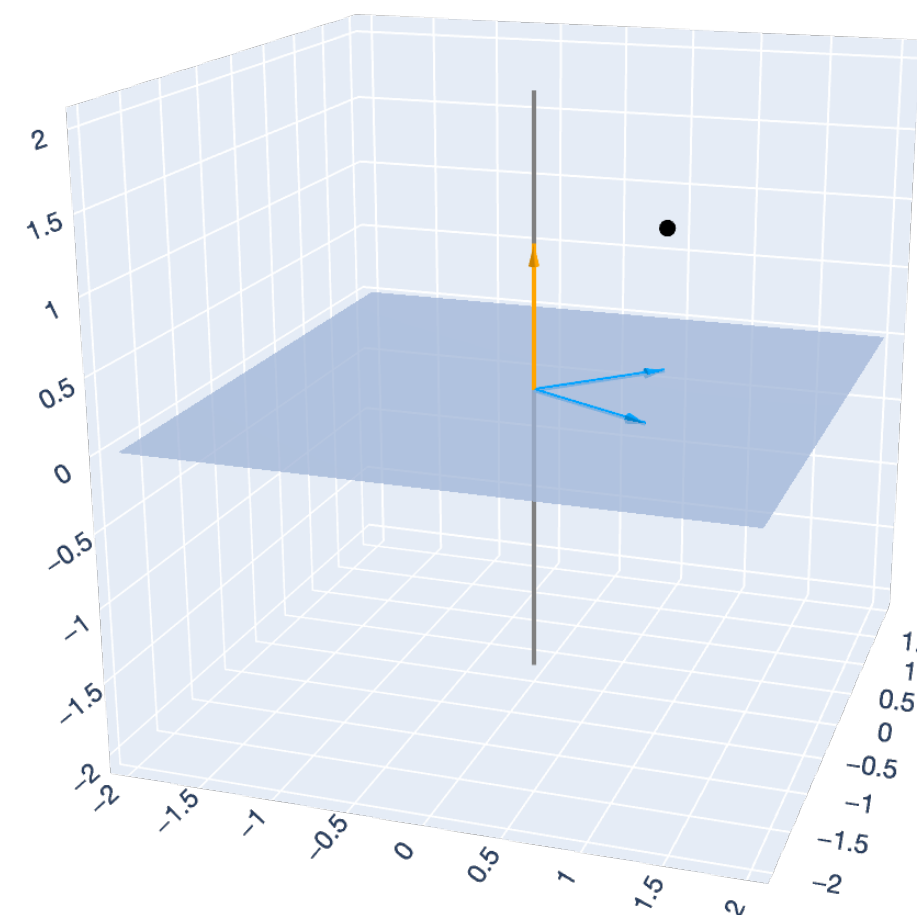
— u1 — u2 — v1 — — — • x

Orthogonal Complement

Definition

Let $A \subseteq \mathbb{R}^n$ be a subspace. The orthogonal complement of A , written A^\perp , is the set of vectors

$$A^\perp := \{ \underline{\mathbf{v}} \in \mathbb{R}^n : \underline{\langle \mathbf{v}, \mathbf{u} \rangle} = 0 \text{ for all } \underline{\mathbf{u}} \in A \}.$$

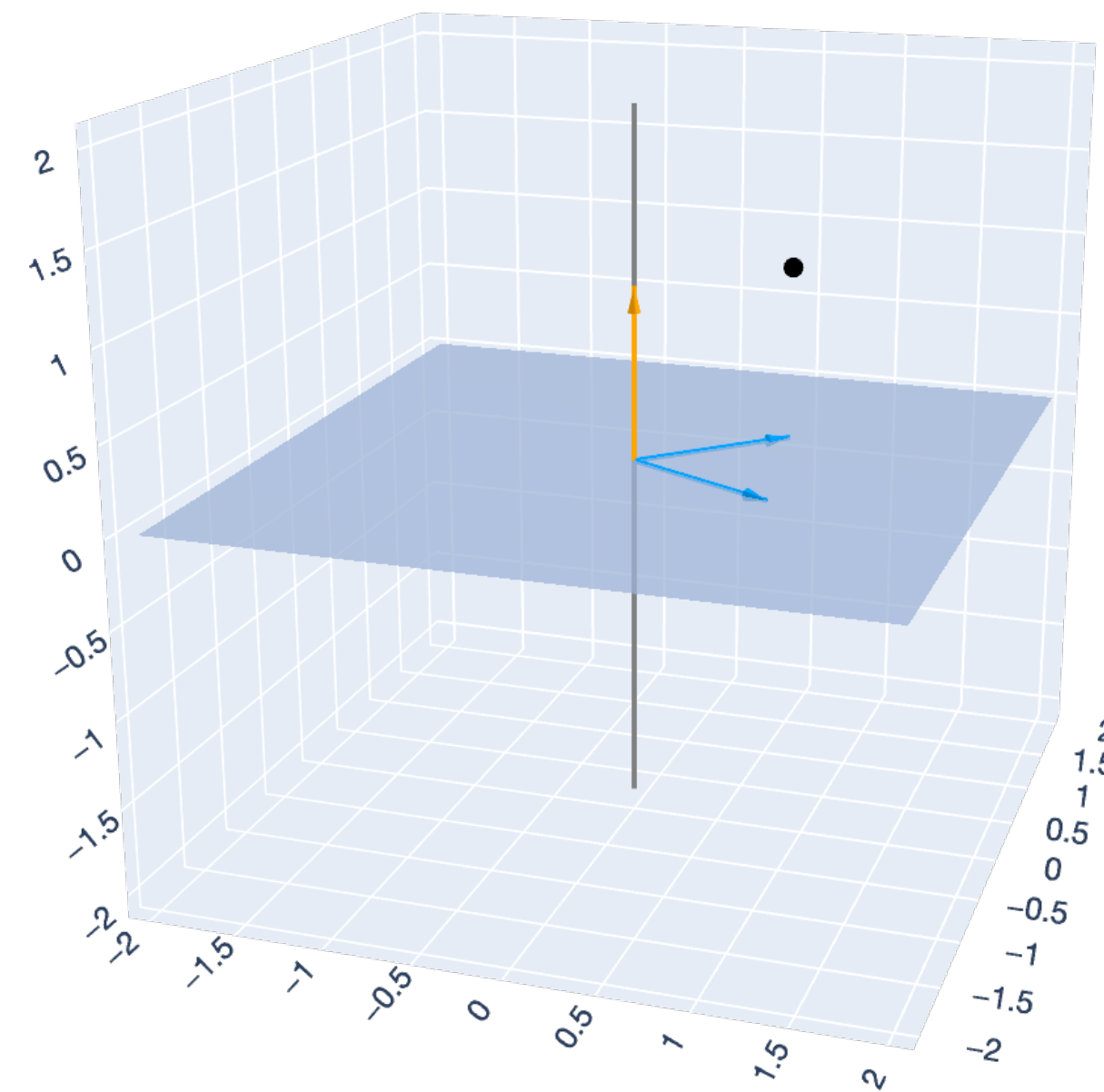


Orthogonal Complement

Dimension

$$A^\perp \subseteq \mathbb{R}^n$$

For any subspace $A \subseteq \mathbb{R}^n$ with $\dim(A) = d$, the orthogonal complement A^\perp has $\dim(A^\perp) = n - d$.



— u_1 — u_2 — v_1 — — — x

Orthogonal Complement

Orthogonal Complement and Matrices

Let $\mathbf{a}_1, \dots, \mathbf{a}_d \in \mathbb{R}^n$ be a basis for the subspace $\underline{A} \subseteq \mathbb{R}^n$. Let $\underline{\mathbf{b}_1, \dots, \mathbf{b}_{n-d}}$ be a basis for the orthogonal complement, \underline{A}^\perp .

Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ have columns $\mathbf{a}_1, \dots, \mathbf{a}_d$. Let $\mathbf{B} \in \mathbb{R}^{n \times (n-d)}$ have columns $\mathbf{b}_1, \dots, \mathbf{b}_{n-d}$. Then:

$$\underline{\mathbf{A}^\top \mathbf{B} = \mathbf{0} \text{ and } \mathbf{B}^\top \mathbf{A} = \mathbf{0}.}$$

$$\mathbf{A} = \begin{bmatrix} | & & | \\ \mathbf{a}_1 & \dots & \mathbf{a}_d \\ | & & | \end{bmatrix} \Rightarrow \mathbf{A}^\top = \begin{bmatrix} \text{---} \mathbf{a}_1 \text{---} \\ \vdots \\ \text{---} \mathbf{a}_d \text{---} \end{bmatrix} \begin{bmatrix} | & & | \\ \mathbf{b}_1 & \dots & \mathbf{b}_{n-d} \\ | & & | \end{bmatrix} = \begin{bmatrix} 0 \end{bmatrix}$$

Orthogonal Complement

Orthogonal Complement and Projections

Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ have columns $\mathbf{a}_1, \dots, \mathbf{a}_d$. Let $\mathbf{B} \in \mathbb{R}^{n \times (n-d)}$ have columns $\mathbf{b}_1, \dots, \mathbf{b}_{n-d}$, a basis for the orthogonal complement of $\text{span}(\text{col}(\mathbf{A}))$. Then:

$$\mathbf{A}^\top \mathbf{B} = \mathbf{0} \text{ and } \mathbf{B}^\top \mathbf{A} = \mathbf{0}.$$

We can break down any vector $\mathbf{x} \in \mathbb{R}^n$ into two projections:

$$\mathbf{x} = P_{\mathbf{A}} \mathbf{x} + P_{\mathbf{B}} \mathbf{x}.$$

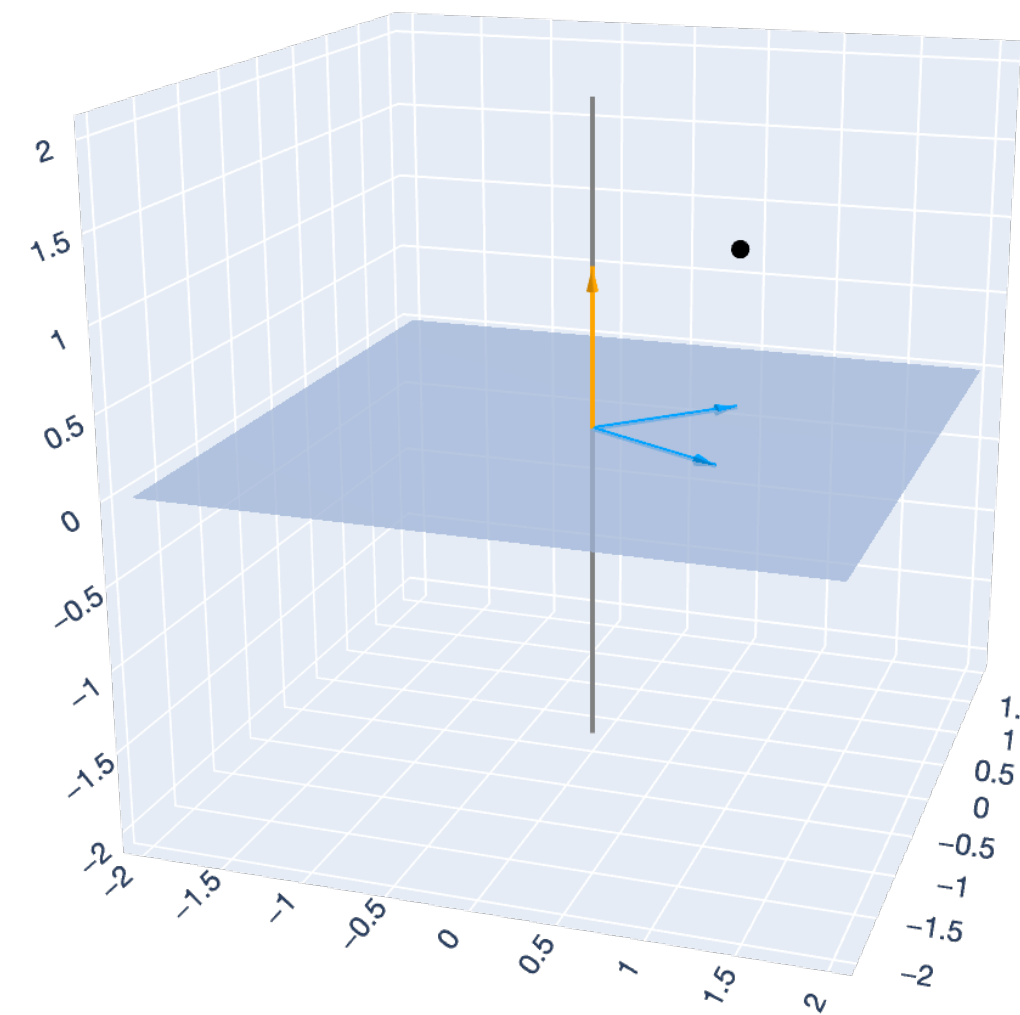
Orthogonal Complement

Orthogonal Complement and Projections

We can break down any vector $\mathbf{x} \in \mathbb{R}^n$ into two projections:

$$\mathbf{x} = P_A \mathbf{x} + P_B \mathbf{x}.$$

Additionally, $\mathbf{I} = P_A + P_B.$



— u1 — u2 — v1 — — — • x

Projection Matrices

Properties

Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ be a matrix and let $\mathbf{B} \in \mathbb{R}^{n \times (n-d)}$ have columns $\mathbf{b}_1, \dots, \mathbf{b}_{n-d}$, a basis for the orthogonal complement of $\text{span}(\text{col}(\mathbf{A}))$.

Prop (Orthogonal Decomposition). For any vector $\mathbf{x} \in \mathbb{R}^n$,

$$\mathbf{x} = P_{\mathbf{A}}\mathbf{x} + P_{\mathbf{B}}\mathbf{x}.$$

Projection Matrices

Properties

Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ be a matrix and let $\mathbf{B} \in \mathbb{R}^{n \times (n-d)}$ have columns $\mathbf{b}_1, \dots, \mathbf{b}_{n-d}$, a basis for the orthogonal complement of $\text{span}(\text{col}(\mathbf{A}))$.

Prop (Orthogonal Decomposition). For any vector $\mathbf{x} \in \mathbb{R}^n$,

$$\mathbf{x} = P_{\mathbf{A}}\mathbf{x} + P_{\mathbf{B}}\mathbf{x}.$$

Prop (Projection and Orthogonal Complement Matrices). $P_{\mathbf{A}} + P_{\mathbf{B}} = \mathbf{I}$.

Projection Matrices

Properties

Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ be a matrix and let $\mathbf{B} \in \mathbb{R}^{n \times (n-d)}$ have columns $\mathbf{b}_1, \dots, \mathbf{b}_{n-d}$, a basis for the orthogonal complement of $\text{span}(\text{col}(\mathbf{A}))$.

Prop (Orthogonal Decomposition). For any vector $\mathbf{x} \in \mathbb{R}^n$,

$$\mathbf{x} = P_{\mathbf{A}}\mathbf{x} + P_{\mathbf{B}}\mathbf{x}.$$

Prop (Projection and Orthogonal Complement Matrices). $P_{\mathbf{A}} + P_{\mathbf{B}} = \mathbf{I}$.

Prop (Projecting twice doesn't do anything). $P_{\mathbf{A}} = P_{\mathbf{A}}P_{\mathbf{A}} = P_{\mathbf{A}}^2$.

Projection Matrices

Properties

Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ be a matrix and let $\mathbf{B} \in \mathbb{R}^{n \times (n-d)}$ have columns $\mathbf{b}_1, \dots, \mathbf{b}_{n-d}$, a basis for the orthogonal complement of $\text{span}(\text{col}(\mathbf{A}))$.

Prop (Orthogonal Decomposition). For any vector $\mathbf{x} \in \mathbb{R}^n$,

$$\mathbf{x} = P_{\mathbf{A}}\mathbf{x} + P_{\mathbf{B}}\mathbf{x}.$$

Prop (Projection and Orthogonal Complement Matrices). $P_{\mathbf{A}} + P_{\mathbf{B}} = \mathbf{I}$.

Prop (Projecting twice doesn't do anything). $P_{\mathbf{A}} = P_{\mathbf{A}}P_{\mathbf{A}} = P_{\mathbf{A}}^2$.

Prop (Projections are symmetric). $P_{\mathbf{A}} = P_{\mathbf{A}}^{\top}$.

$$P_{\mathbf{X}} = \mathbf{X}(\mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{X}^{\top}$$
$$P_{\mathbf{X}}^{\top} = \mathbf{X}(\mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{X}^{\top}$$

Projection Matrices

Properties

Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ be a matrix and let $\mathbf{B} \in \mathbb{R}^{n \times (n-d)}$ have columns $\mathbf{b}_1, \dots, \mathbf{b}_{n-d}$, a basis for the orthogonal complement of $\text{span}(\text{col}(\mathbf{A}))$.

Prop (Orthogonal Decomposition). For any vector $\mathbf{x} \in \mathbb{R}^n$,

$$\mathbf{x} = P_{\mathbf{A}}\mathbf{x} + P_{\mathbf{B}}\mathbf{x}.$$

Prop (Projection and Orthogonal Complement Matrices). $P_{\mathbf{A}} + P_{\mathbf{B}} = \mathbf{I}$.

Prop (Projecting twice doesn't do anything). $P_{\mathbf{A}} = P_{\mathbf{A}}P_{\mathbf{A}} = P_{\mathbf{A}}^2$.

Prop (Projections are symmetric). $P_{\mathbf{A}} = P_{\mathbf{A}}^{\top}$.

Prop (1D projection formula). For the one-dimensional subspace associated to the vector $\mathbf{a} \in \mathbb{R}^n$, the projection matrix is:

$$\boxed{P_{\mathbf{a}} = \frac{\mathbf{a}\mathbf{a}^{\top}}{\mathbf{a}^{\top}\mathbf{a}}}$$

$\Rightarrow \mathbf{X}(\mathbf{X}^{\top}\mathbf{X})^{-1}\mathbf{X}^{\top} = \frac{\mathbf{a}\mathbf{a}^{\top}}{\mathbf{a}^{\top}\mathbf{a}}$

Handwritten notes: $\mathbf{x} \in \mathbb{R}^n$ and $P_{\mathbf{a}}\vec{\mathbf{x}} = \frac{\mathbf{a}\mathbf{a}^{\top}\vec{\mathbf{x}}}{\mathbf{a}^{\top}\mathbf{a}}$

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

Singular Value Decomposition

1D Intuition and Derivation

Singular Value Decomposition (SVD)

1D Picture

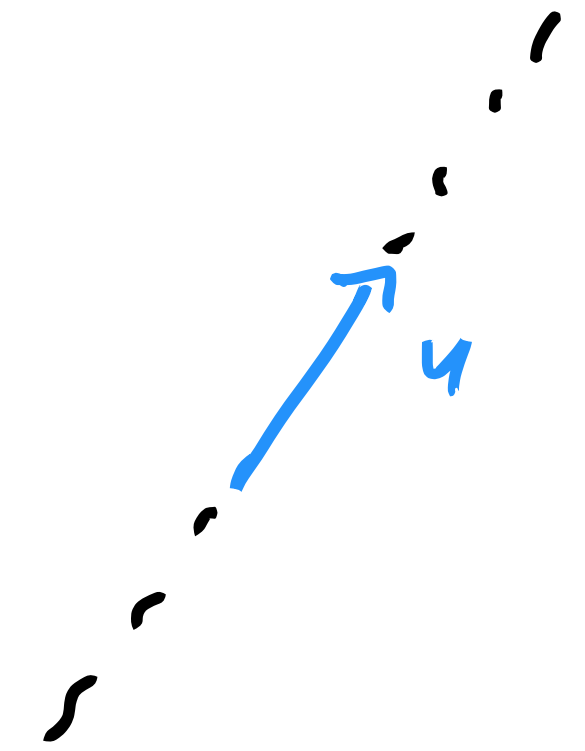
Observe data $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$, with matrix of *training samples* $\mathbf{X} \in \mathbb{R}^{n \times d}$:

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

Goal: Find the best one-dimensional subspace $\mathcal{U} \subseteq \mathbb{R}^n$ that fits the points.

Singular Value Decomposition (SVD)

1D Picture



Observe data $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$, with matrix of *training samples* $\mathbf{X} \in \mathbb{R}^{n \times d}$:

$$\mathbf{X} = \begin{bmatrix} \uparrow & & \uparrow \\ \mathbf{x}_1 & \dots & \mathbf{x}_d \\ \downarrow & & \downarrow \end{bmatrix} = \begin{bmatrix} \leftarrow & \mathbf{x}_1^T & \rightarrow \\ \vdots & & \vdots \\ \leftarrow & \mathbf{x}_n^T & \rightarrow \end{bmatrix} \cdot \begin{matrix} \mathbf{w} \in \mathbb{R}^d \\ \mathbf{x}_0 \in \mathbb{R}^d \\ \boxed{\mathbf{w}^T \mathbf{x}_0 = \hat{y}_0} \end{matrix}$$

Goal: Find the best one-dimensional subspace $\mathcal{U} \subseteq \mathbb{R}^n$ that fits the points.

A one-dimensional subspace is determined by a single vector $\mathbf{u} \in \mathbb{R}^n$:

$$\mathcal{U} = \{c\mathbf{u} : c \in \mathbb{R}\}.$$

Singular Value Decomposition (SVD)

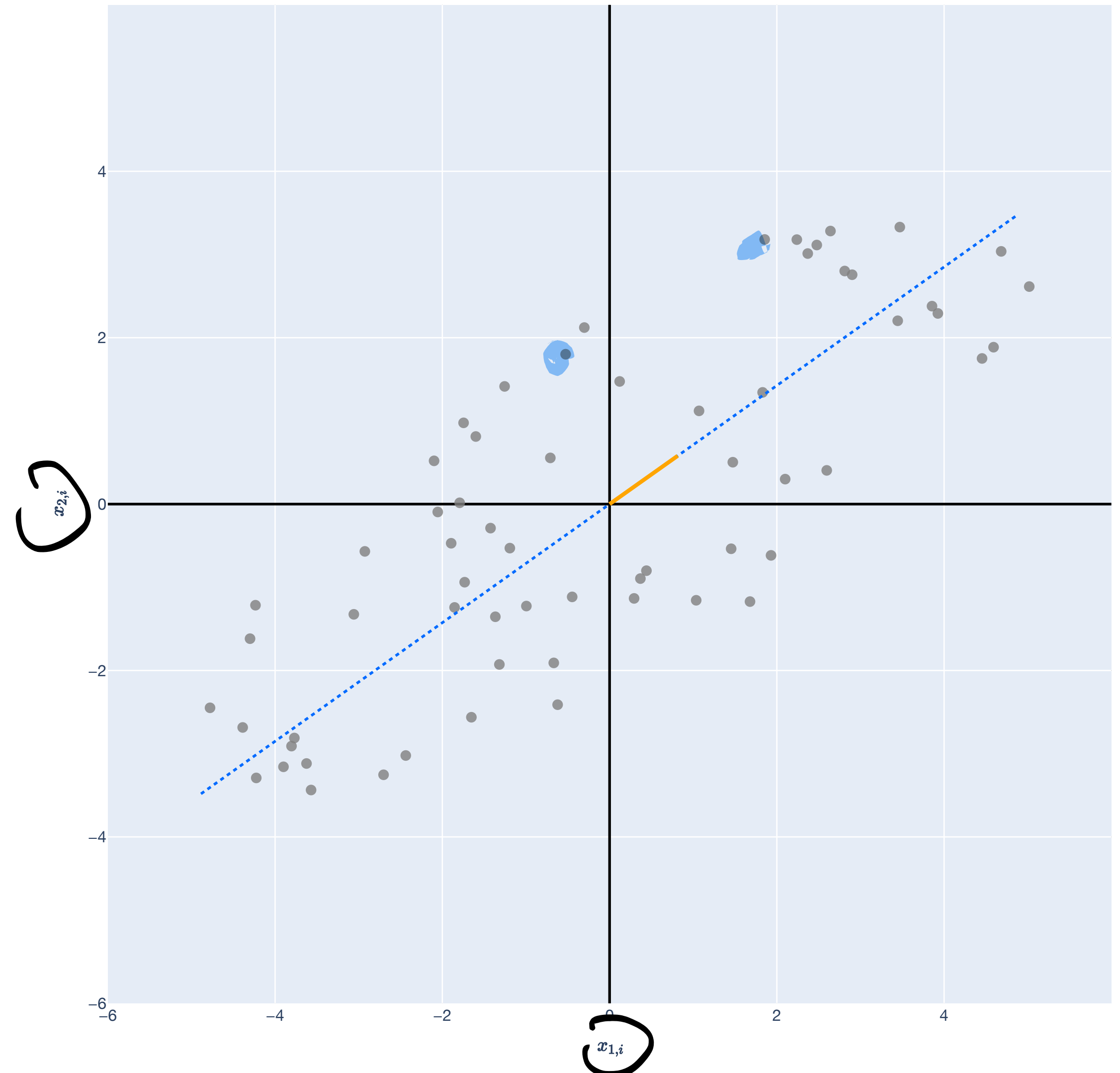
1D Picture

Observe data $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$.

Goal: Find the best one-dimensional subspace $\mathcal{U} \subseteq \mathbb{R}^n$ that fits the points.

$$X \in \mathbb{R}^{2 \times 50}$$

$$\begin{bmatrix} | & | & \dots & | & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_{49} & \mathbf{x}_{50} \\ | & | & \dots & | & | \end{bmatrix}$$



Singular Value Decomposition (SVD)

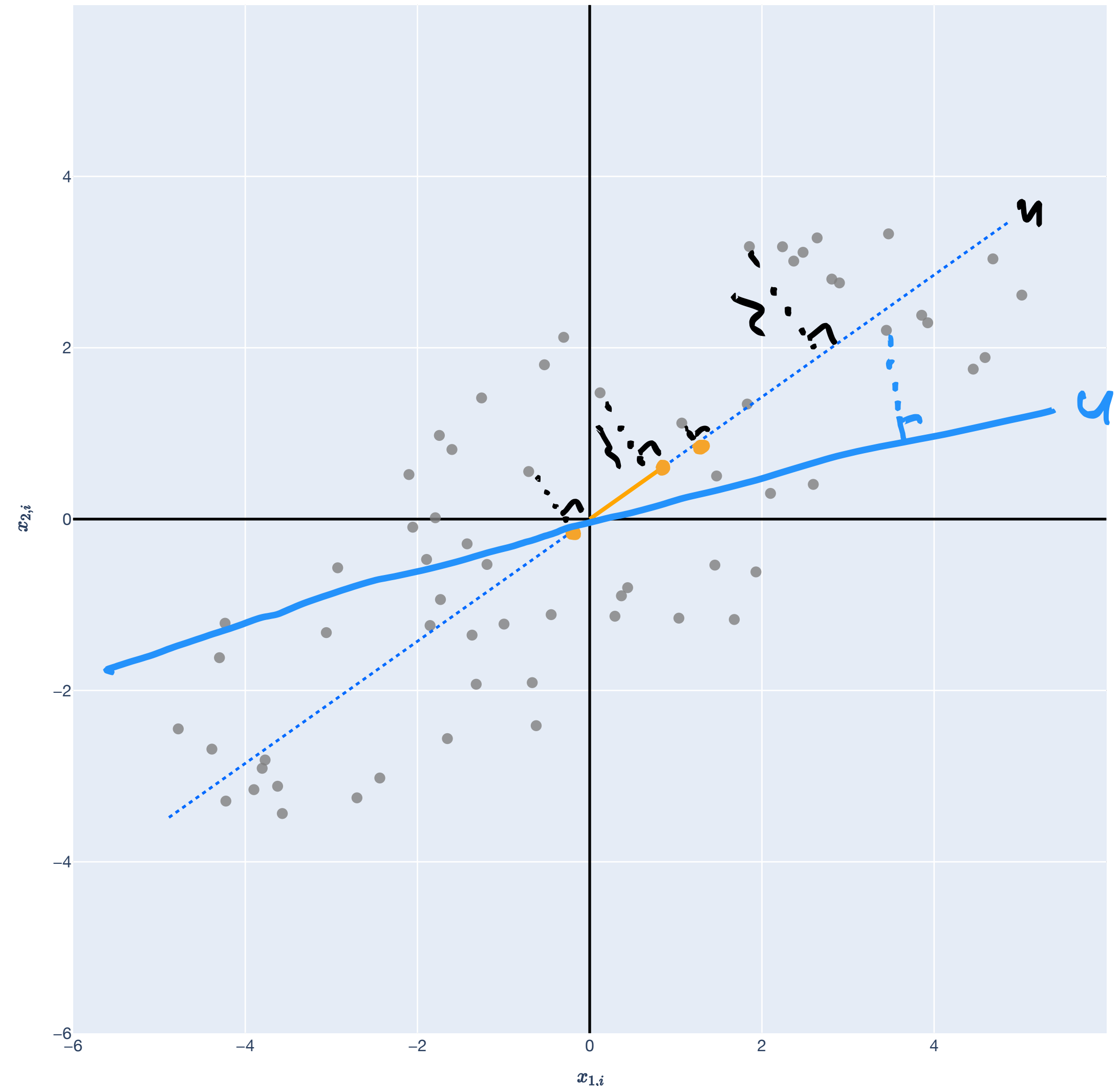
1D Picture

Observe data $\underbrace{\mathbf{x}_1, \dots, \mathbf{x}_d}_{\mathcal{X}} \in \mathbb{R}^n$.

Goal: Find the best one-dimensional subspace $\mathcal{U} \subseteq \mathbb{R}^n$ that fits the points.

How? Find $\mathbf{u} \in \mathbb{R}^n$ that minimizes the sum of squared projection distances:

$$\arg \min_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \underbrace{\|\mathbf{x}_i - \Pi_{\mathbf{u}}(\mathbf{x}_i)\|^2}_{\text{distance squared}}$$



Comparison with OLS

1D Pictures

OLS: Observe data $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$.

Goal: Find best linear combination $\hat{\mathbf{w}} \in \mathbb{R}^d$ of $\mathbf{x}_1, \dots, \mathbf{x}_d$ such that

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

BFS: Observe data $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$.

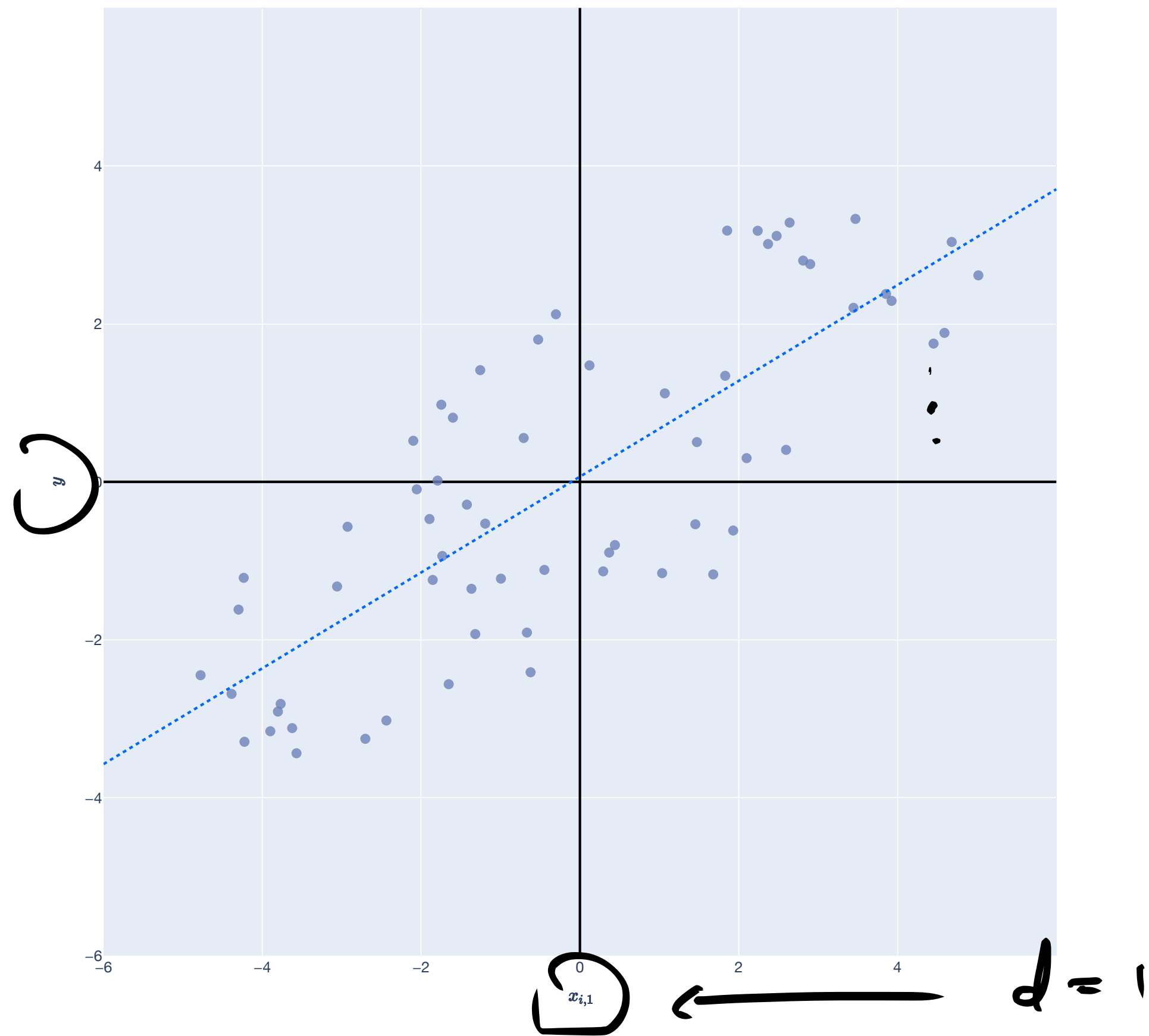
Goal: Find one-dimensional subspace determined by $\mathbf{u} \in \mathbb{R}^n$ such that

$$\arg \min_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \|\mathbf{x}_i - \Pi_{\mathbf{u}}(\mathbf{x}_i)\|^2$$

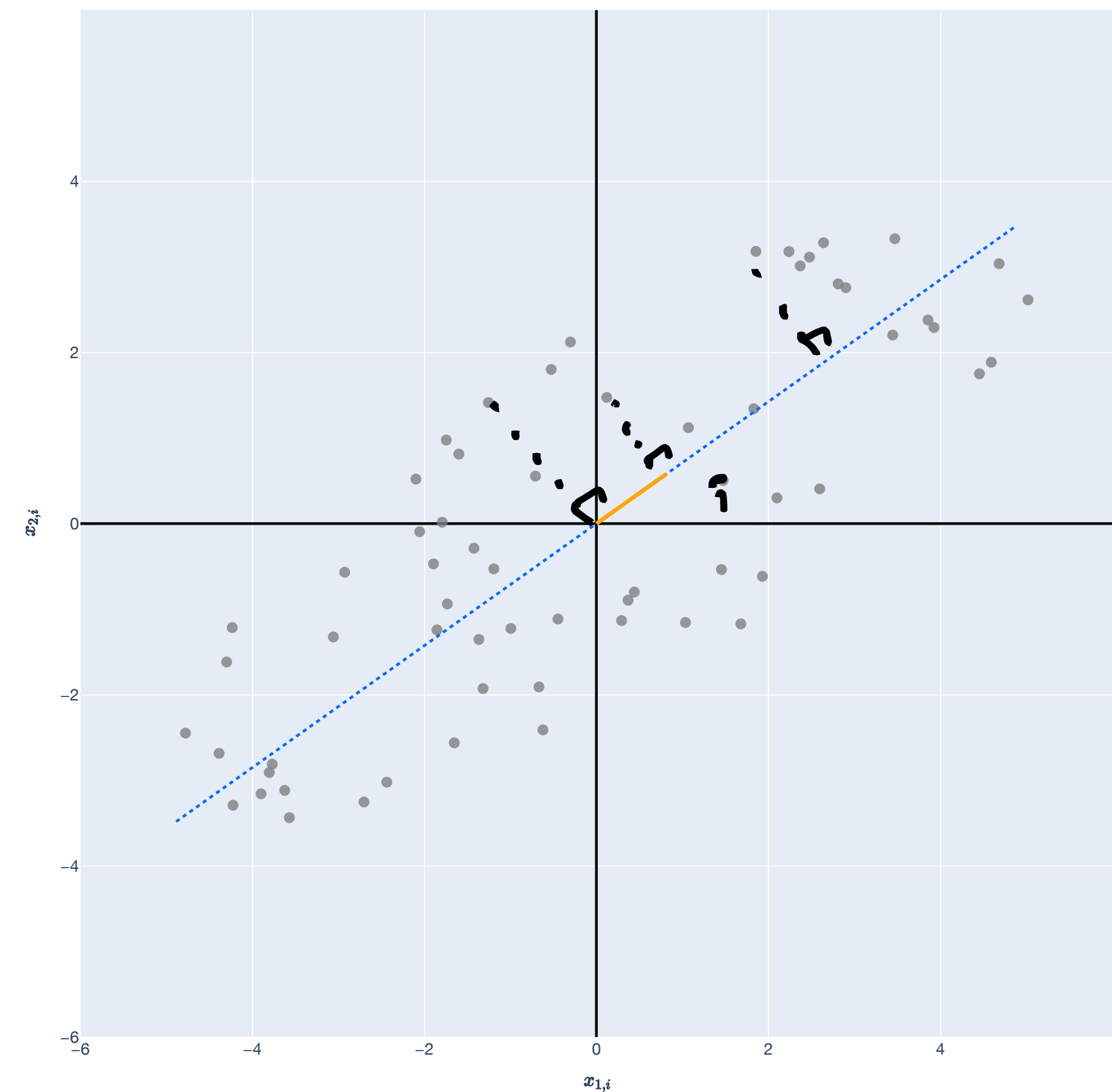
Comparison with OLS

1D Pictures

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



$$\arg \min_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \|\mathbf{x}_i - \underbrace{\Pi_{\mathbf{u}}(\mathbf{x}_i)}\|^2$$



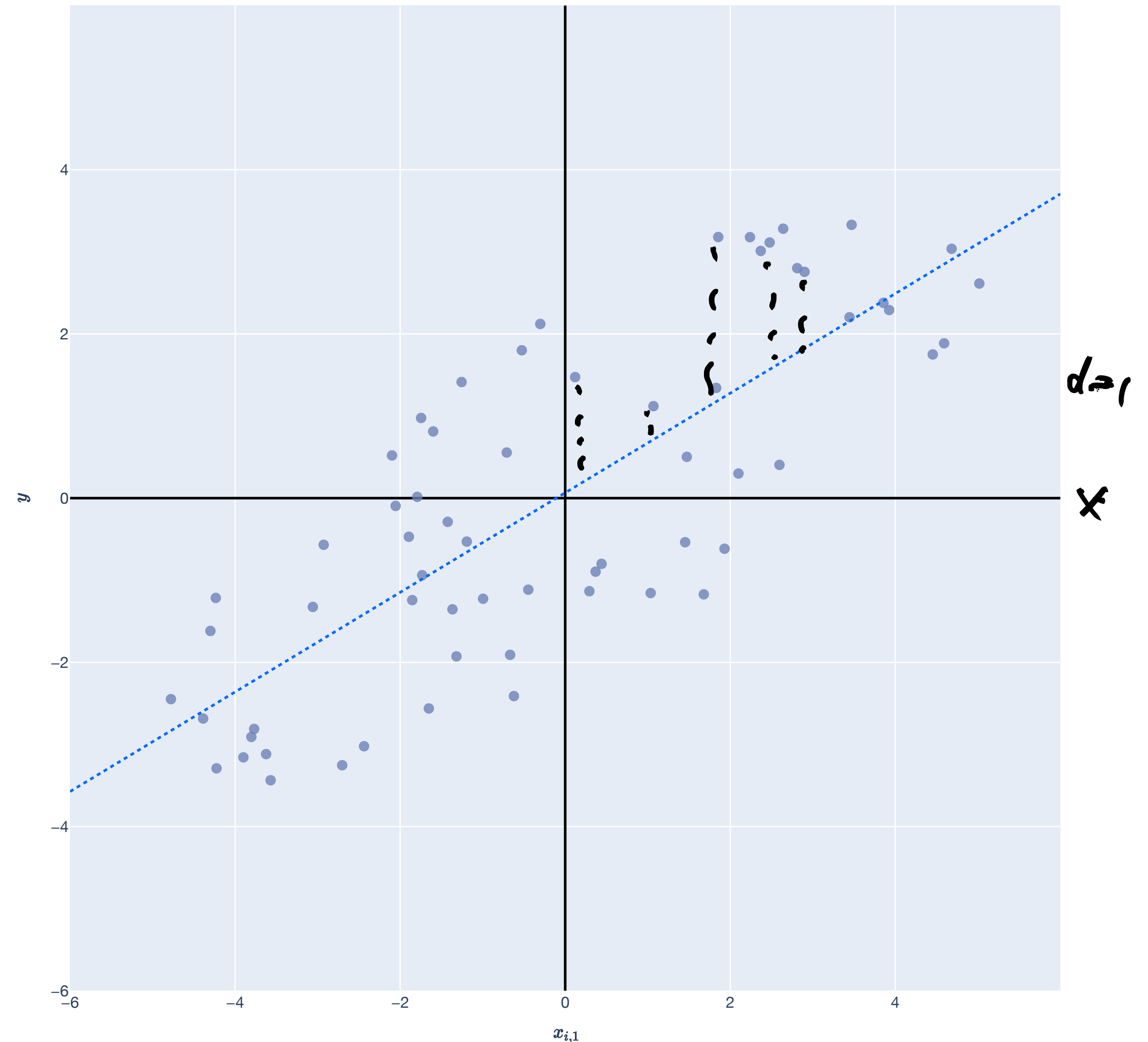
Comparison with OLS

1D Pictures

OLS: Observe data $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$.

Goal: Find best linear combination $\hat{\mathbf{w}} \in \mathbb{R}^d$ of $\mathbf{x}_1, \dots, \mathbf{x}_d$ such that

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



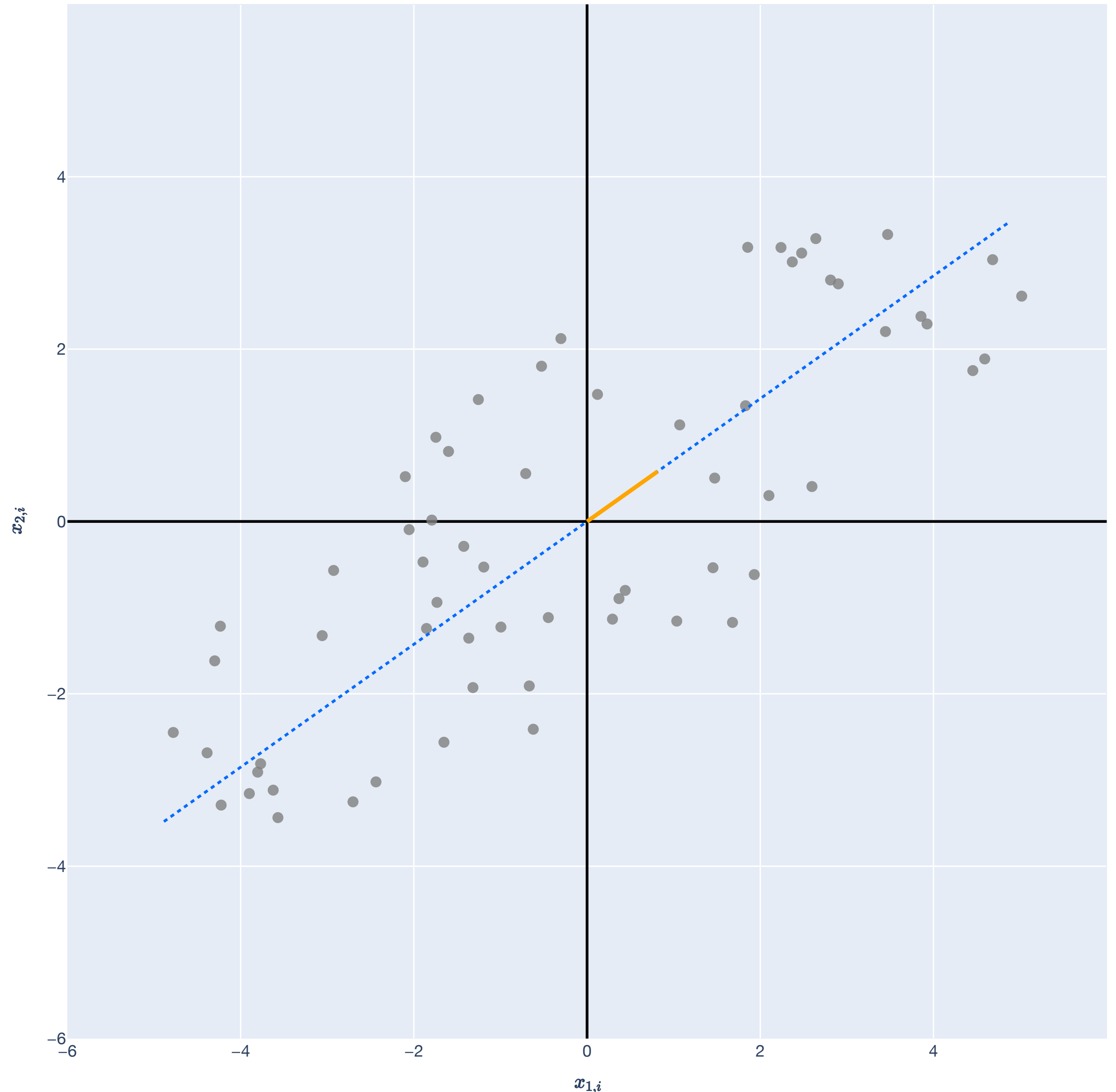
Comparison with OLS

1D Pictures

BFS: Observe data $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$.

Goal: Find one-dimensional subspace determined by $\mathbf{u} \in \mathbb{R}^n$ such that

$$\arg \min_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \|\mathbf{x}_i - \Pi_{\mathbf{u}}(\mathbf{x}_i)\|^2$$



Singular Value Decomposition (SVD)

Deriving 1D SVD


Find $\mathbf{u} \in \mathbb{R}^n$ that minimizes the sum of squared projection distances:

$$\arg \min_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \|\mathbf{x}_i - \Pi_{\mathbf{u}}(\mathbf{x}_i)\|^2.$$

Singular Value Decomposition (SVD)

Deriving 1D SVD

Find $\mathbf{u} \in \mathbb{R}^n$ that minimizes the sum of squared projection distances:

$$\arg \min_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \|\mathbf{x}_i - \underline{\Pi_{\mathbf{u}}(\mathbf{x}_i)}\|^2 = \sum_{i=1}^d \|\mathbf{x}_i - \underline{P_{\mathbf{u}}\mathbf{x}_i}\|^2.$$


Singular Value Decomposition (SVD)

Deriving 1D SVD

Find $\mathbf{u} \in \mathbb{R}^n$ that minimizes the sum of squared projection distances:

$$\arg \min_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \|\mathbf{x}_i - \Pi_{\mathbf{u}}(\mathbf{x}_i)\|^2 = \sum_{i=1}^d \|\mathbf{x}_i - P_{\mathbf{u}}\mathbf{x}_i\|^2.$$

What's $\|\mathbf{x}_i - P_{\mathbf{u}}\mathbf{x}_i\|^2$?

Singular Value Decomposition (SVD)

Deriving 1D SVD

$$\mathbf{x}_i \in \mathbb{R}^n$$

$$\frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \in \mathbb{R}^n$$

$$\begin{bmatrix} 1 \\ \vdots \\ u \\ \vdots \\ 1 \end{bmatrix} \begin{bmatrix} -u & - \end{bmatrix}$$

$n \times 1$ $1 \times n$ \rightarrow $n \times n$

Consider any $i \in [d]$. Then,

$$\|\mathbf{x}_i - P_{\mathbf{u}}\mathbf{x}_i\|^2 = \left\| \mathbf{x}_i - \left(\frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right) \mathbf{x}_i \right\|^2$$

(Prop: 1D projection formula)

Singular Value Decomposition (SVD)

Deriving 1D SVD

Consider any $i \in [d]$. Then,

$$\begin{aligned}\|\mathbf{x}_i - P_{\mathbf{u}}\mathbf{x}_i\|^2 &= \left\| \mathbf{x}_i - \left(\frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right) \mathbf{x}_i \right\|^2 \\ &= \left\| \underbrace{\left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right)}_{P_{\mathbf{u}^\perp}} \mathbf{x}_i \right\|^2\end{aligned}$$

$$P_{\mathbf{u}} + P_{\mathbf{u}^\perp} = \mathbf{I}$$

$$\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} = \mathbf{I} - P_{\mathbf{u}}$$

(Prop: 1D projection formula)

(Prop: Projection and Orthogonal Complement Matrices)

Singular Value Decomposition (SVD)

Deriving 1D SVD

Consider any $i \in [d]$. Then,

$$\|\mathbf{x}_i - P_{\mathbf{u}}\mathbf{x}_i\|^2 = \left\| \mathbf{x}_i - \left(\frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right) \mathbf{x}_i \right\|^2 \quad (\text{Prop: 1D projection formula})$$

$$= \left\| \left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right) \mathbf{x}_i \right\|^2$$
$$= \mathbf{x}_i^\top \underbrace{\left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right)^\top}_{\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}}} \underbrace{\left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right)}_{\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}}} \mathbf{x}_i$$

(Prop: Projection and Orthogonal Complement Matrices)

$$\|A\mathbf{x}\|^2 = (A\mathbf{x})^\top A\mathbf{x} = \mathbf{x}^\top \underbrace{A^\top A}_{\mathbf{I}}$$

Singular Value Decomposition (SVD) $(AB)^T = B^T A^T$

Deriving 1D SVD $(A^{-1})^T = (A^T)^{-1}$

Consider any $i \in [d]$. Then,

$$\begin{aligned} \|x_i - P_u x_i\|^2 &= \left\| x_i - \left(\frac{uu^T}{u^T u} \right) x_i \right\|^2 \\ &= \left\| \left(I - \frac{uu^T}{u^T u} \right) x_i \right\|^2 \\ &= x_i^T \left(I - \frac{uu^T}{u^T u} \right)^T \left(I - \frac{uu^T}{u^T u} \right) x_i \\ &= x_i^T \left(I - \frac{uu^T}{u^T u} \right)^2 x_i \end{aligned}$$

$$\begin{aligned} P &= X (X^T X)^{-1} X^T \\ P^T &= \left(X (X^T X)^{-1} X^T \right)^T \\ &= (X^T)^T \left(X (X^T X)^{-1} \right)^T \\ &\quad \text{(Prop: 1D projection formula)} \\ &= X \left((X^T X)^{-1} \right)^T X^T \\ &= X \left(\underbrace{(X^T X)^T} \right)^{-1} X^T = X (X^T X)^{-1} X^T \end{aligned}$$

(Prop: Projection and Orthogonal Complement Matrices)

$$I - \frac{uu^T}{u^T u} = \boxed{P_{u^\perp}}$$

$$\underline{P^T = P}$$

(Prop: Projections are symmetric)

Singular Value Decomposition (SVD)

Deriving 1D SVD

$$u \in \mathbb{R}^n$$

- u spans a one dim. subspace.
- P_{u^\perp} projects onto a $n-1$ dim. subspace

Consider any $i \in [d]$. Then,

$$\|x_i - P_u x_i\|^2 = \left\| x_i - \left(\frac{uu^T}{u^T u} \right) x_i \right\|^2$$

(Prop: 1D projection formula)

$$= \left\| \left(I - \frac{uu^T}{u^T u} \right) x_i \right\|^2$$

(Prop: Projection and Orthogonal Complement Matrices)

$$= x_i^T \left(I - \frac{uu^T}{u^T u} \right)^T \left(I - \frac{uu^T}{u^T u} \right) x_i$$

$$P^2 = P$$

$$= x_i^T \left(I - \frac{uu^T}{u^T u} \right)^2 x_i$$

(Prop: Projections are symmetric)

$$= x_i^T \left(I - \frac{uu^T}{u^T u} \right) x_i$$

P_{u^\perp}

$$\left\langle x_i^T P_{u^\perp} x_i \right\rangle$$

(Prop: Projecting twice doesn't do anything)

Singular Value Decomposition (SVD)

Deriving 1D SVD

Therefore, for any $i \in [d]$,

$$\| \mathbf{x}_i - P_{\mathbf{u}} \mathbf{x}_i \|^2 = \left(\mathbf{x}_i^\top \left(\mathbf{I} - \frac{\mathbf{u} \mathbf{u}^\top}{\mathbf{u}^\top \mathbf{u}} \right) \mathbf{x}_i \right)$$

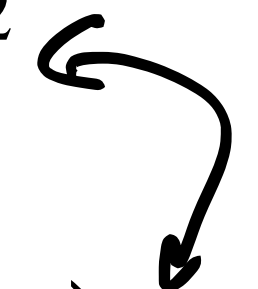
Singular Value Decomposition (SVD)

Deriving 1D SVD

Therefore, for any $i \in [d]$,

$$\|\mathbf{x}_i - P_{\mathbf{u}}\mathbf{x}_i\|^2 = \mathbf{x}_i^\top \left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right) \mathbf{x}_i$$

Find $\mathbf{u} \in \mathbb{R}^n$ that minimizes the sum of squared projection distances:

$$\begin{aligned} \arg \min_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \|\mathbf{x}_i - \Pi_{\mathbf{u}}(\mathbf{x}_i)\|^2 &= \sum_{i=1}^d \|\mathbf{x}_i - P_{\mathbf{u}}\mathbf{x}_i\|^2 \\ &= \left[\sum_{i=1}^d \mathbf{x}_i^\top \left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right) \mathbf{x}_i \right] \end{aligned}$$



Singular Value Decomposition (SVD)

Deriving 1D SVD

Therefore, for any $i \in [d]$,

$$\|\mathbf{x}_i - P_{\mathbf{u}}\mathbf{x}_i\|^2 = \mathbf{x}_i^\top \left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right) \mathbf{x}_i$$

Find $\mathbf{u} \in \mathbb{R}^n$ that minimizes the sum of squared projection distances:

$$\begin{aligned} \arg \min_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \|\mathbf{x}_i - \Pi_{\mathbf{u}}(\mathbf{x}_i)\|^2 &= \sum_{i=1}^d \|\mathbf{x}_i - P_{\mathbf{u}}\mathbf{x}_i\|^2 \\ &= \sum_{i=1}^d \mathbf{x}_i^\top \left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right) \mathbf{x}_i \\ &= \sum_{i=1}^d \mathbf{x}_i^\top \mathbf{x}_i - \mathbf{x}_i^\top \left(\frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right) \mathbf{x}_i \end{aligned}$$


Singular Value Decomposition (SVD)

Deriving 1D SVD

Therefore, for any $i \in [d]$,

$$\|\mathbf{x}_i - P_{\mathbf{u}}\mathbf{x}_i\|^2 = \mathbf{x}_i^T \left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^T}{\mathbf{u}^T\mathbf{u}} \right) \mathbf{x}_i.$$

Find $\mathbf{u} \in \mathbb{R}^n$ that minimizes the sum of squared projection distances:

$$\mathbf{u} = \arg \min_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \mathbf{x}_i^T \mathbf{x}_i - \mathbf{x}_i^T \left(\frac{\mathbf{u}\mathbf{u}^T}{\mathbf{u}^T\mathbf{u}} \right) \mathbf{x}_i$$

$$\Leftrightarrow \arg \max_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \mathbf{x}_i^T \left(\frac{\mathbf{u}\mathbf{u}^T}{\mathbf{u}^T\mathbf{u}} \right) \mathbf{x}_i$$

$P_{\mathbf{u}}$

$$\sum_{i=1}^d - \mathbf{x}_i^T \left(\frac{\mathbf{u}\mathbf{u}^T}{\mathbf{u}^T\mathbf{u}} \right) \mathbf{x}_i$$
$$= - \sum_{i=1}^d \mathbf{x}_i^T \left(\frac{\mathbf{u}\mathbf{u}^T}{\mathbf{u}^T\mathbf{u}} \right) \mathbf{x}_i$$

$$\sum_{i=1}^d \mathbf{x}_i^T P_{\mathbf{u}} \mathbf{x}_i$$

Singular Value Decomposition (SVD)

Deriving 1D SVD

Therefore, for any $i \in [d]$,

$$\|\mathbf{x}_i - P_{\mathbf{u}}\mathbf{x}_i\|^2 = \mathbf{x}_i^\top \left(\mathbf{I} - \frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right) \mathbf{x}_i.$$

Find $\mathbf{u} \in \mathbb{R}^n$ that minimizes the sum of squared projection distances:

$$\mathbf{u} = \arg \min_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \mathbf{x}_i^\top \mathbf{x}_i - \mathbf{x}_i^\top \left(\frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right) \mathbf{x}_i$$

$$\iff \arg \max_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \mathbf{x}_i^\top \left(\frac{\mathbf{u}\mathbf{u}^\top}{\mathbf{u}^\top\mathbf{u}} \right) \mathbf{x}_i = \arg \max_{\mathbf{u} \in \mathbb{R}^n} \frac{\mathbf{u}^\top \mathbf{X} \mathbf{X}^\top \mathbf{u}}{\mathbf{u}^\top \mathbf{u}} \text{ (AKA: squared operator norm of } \mathbf{X}, \text{ i.e. } \|\mathbf{X}\|_{op}^2 \text{).}$$

Singular Value Decomposition (SVD)

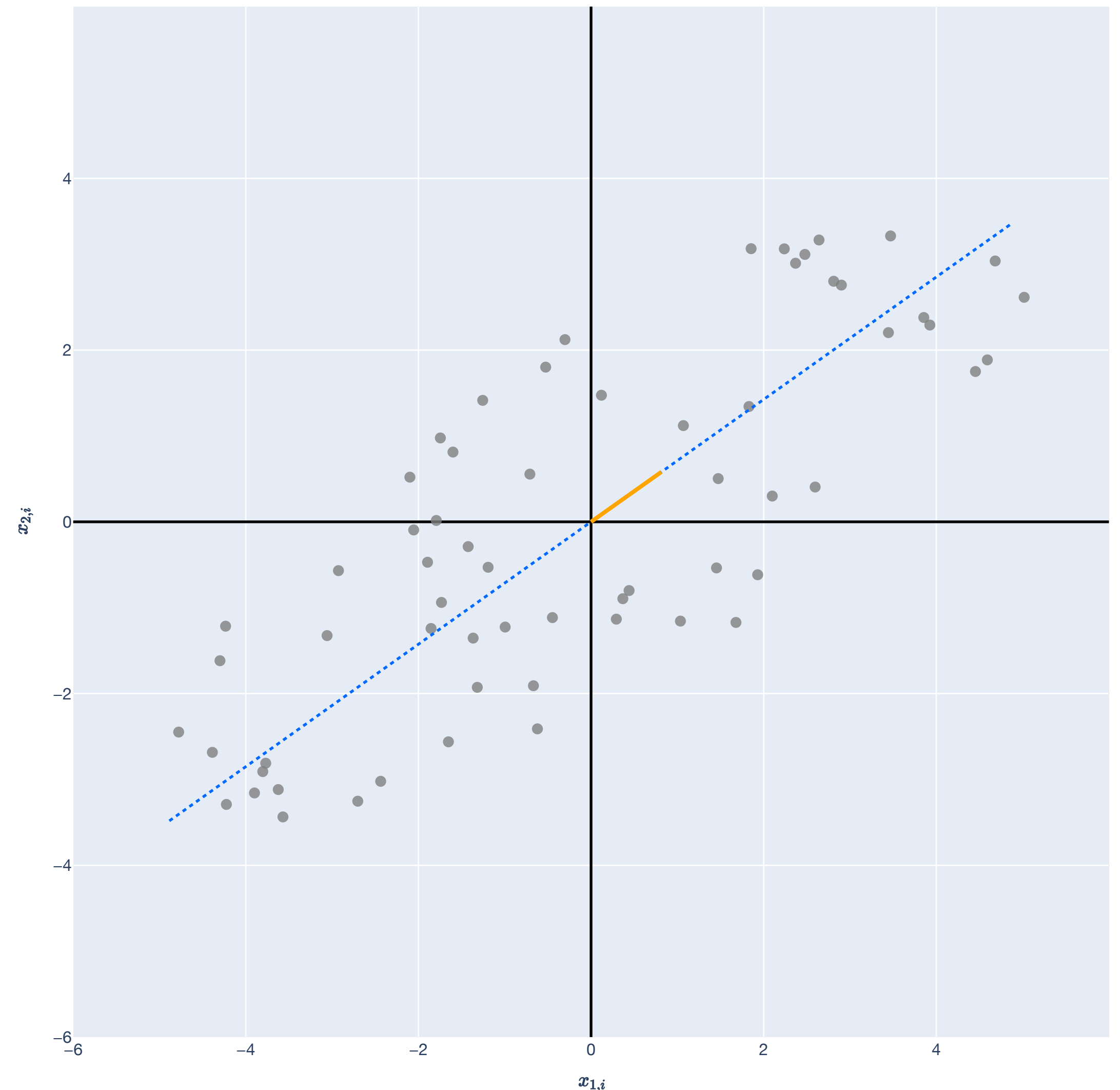
1D Picture

Observe data $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$.

Goal: Find the best one-dimensional subspace $\mathcal{U} \subseteq \mathbb{R}^n$ that fits the points.

How? Find $\mathbf{u} \in \mathbb{R}^n$ that minimizes the sum of squared projection distances:

$$\arg \min_{\mathbf{u} \in \mathbb{R}^n} \sum_{i=1}^d \|\mathbf{x}_i - \Pi_{\mathbf{u}}(\mathbf{x}_i)\|^2.$$



Singular Value Decomposition (SVD)

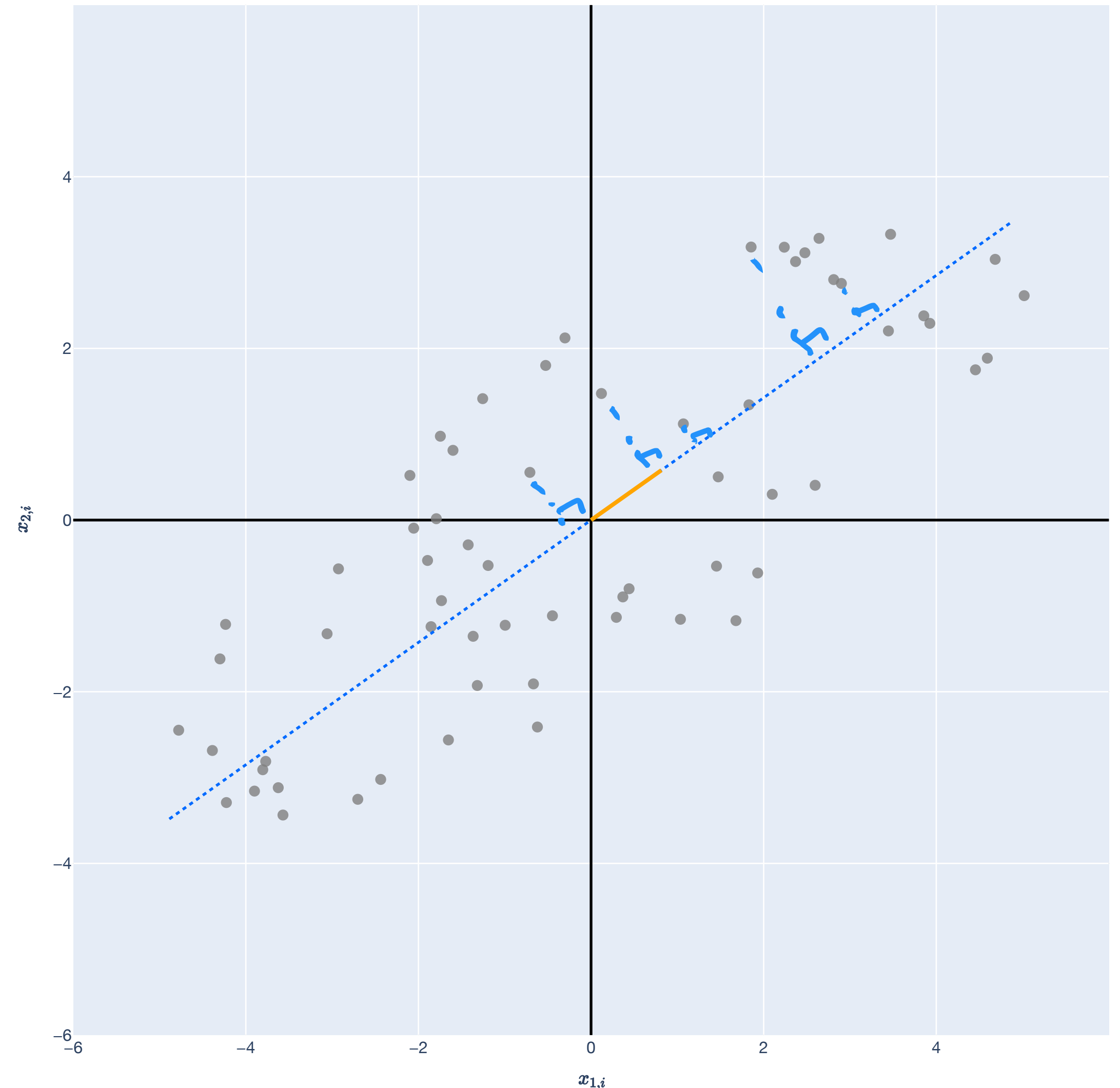
1D Picture

Observe data $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$.

Goal: Find the best one-dimensional subspace $\mathcal{U} \subseteq \mathbb{R}^n$ that fits the points.

How? Find $\mathbf{u} \in \mathbb{R}^n$ that minimizes the sum of squared projection distances:

$$\arg \max_{\mathbf{u} \in \mathbb{R}^n} \frac{\mathbf{u}^\top \mathbf{X} \mathbf{X}^\top \mathbf{u}}{\mathbf{u}^\top \mathbf{u}}.$$



$$P = X(X^T X)^{-1} X^T$$

Singular Value Decomposition (SVD)

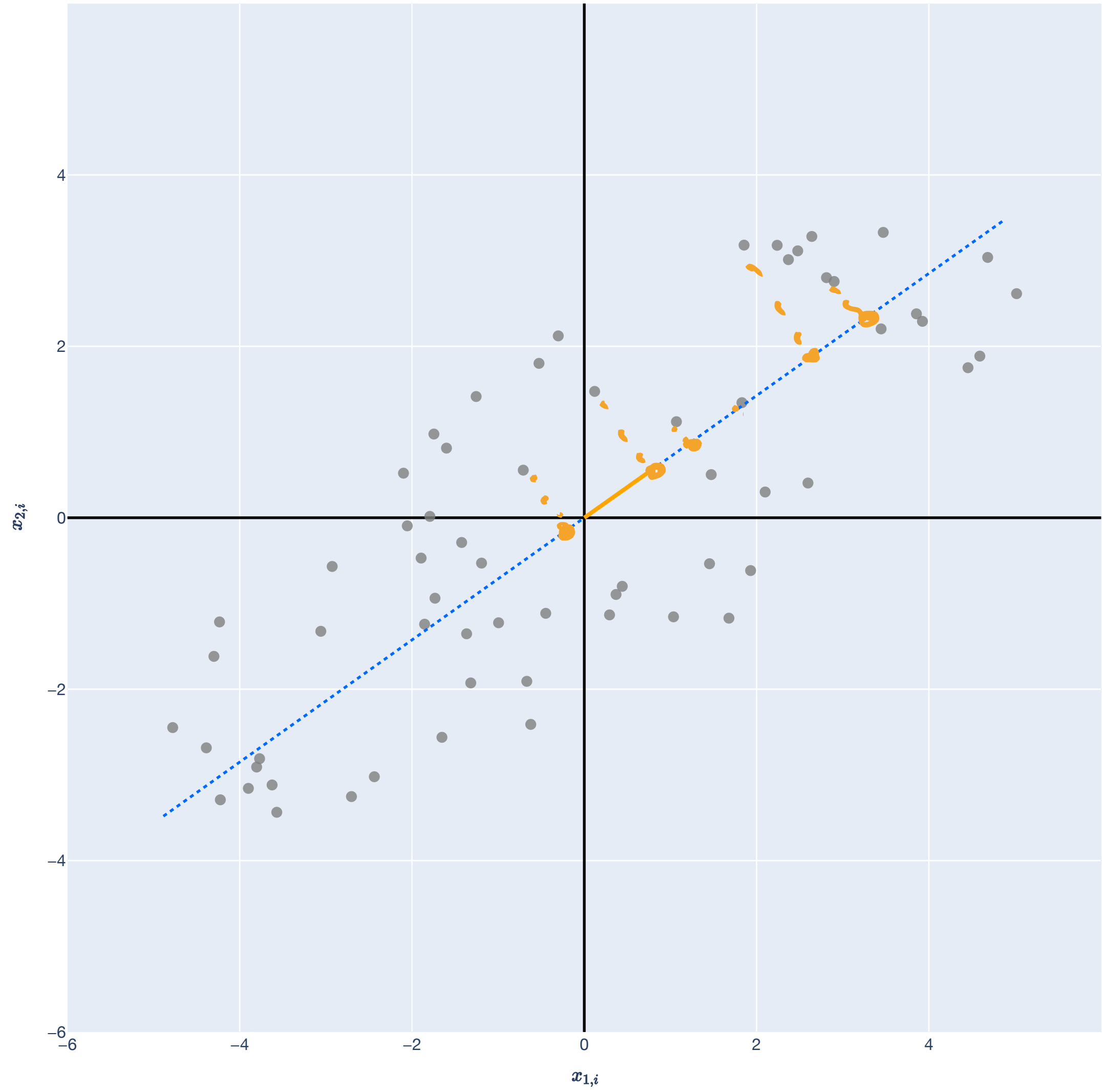
1D Picture

Find $\mathbf{u} \in \mathbb{R}^n$ that minimizes the sum of squared projection distances:

$$\arg \max_{\mathbf{u} \in \mathbb{R}^n} \frac{\mathbf{u}^T \mathbf{X} \mathbf{X}^T \mathbf{u}}{\mathbf{u}^T \mathbf{u}}$$

The vector $\mathbf{u} \in \mathbb{R}^n$ that achieves this maximum is the 1st left singular vector.

The value $\frac{\mathbf{u}^T \mathbf{X} \mathbf{X}^T \mathbf{u}}{\mathbf{u}^T \mathbf{u}}$ is σ_1^2 , the squared 1st singular value of \mathbf{X} .



Singular Value Decomposition

Definition of Full SVD and Compact SVD

Singular Value Decomposition (SVD)

Building up the SVD

Observe data $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$. Consider the following procedure...

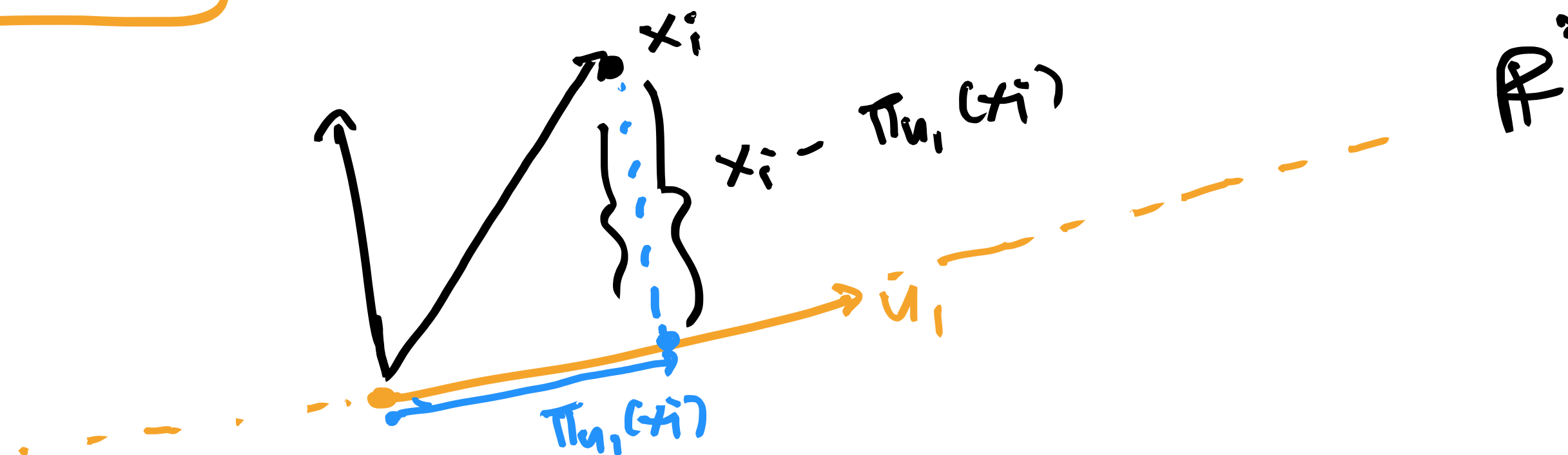
For $t = 1, 2, \dots, n \dots$

1. Find $\mathbf{u}_1 \in \mathbb{R}^n$, the best one-dimensional subspace fit to $\mathbf{x}_1, \dots, \mathbf{x}_d$.

$\operatorname{argmax}_u \frac{u^T X X^T u}{u^T u}$

• Let $\mathbf{x}_i^{(1)} = \mathbf{x}_i - \Pi_{\mathbf{u}_1}(\mathbf{x}_i)$.

← Part of each \mathbf{x}_i "unexplained" by \vec{u}_1 .



Singular Value Decomposition (SVD)

Building up the SVD

Observe data $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$. Consider the following procedure...

For $t = 1, 2, \dots, n \dots$

1. Find $\mathbf{u}_1 \in \mathbb{R}^n$, the best one-dimensional subspace fit to $\mathbf{x}_1, \dots, \mathbf{x}_d$.

- Let $\mathbf{x}_i^{(1)} = \mathbf{x}_i - \Pi_{\mathbf{u}_1}(\mathbf{x}_i)$.

2. Find $\mathbf{u}_2 \in \mathbb{R}^n$, the best one-dimensional subspace fit to $\mathbf{x}_1^{(1)}, \dots, \mathbf{x}_d^{(1)}$.

- Let $\mathbf{x}_i^{(2)} = \mathbf{x}_i^{(1)} - \Pi_{\mathbf{u}_2}(\mathbf{x}_i^{(1)}) = \underbrace{\mathbf{x}_i}_{\text{original}} - \underbrace{\Pi_{\mathbf{u}_1}(\mathbf{x}_i)}_{\text{projection onto } \mathbf{u}_1} - \underbrace{\Pi_{\mathbf{u}_2}(\mathbf{x}_i^{(1)})}_{\text{projection onto } \mathbf{u}_2}$.

Apply to modified $\mathbf{x}_1^{(1)}, \dots, \mathbf{x}_d^{(1)}$

Singular Value Decomposition (SVD)

Building up the SVD

Observe data $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^n$. Consider the following procedure...

For $t = 1, 2, \dots, n$...

1. Find $\mathbf{u}_1 \in \mathbb{R}^n$, the best one-dimensional subspace fit to $\mathbf{x}_1, \dots, \mathbf{x}_d$.
 - Let $\mathbf{x}_i^{(1)} = \mathbf{x}_i - \Pi_{\mathbf{u}_1}(\mathbf{x}_i)$.
2. Find $\mathbf{u}_2 \in \mathbb{R}^n$, the best one-dimensional subspace fit to $\mathbf{x}_1^{(1)}, \dots, \mathbf{x}_d^{(1)}$.
 - Let $\mathbf{x}_i^{(2)} = \mathbf{x}_i^{(1)} - \Pi_{\mathbf{u}_2}(\mathbf{x}_i^{(1)}) = \mathbf{x}_i - \Pi_{\mathbf{u}_1}(\mathbf{x}_i) - \Pi_{\mathbf{u}_2}(\mathbf{x}_i)$.
3. Find $\mathbf{u}_3 \in \mathbb{R}^n$, the best one-dimensional subspace fit to $\mathbf{x}_1^{(2)}, \dots, \mathbf{x}_d^{(2)}$...

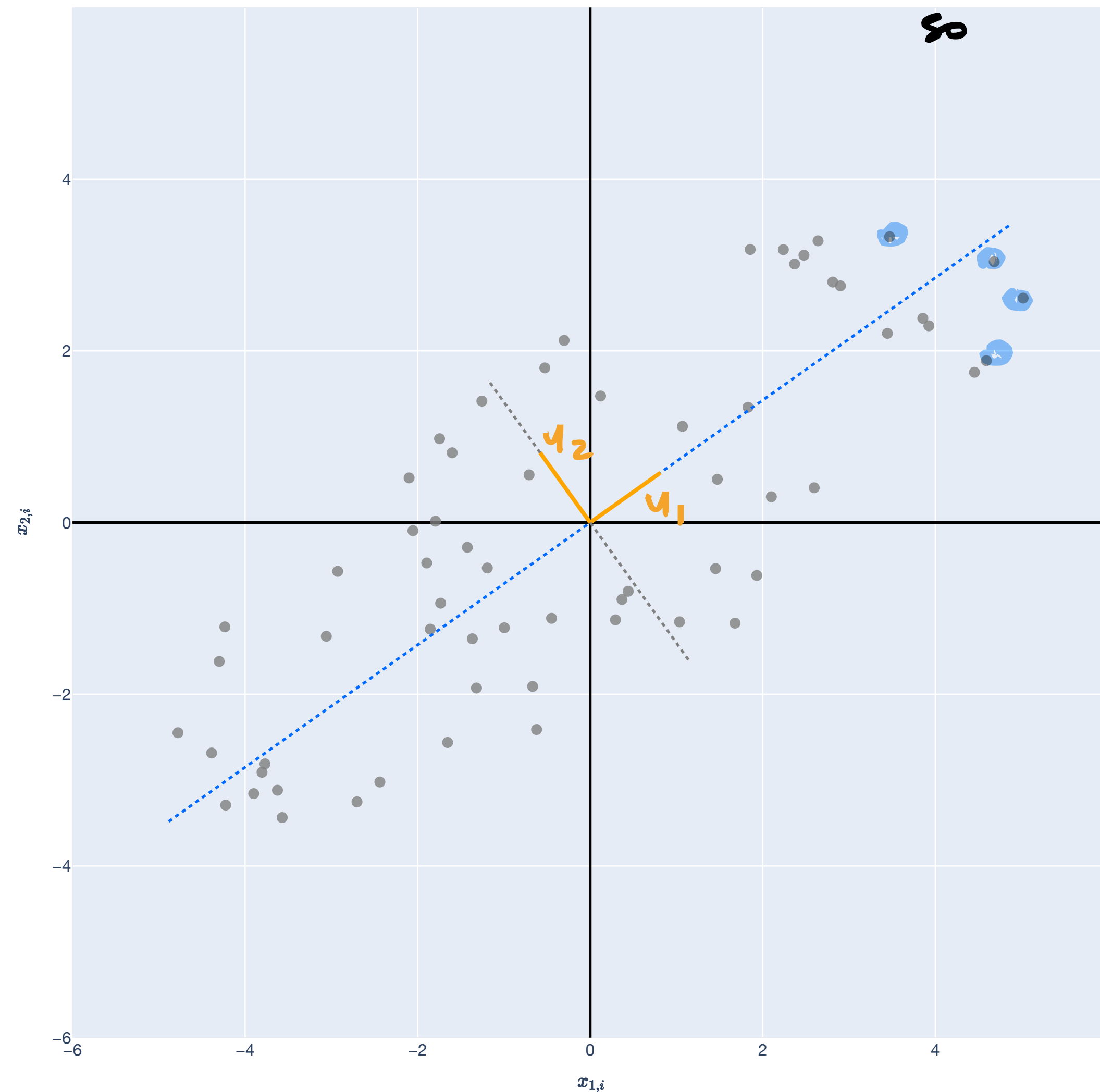
Singular Value Decomposition (SVD)

Building up the SVD

$$n = 2$$
$$d = 50$$

$$X = \left[\begin{array}{c|c} \text{blue} & \text{grey} \end{array} \right]$$

Observe data $\mathbf{x}_1, \dots, \mathbf{x}_d \in \mathbb{R}^2$.
Then, \mathbf{u}_1 and \mathbf{u}_2 look like...



Singular Value Decomposition (SVD)

Building up the SVD

Find $\mathbf{u}_t \in \mathbb{R}^n$, the best one-dimensional subspace fit to:

$$\mathbf{x}_i - \sum_{k=1}^{t-1} \Pi_{\mathbf{u}_k}(\mathbf{x}_i).$$

These are the n left singular vectors of $\mathbf{X} \in \mathbb{R}^{n \times d}$.

$$\mathbf{u}_1, \dots, \mathbf{u}_n$$

The n left singular vectors are orthogonal, by construction (left singular vector \mathbf{u}_k is in the orthogonal complement of $\mathbf{u}_1, \dots, \mathbf{u}_{k-1}$).

Singular Value Decomposition (SVD)

Definition of the Full SVD

Consider any matrix $X \in \mathbb{R}^{n \times d}$. By the full singular value decomposition (SVD), there exist matrices U, Σ, V such that

$$U = \begin{bmatrix} | & & | \\ u_1 & \dots & u_n \\ | & & | \end{bmatrix}$$

$$\underbrace{X}_{n \times d} = \underbrace{U}_{n \times n} \underbrace{\Sigma}_{n \times d} \underbrace{V^T}_{d \times d} \quad \leftarrow \text{"factored"}$$

$\sigma_1, \dots, \sigma_d \rightarrow$ values $\frac{u^T X X^T u}{u^T u}$

The columns of $U \in \mathbb{R}^{n \times n}$ are the left singular vectors and U is orthogonal, i.e. $U^T U = \underline{U U^T} = \underline{I}$.
Because u_1, \dots, u_n are orthogonal.

The columns of $V \in \mathbb{R}^{d \times d}$ are the right singular vectors and V is orthogonal, i.e. $V^T V = V V^T = I$.

$$V = \begin{bmatrix} | & & | \\ v_1 & \dots & v_d \\ | & & | \end{bmatrix}$$

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

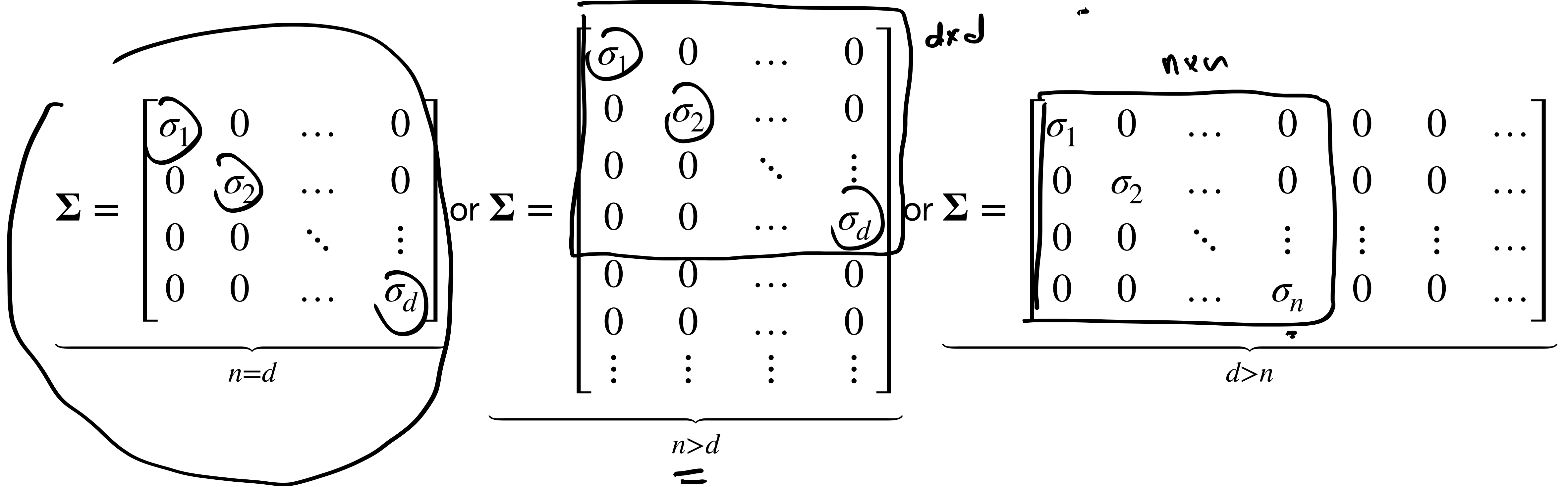
$$\begin{bmatrix} \sigma_1 & & & \\ & \dots & & \\ & & \sigma_d & \\ & & & 0 \end{bmatrix}$$

$$\text{rank}(X) = \# \text{ singular values } \sigma_i > 0$$

Singular Value Decomposition (SVD)

Shape of the Σ Matrix

$\Sigma \in \mathbb{R}^{n \times d}$ is a diagonal matrix with singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$, with $r \leq \min\{n, d\}$.



Interpreting the SVD

Example in \mathbb{R}^2

$$X = \left[\begin{array}{c|c|c|c} \begin{array}{c} | \\ x_1 \\ | \end{array} & \begin{array}{c} | \\ x_2 \\ | \end{array} & \dots & \begin{array}{c} | \\ x_{212} \\ | \end{array} \end{array} \right] \left. \vphantom{\begin{array}{c} | \\ x_1 \\ | \end{array}} \right\} \begin{array}{l} n=2 \\ d=2(2). \end{array}$$

Let $\mathbf{x}_1, \dots, \mathbf{x}_{212} \in \mathbb{R}^2$. The SVD is given by:

$$\underbrace{\mathbf{X}}_{2 \times 212} = \underbrace{\mathbf{U}}_{2 \times 2} \underbrace{\mathbf{\Sigma}}_{2 \times 212} \underbrace{\mathbf{V}^T}_{212 \times 212}$$

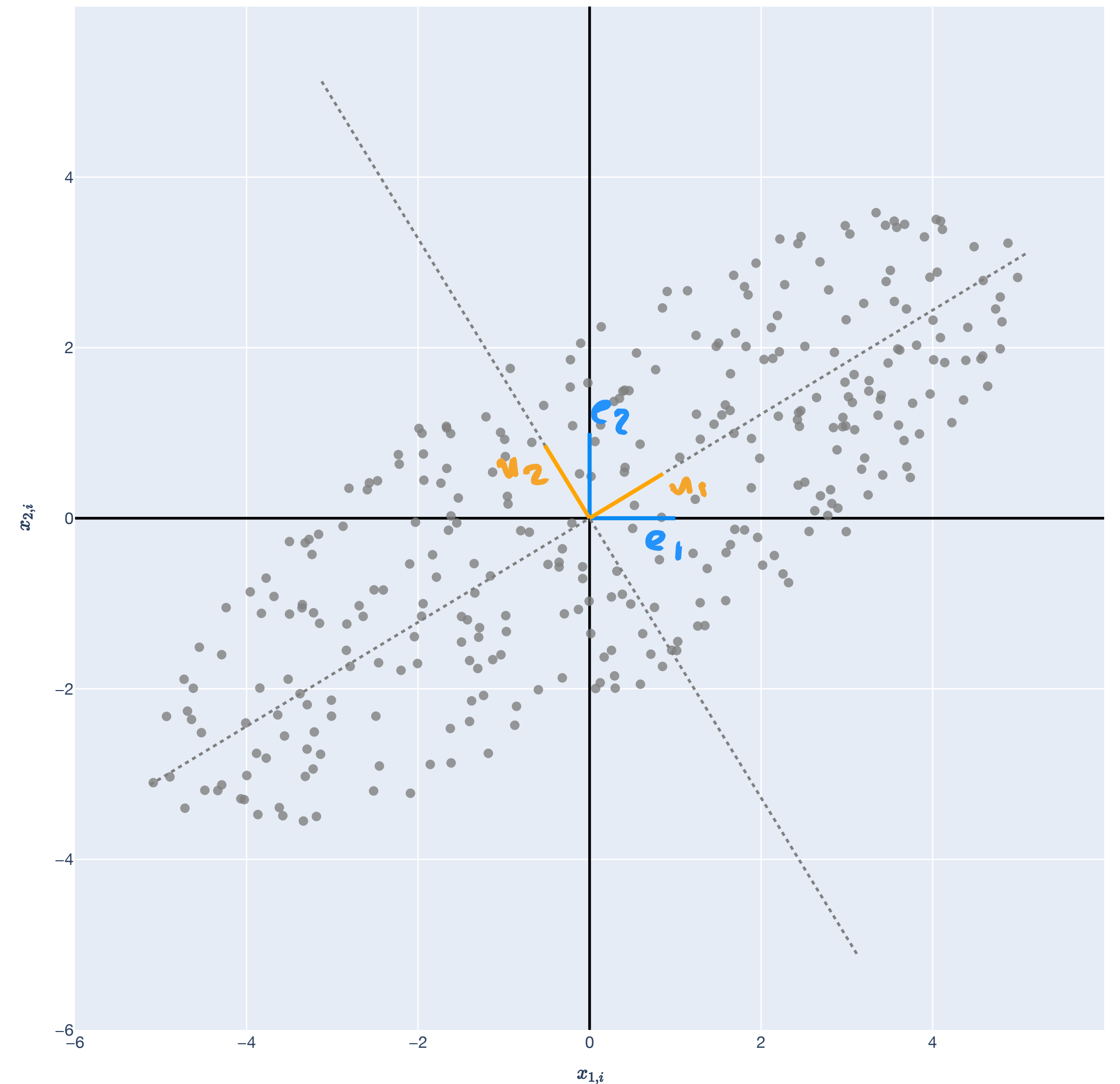
Left Singular Vectors

Interpreting the U matrix

$$\underbrace{\mathbf{X}}_{2 \times 212} = \underbrace{\mathbf{U}}_{2 \times 2} \underbrace{\mathbf{\Sigma}}_{2 \times 212} \underbrace{\mathbf{V}^T}_{212 \times 212}$$

The columns $\mathbf{u}_1, \mathbf{u}_2 \in \mathbb{R}^2$ of \mathbf{U} are an orthonormal basis for $\text{span}(\text{col}(\mathbf{X}))$.

$$\begin{bmatrix} | & | \\ \mathbf{u}_1 & \mathbf{u}_2 \\ | & | \end{bmatrix}_{2 \times 2}$$



Singular Values

Interpreting the Σ matrix

$$\underbrace{\mathbf{X}}_{2 \times 212} \underbrace{\vec{y}}_{212} = \underbrace{\mathbf{U}}_{2 \times 2} \underbrace{\mathbf{\Sigma}}_{2 \times 212} \underbrace{\mathbf{V}^T}_{212 \times 212} \underbrace{\vec{y}}_{212}$$

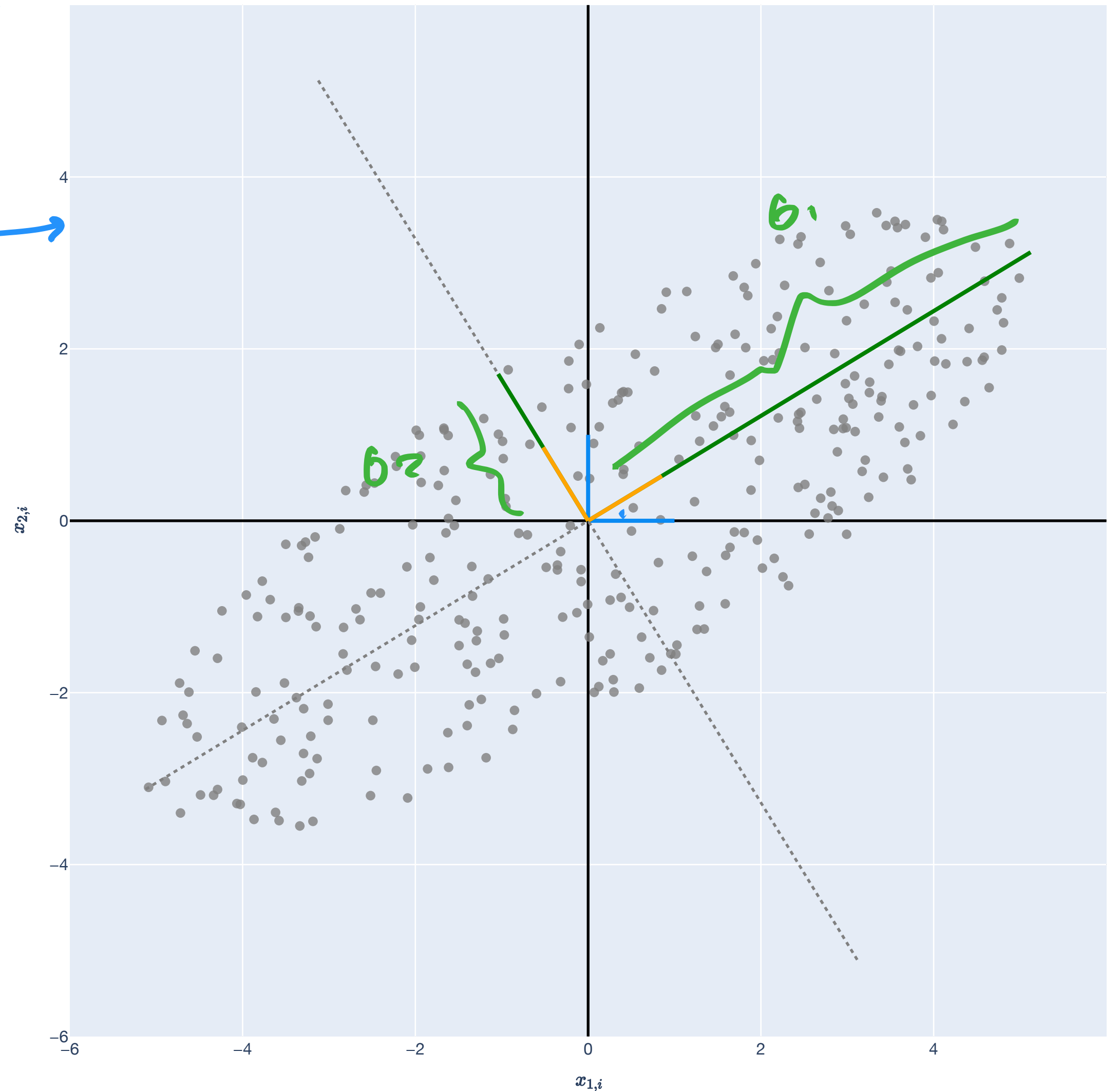
$\vec{w} = \mathbf{U}^T \vec{y}$

$$\mathbf{U}\mathbf{\Sigma} = \begin{bmatrix} u_1 & u_2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \sigma_1 & 0 & \vdots & 0 & 0 \\ 0 & \sigma_2 & \vdots & 0 & 0 \end{bmatrix}$$

$$= \begin{bmatrix} \sigma_1 u_1 & \sigma_2 u_2 & \vdots & 0 & 0 \\ \sigma_1 u_1 & \sigma_2 u_2 & \vdots & 0 & 0 \end{bmatrix}$$

The singular values $\sigma_1, \sigma_2 > 0$ represent how to scale \mathbf{u}_1 and \mathbf{u}_2 to “fit” all the data.

They represent the relative “strength” of \mathbf{u}_1 and \mathbf{u}_2 in explaining the data.



Right Singular Vectors

Interpreting the V matrix

$$\underbrace{\mathbf{X}}_{2 \times 212} = \underbrace{\mathbf{U}}_{2 \times 2} \underbrace{\mathbf{\Sigma}}_{2 \times 212} \underbrace{\mathbf{V}^T}_{212 \times 212}$$

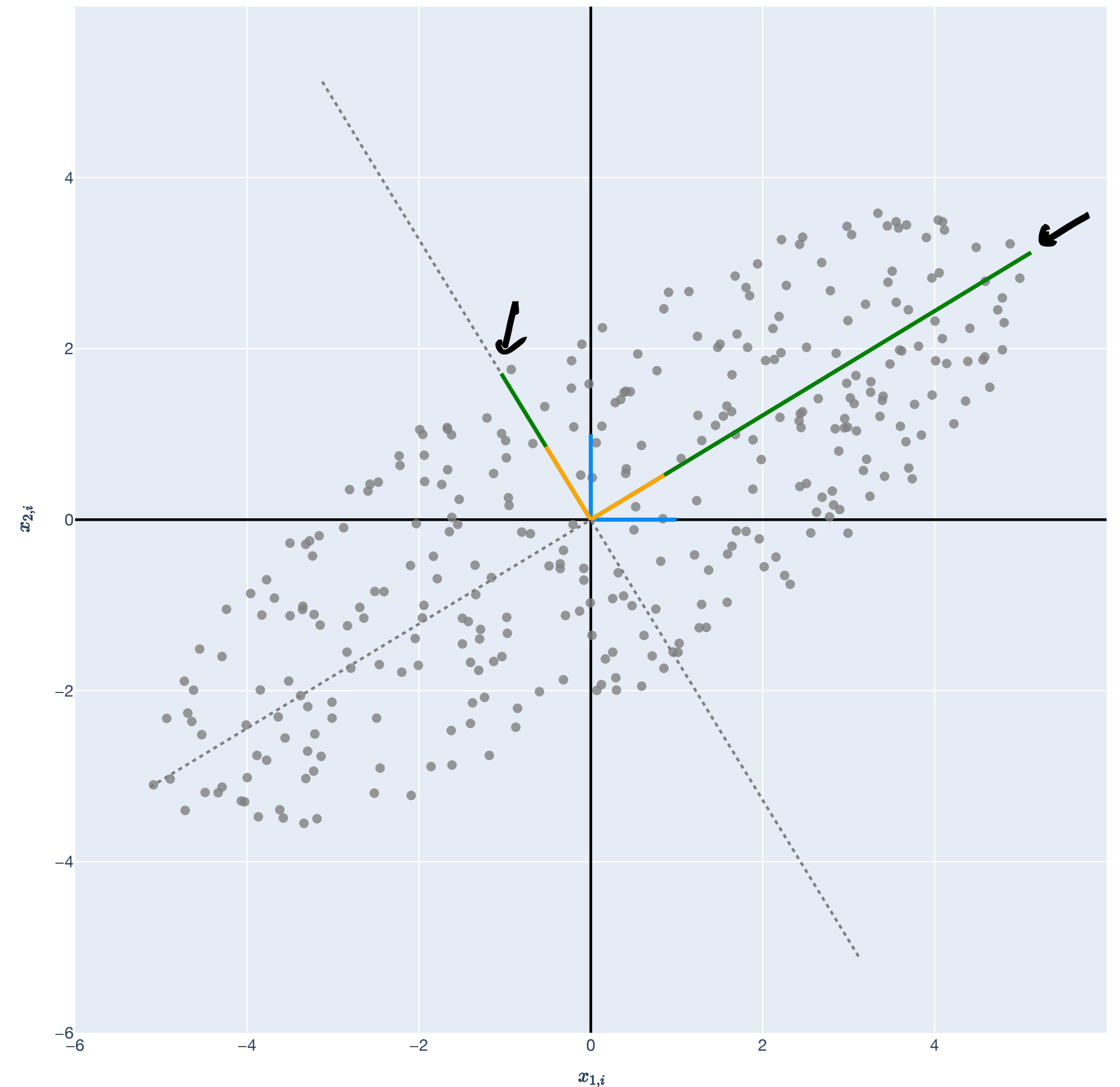
The rows of \mathbf{V} give the coordinates for each point under the basis $\sigma_1 \mathbf{u}_1, \sigma_2 \mathbf{u}_2$.

Specifically, for $j \in [d]$,

$$\boxed{\mathbf{x}_j} = \underline{v_{1j} \sigma_1 \mathbf{u}_1 + v_{2j} \sigma_2 \mathbf{u}_2}.$$

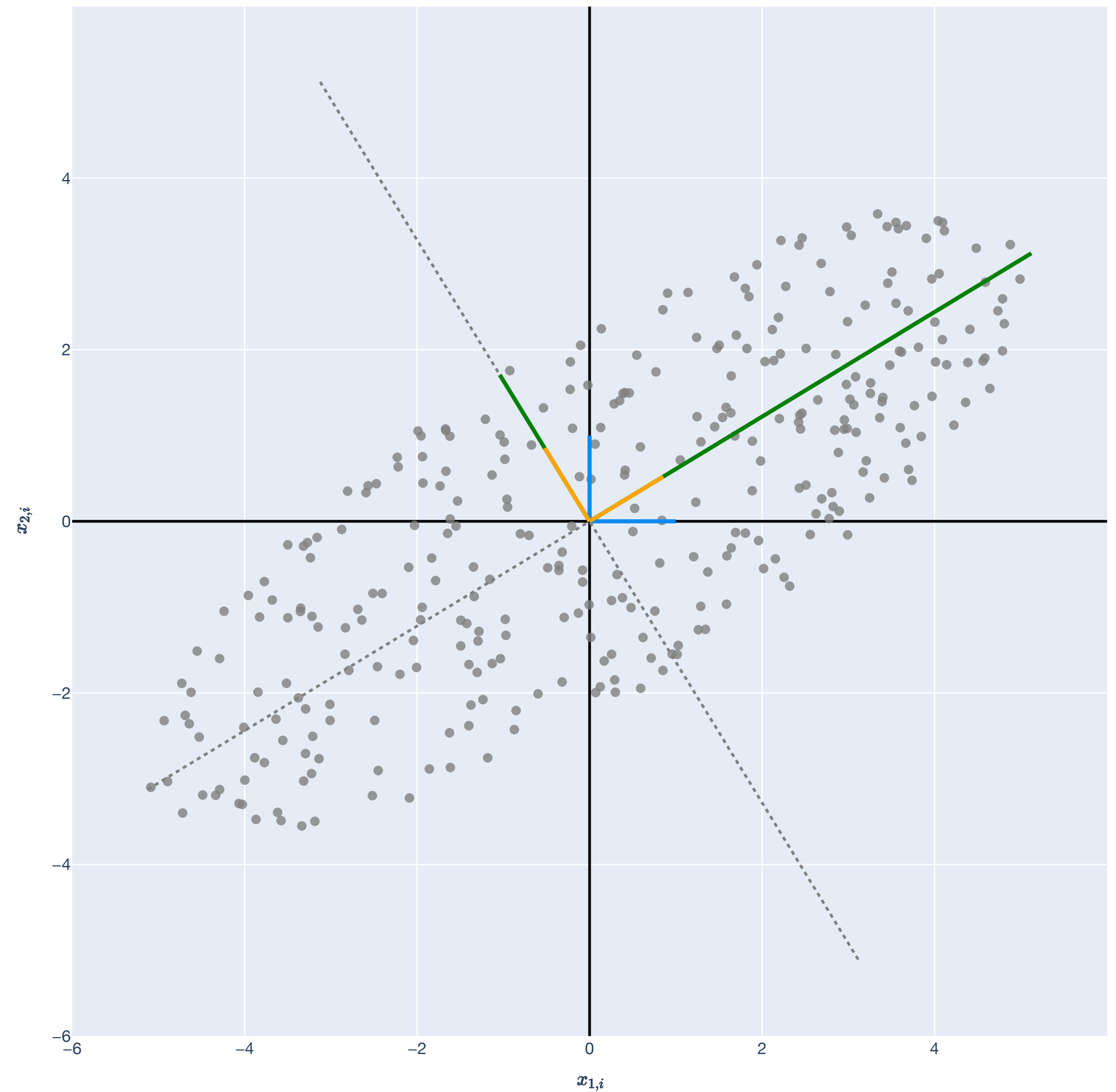
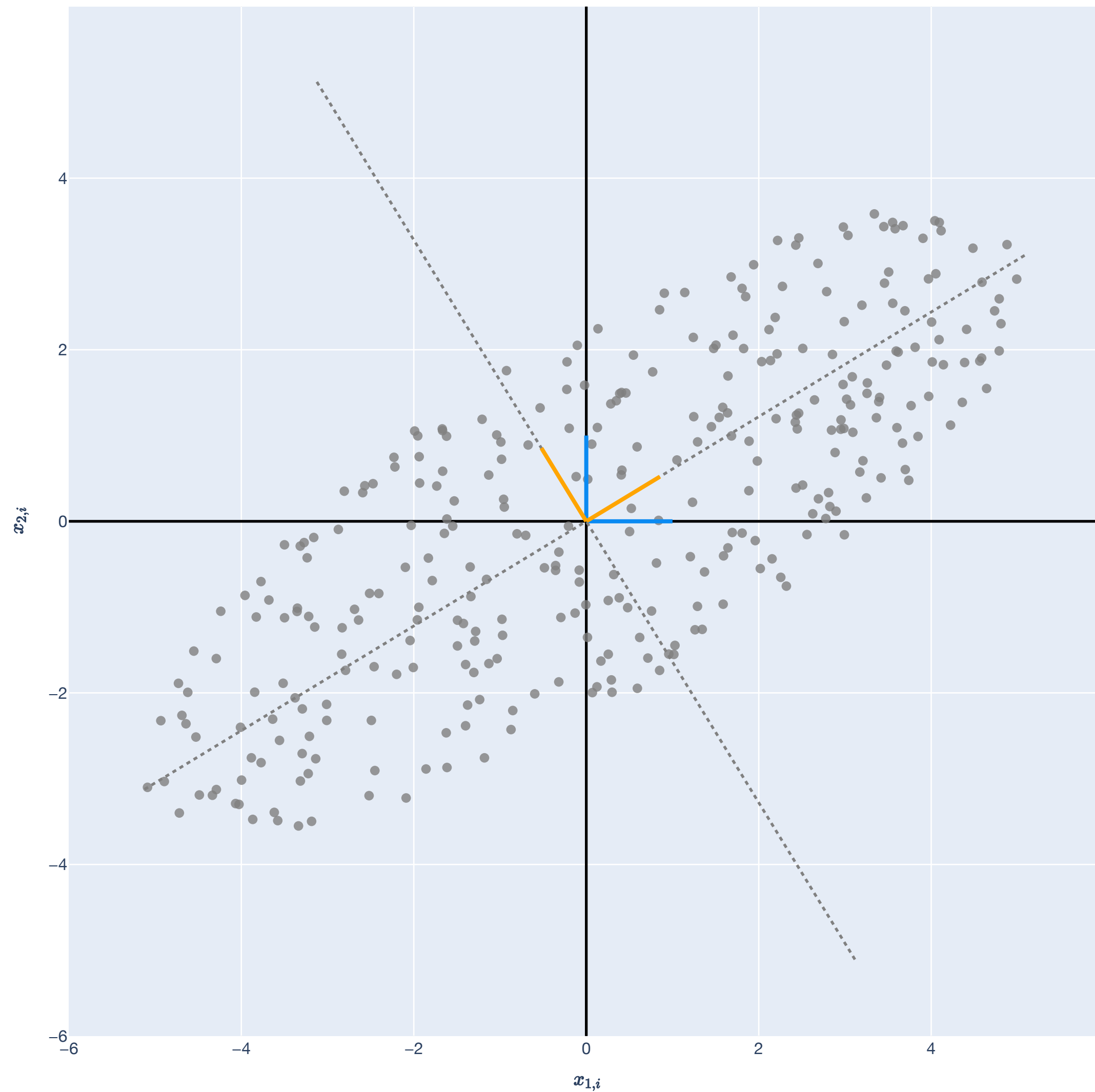
$$\mathbf{V} = \begin{bmatrix} \vdots & v_{1j} & \vdots \\ \vdots & v_{2j} & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}$$

doesn't matter.



Interpretation of the SVD

Full Interpretation of the SVD

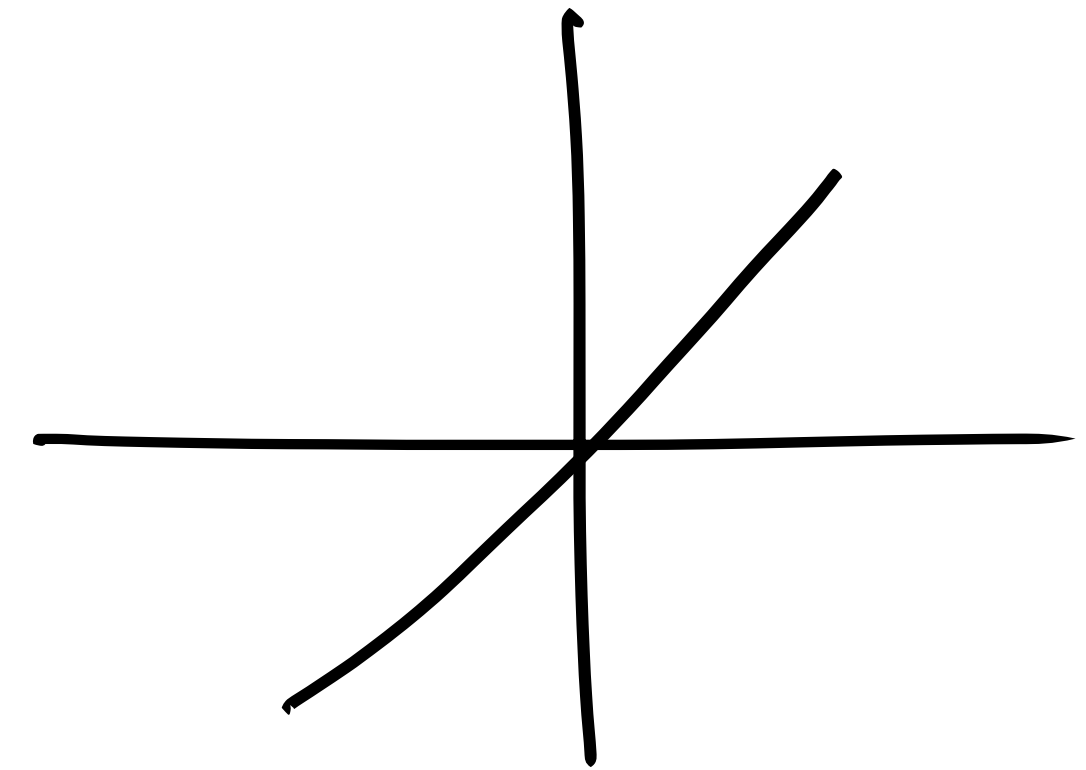


Singular Value Decomposition (SVD)

Example of SVD

$U, \Sigma, V^T = \text{np.linalg.svd}(X)$

$$X = \begin{matrix} 3 \times 3 & & n=3 \\ & & d=3 \end{matrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 10 \end{bmatrix}$$



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

$$U = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 10 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$$V^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

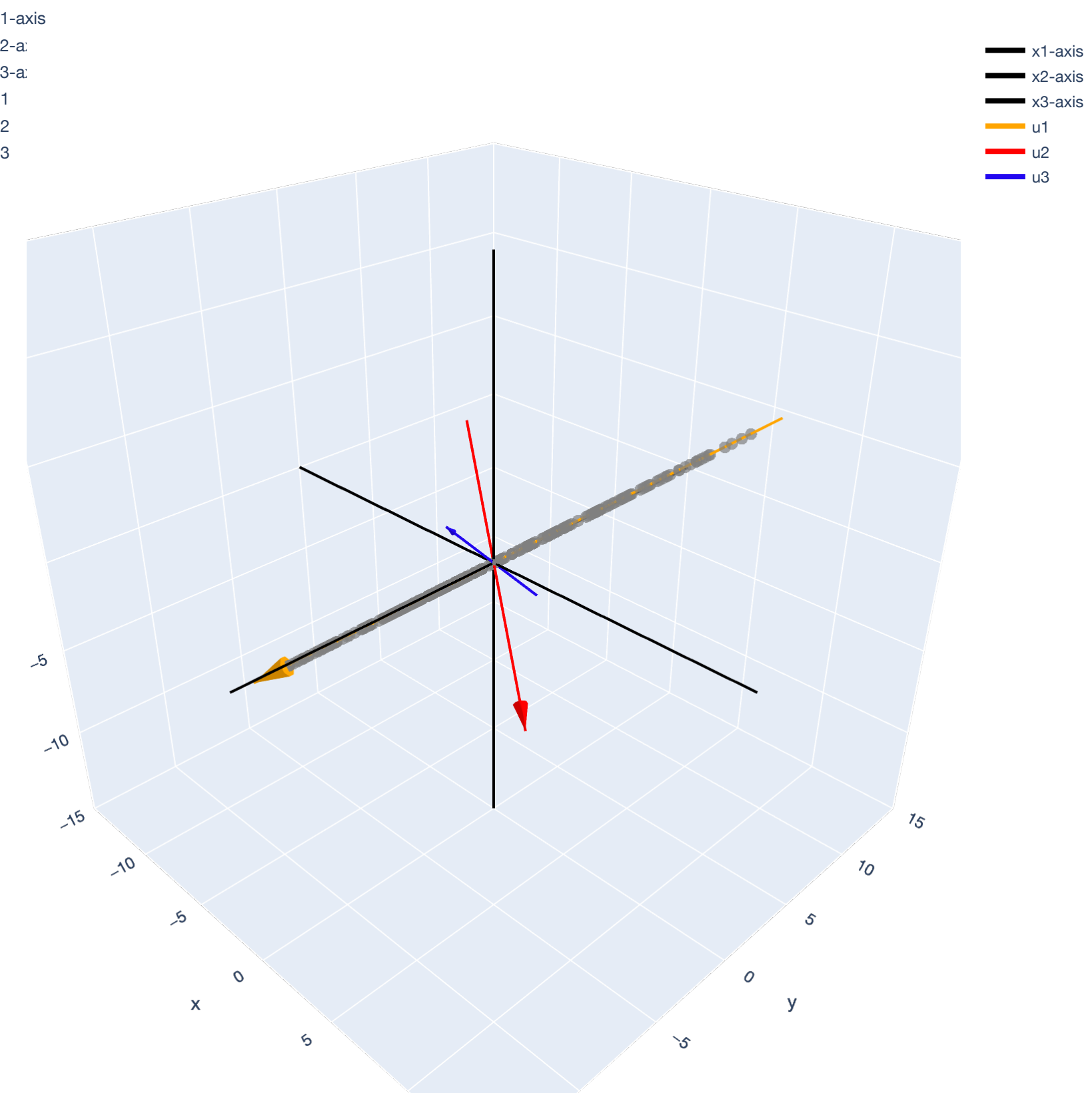
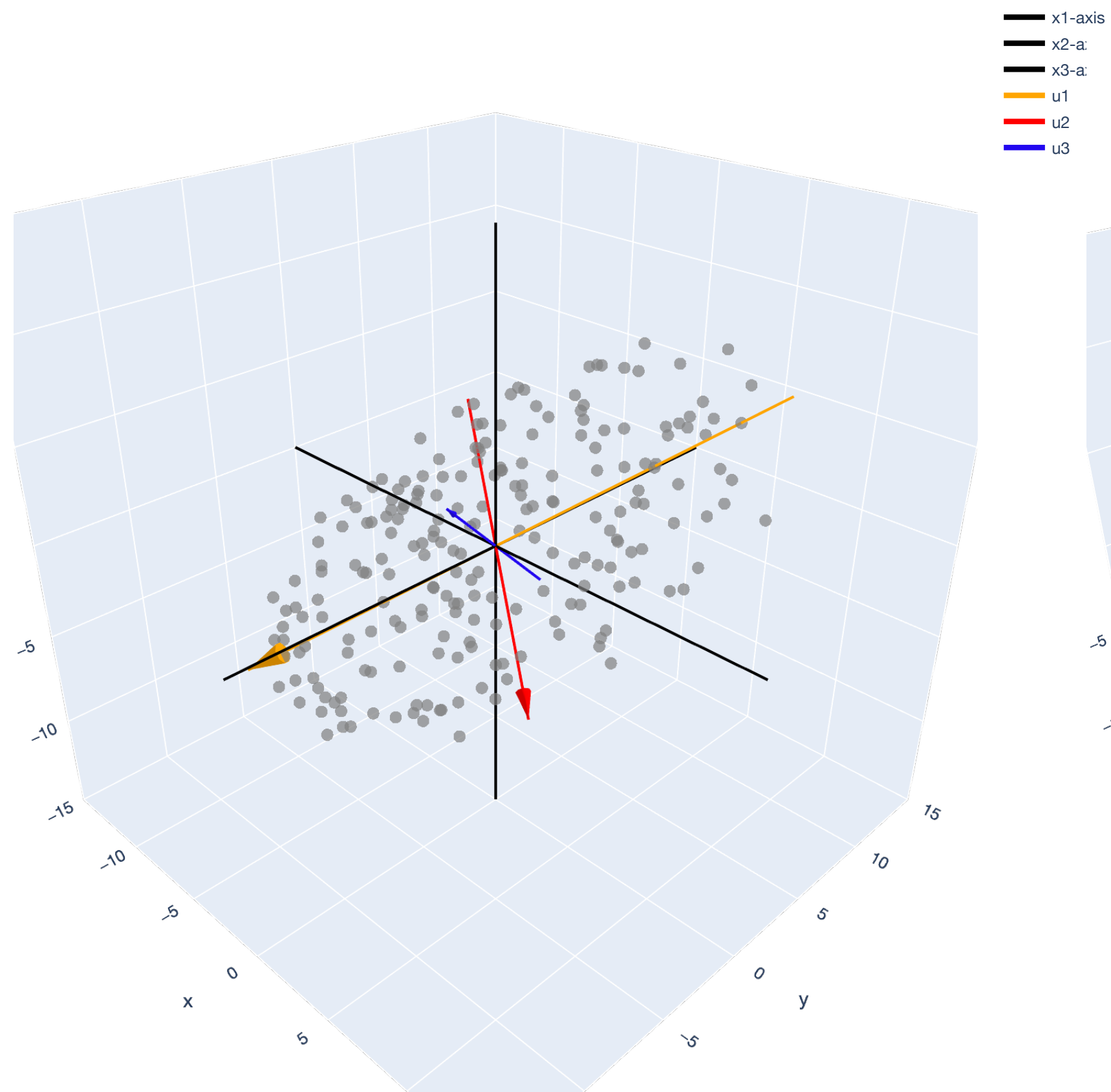
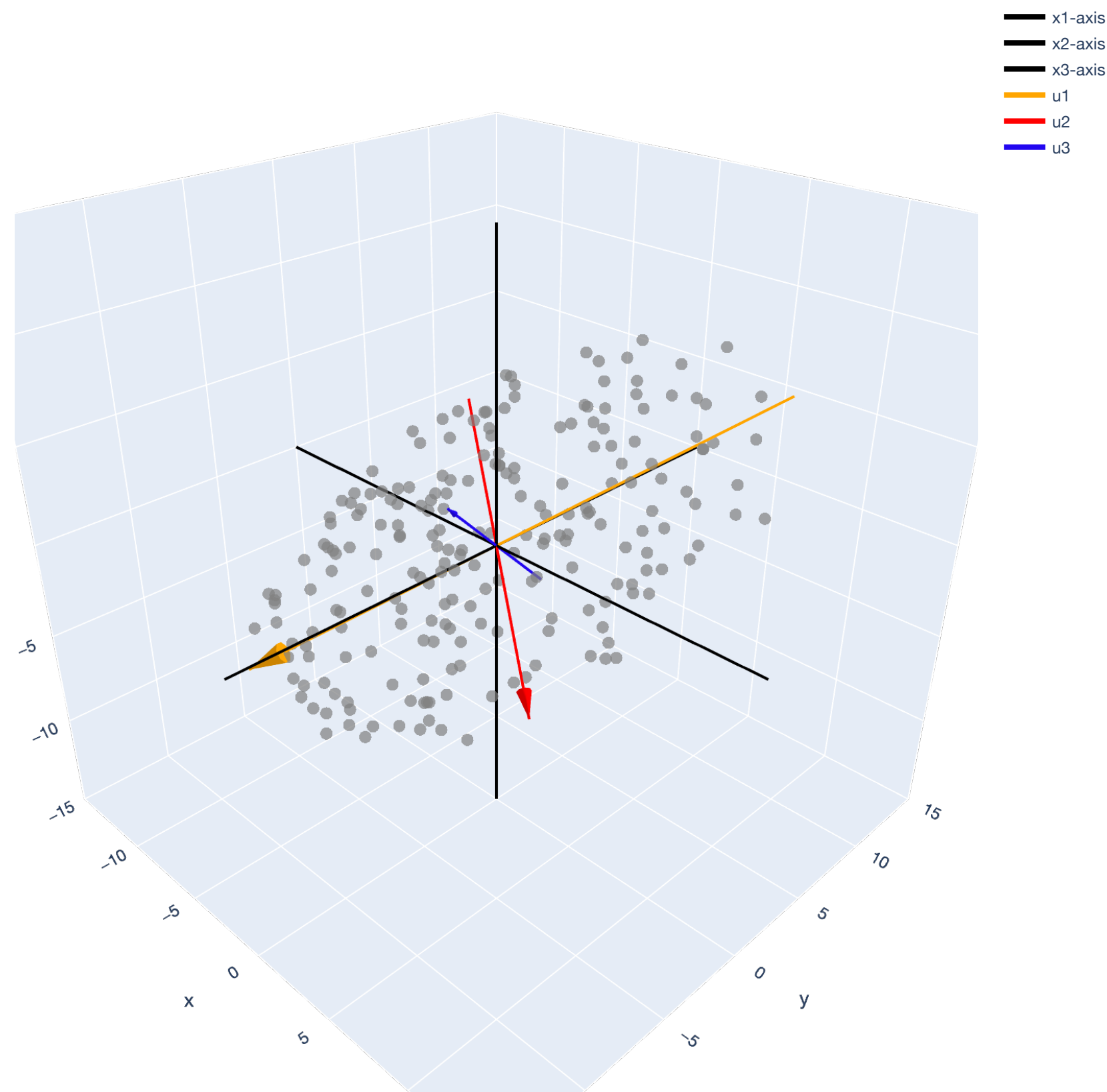
$$X = U \Sigma V^T = \begin{matrix} u_1 & u_2 & u_3 \\ \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 10 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \end{matrix}$$

The diagram shows the decomposition of matrix X into three matrices: U, Σ, and V^T. The columns of U are labeled u₁, u₂, and u₃. The diagonal elements of Σ are 10, 5, and 1. The rows of V^T are labeled 1, 0, 0, 0, 1, 0, 0, 0, 1. The matrices U, Σ, and V^T are shown as products of matrices with their respective dimensions and elements.

Singular Value Decomposition (SVD)

Example in \mathbb{R}^3

$$X = \begin{bmatrix} | & & | \\ x_1 & \dots & x_{2 \times 2} \\ | & & | \end{bmatrix} \quad \begin{matrix} n=3 \\ d=2 \end{matrix}$$



Singular Value Decomposition (SVD)

Definition of the Compact SVD

n.p. linear. svd (full = False)

Consider any matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ with rank $r \leq \min\{n, d\}$. By the compact singular value decomposition (SVD), there exist matrices \mathbf{U} , $\mathbf{\Sigma}$, \mathbf{V} such that

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

diagonal .

The columns of $\mathbf{U} \in \mathbb{R}^{n \times r}$ are the left singular vectors and $\mathbf{U}^T \mathbf{U} = \mathbf{I}$. They form an orthonormal basis for $\text{span}(\text{col}(\mathbf{X}))$, the columnspace of \mathbf{X} .

The columns of $\mathbf{V} \in \mathbb{R}^{r \times d}$ are the right singular vectors and $\mathbf{V}^T \mathbf{V} = \mathbf{I}$. They form an orthonormal basis for $\text{span}(\text{row}(\mathbf{X}))$, the row space of \mathbf{X} .

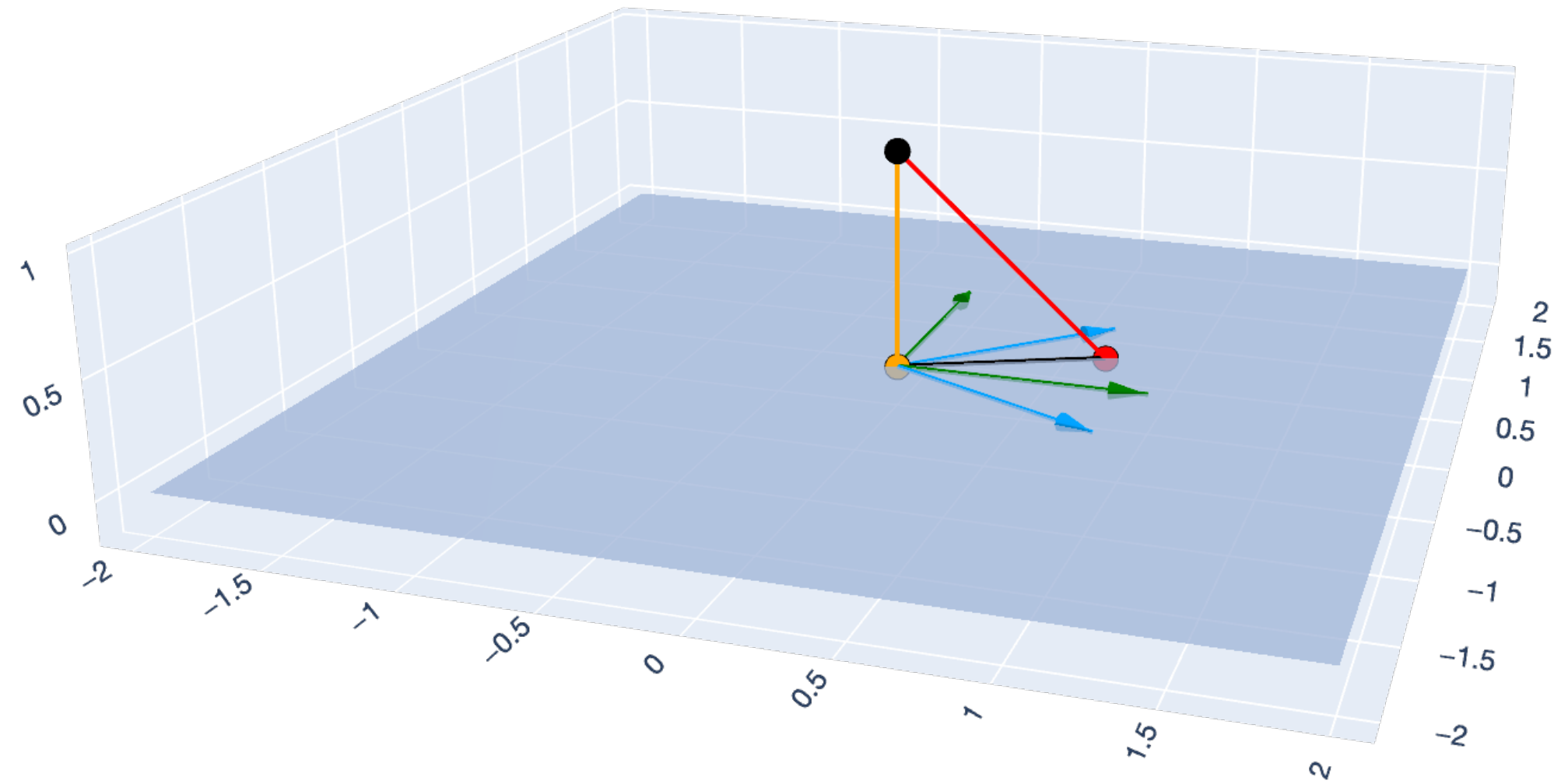
$\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$ is a diagonal matrix with singular values $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$ on the diagonal.

Best.

How to find a ~~good~~ orthogonal basis?

numpy

np.linalg.svd(x)



x1 x2 u1 u2 $y - \hat{y}$ $\tilde{y} - \hat{y}$ $\tilde{y} - y$ y \hat{y} \tilde{y}

Least Squares

OLS with Orthogonal Basis

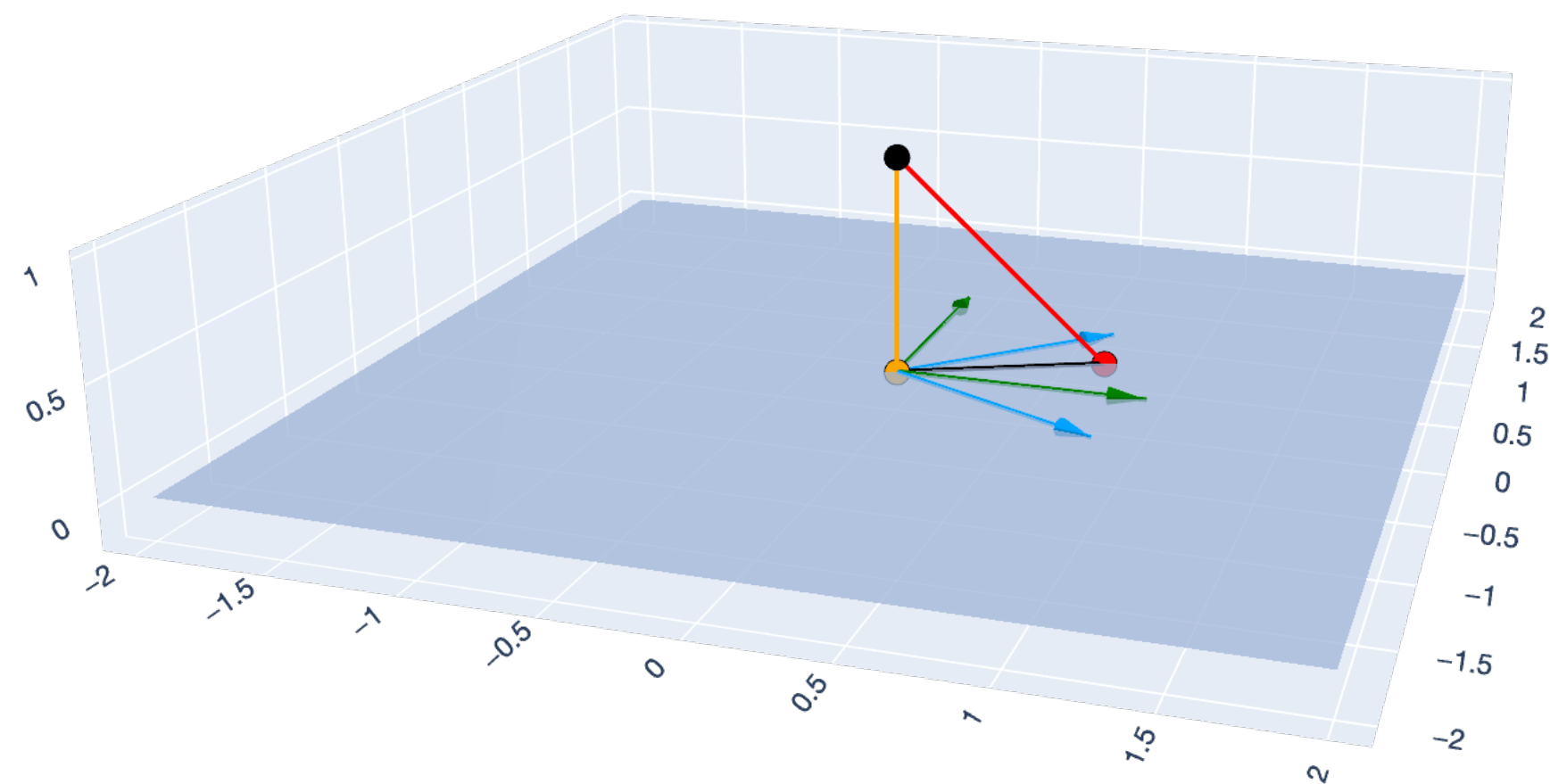
$$\hat{\mathbf{w}} = \underline{(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}$$



$$\hat{\mathbf{w}}_{onb} = \underline{\mathbf{U}^\top \mathbf{y}}$$

$$\hat{\mathbf{y}} = \Pi_{\mathcal{X}}(\mathbf{y}) = \underline{\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}}$$

$$\hat{\mathbf{y}} = \Pi_{\mathcal{X}}(\mathbf{y}) = \underline{\mathbf{U}\mathbf{U}^\top \mathbf{y}}$$



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

Least Squares

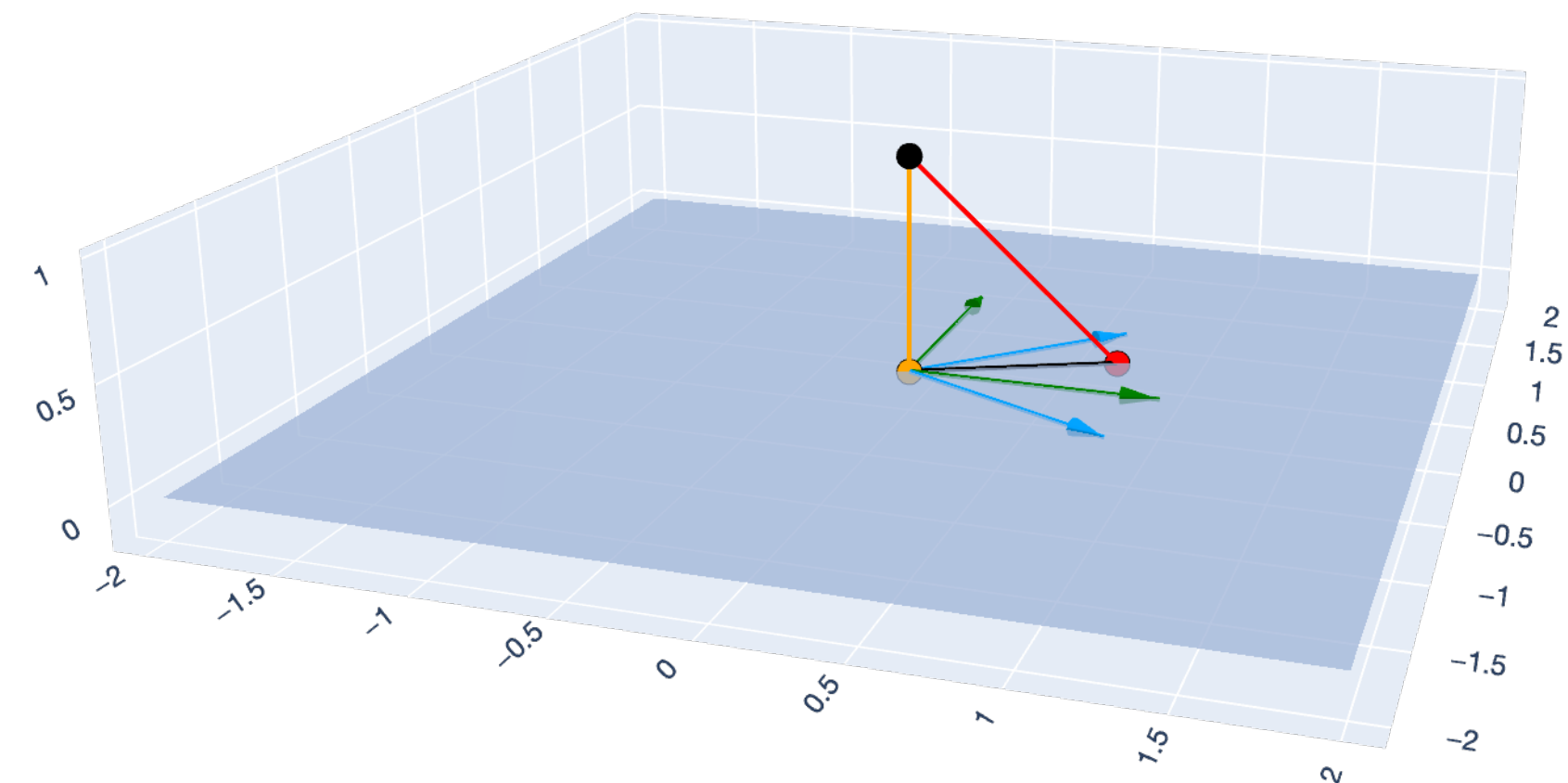
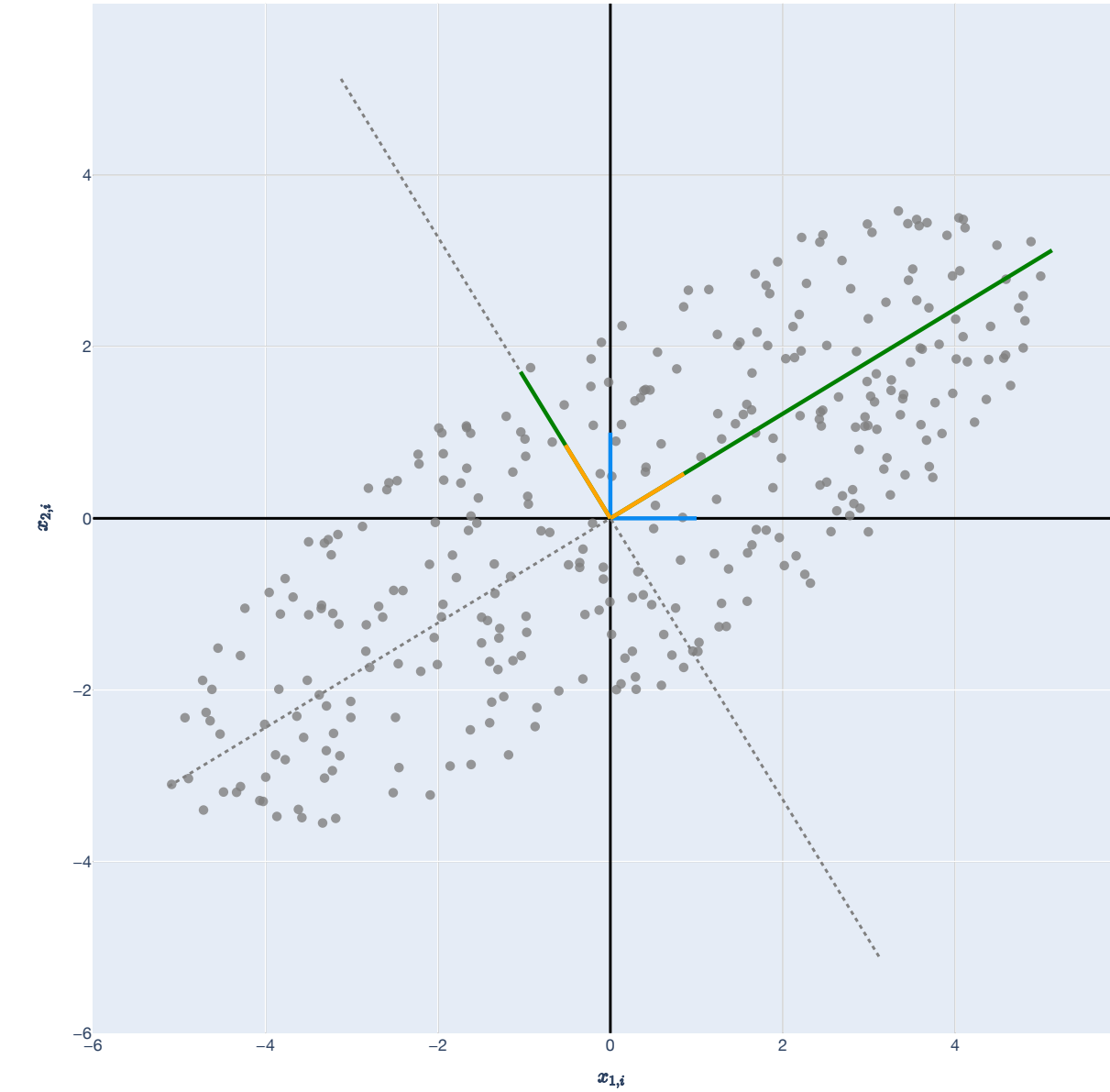
OLS with Orthogonal Basis

Prop (OLS using the ONB from Compact SVD).

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and let $\mathcal{X} = \text{span}(\text{col}(\mathbf{X}))$ be a subspace, with dimension $\dim(\mathcal{X}) = \text{rank}(\mathbf{X}) = r$.

Then, if the compact SVD of $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ then the columns $\mathbf{u}_1, \dots, \mathbf{u}_r \in \mathbb{R}^n$ of \mathbf{U} are an ONB for \mathcal{X} and, hence, for any $\mathbf{y} \in \mathbb{R}^n$, the projection of \mathbf{y} onto \mathcal{X} is given by:

$$\Pi_{\mathcal{X}}(\mathbf{y}) = \mathbf{U}\mathbf{U}^T\mathbf{y}.$$



— x_1 — x_2 — u_1 — u_2 — $y - \hat{y}$ — $\tilde{y} - \hat{y}$ — $\tilde{y} - y$ • y • \hat{y} • \tilde{y}

Singular Value Decomposition

Application: Low-rank Approximation

Rank- k Approximation

Idea

In many applications, it is useful to *approximate* a matrix. The *rank* of a matrix represents how many linearly independent columns (or rows) make up a matrix (i.e. how much “novel information” the matrix contains).

We might approximate a matrix \mathbf{X} with $r = \text{rank}(\mathbf{X})$ by asking:

What's the closest rank- k matrix (with $k \ll r$) to \mathbf{X} ?

Rank- k Approximation

Idea

In many applications, it is useful to *approximate* a matrix. The *rank* of a matrix represents how many linearly independent columns (or rows) make up a matrix (i.e. how much “novel information” the matrix contains).

We might approximate a matrix \mathbf{X} with $r = \text{rank}(\mathbf{X})$ by asking:

What's the closest rank- k matrix (with $k \ll r$) to \mathbf{X} ?

One notion of “close” for matrices is the [Frobenius norm](#):

$$\|\mathbf{X}\|_F := \sqrt{\sum_{i=1}^n \sum_{j=1}^d X_{ij}^2}.$$

Rank- k Approximation

Statement

Eckart - Young

Theorem (Rank- k Approximation). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$. Let $\hat{\mathbf{X}}_k \in \mathbb{R}^{n \times d}$ be the rank- k approximation of \mathbf{X} in Frobenius norm:

$$\hat{\mathbf{X}}_k = \arg \min_{\hat{\mathbf{X}} \in \mathbb{R}^{n \times d}} \|\mathbf{X} - \hat{\mathbf{X}}\|_F,$$

such that $\text{rank}(\hat{\mathbf{X}}) = k$.

"closest"

$$\mathbf{V} = \begin{bmatrix} | & & & | \\ \mathbf{v}_1 & \dots & \mathbf{v}_r & \\ | & & & | \end{bmatrix} \quad \mathbf{V}_k \in \mathbb{R}^{n \times k}$$

Then, if $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ is the compact SVD of \mathbf{X} with $\mathbf{U}_k \in \mathbb{R}^{n \times k}$, $\mathbf{\Sigma} \in \mathbb{R}^{k \times k}$, and $\mathbf{V} \in \mathbb{R}^{d \times k}$ are truncated matrices of \mathbf{U} , $\mathbf{\Sigma}$, and \mathbf{V} , respectively, then

$$\hat{\mathbf{X}}_k = \mathbf{U}_k \mathbf{\Sigma}_k \mathbf{V}_k^T \text{ and } \|\mathbf{X} - \hat{\mathbf{X}}_k\|_F^2 = \sum_{i=k+1}^r \sigma_i^2.$$

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_1 & & \\ & \dots & \\ & & \sigma_r \\ & & & \dots \end{bmatrix}$$

Rank- k Approximation

Outer Product Interpretation

NEXT TIME

The (compact) SVD of a matrix can also be written as a sum of rank-1 matrices.

$$\mathbf{X} = \underbrace{\sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top}_{n \times d} + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^\top + \dots + \sigma_r \mathbf{u}_r \mathbf{v}_r^\top.$$

In this way, the rank- k approximation $\hat{\mathbf{X}}_k$ can be written as truncating this sum:

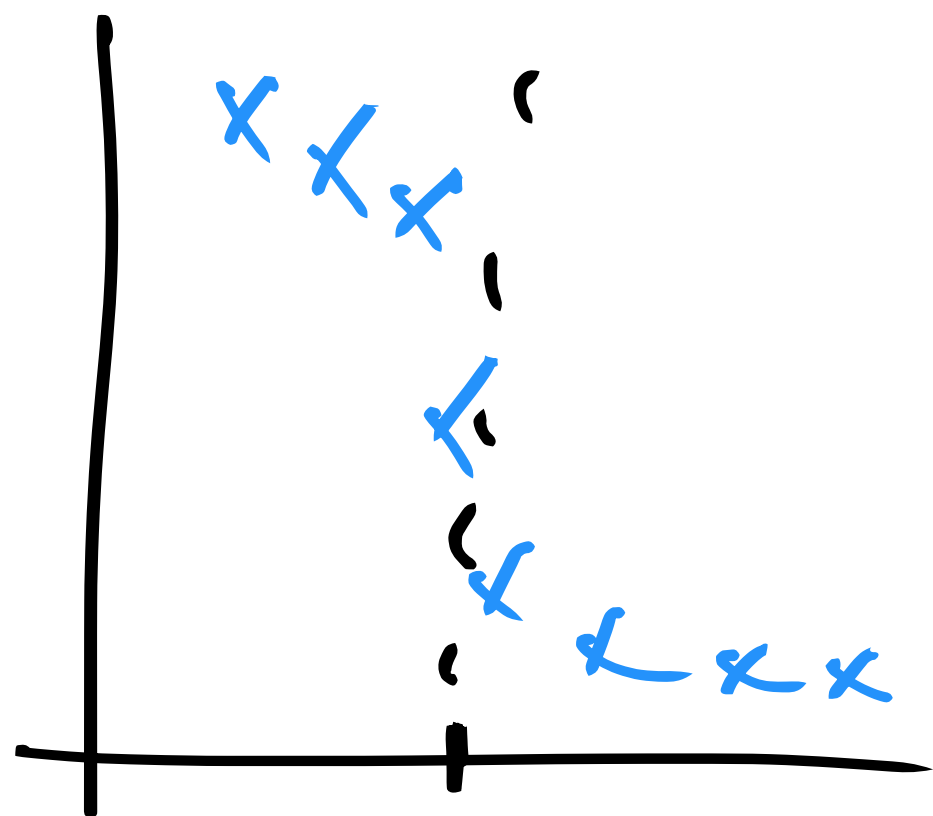
$$\hat{\mathbf{X}}_k = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^\top + \dots + \sigma_k \mathbf{u}_k \mathbf{v}_k^\top.$$

Rank-k Approximation

Example

$k=2$; $X = \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 100 & 0 \\ 0 & 90 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}}_{n \times d}$

Consider the 4×4 matrix



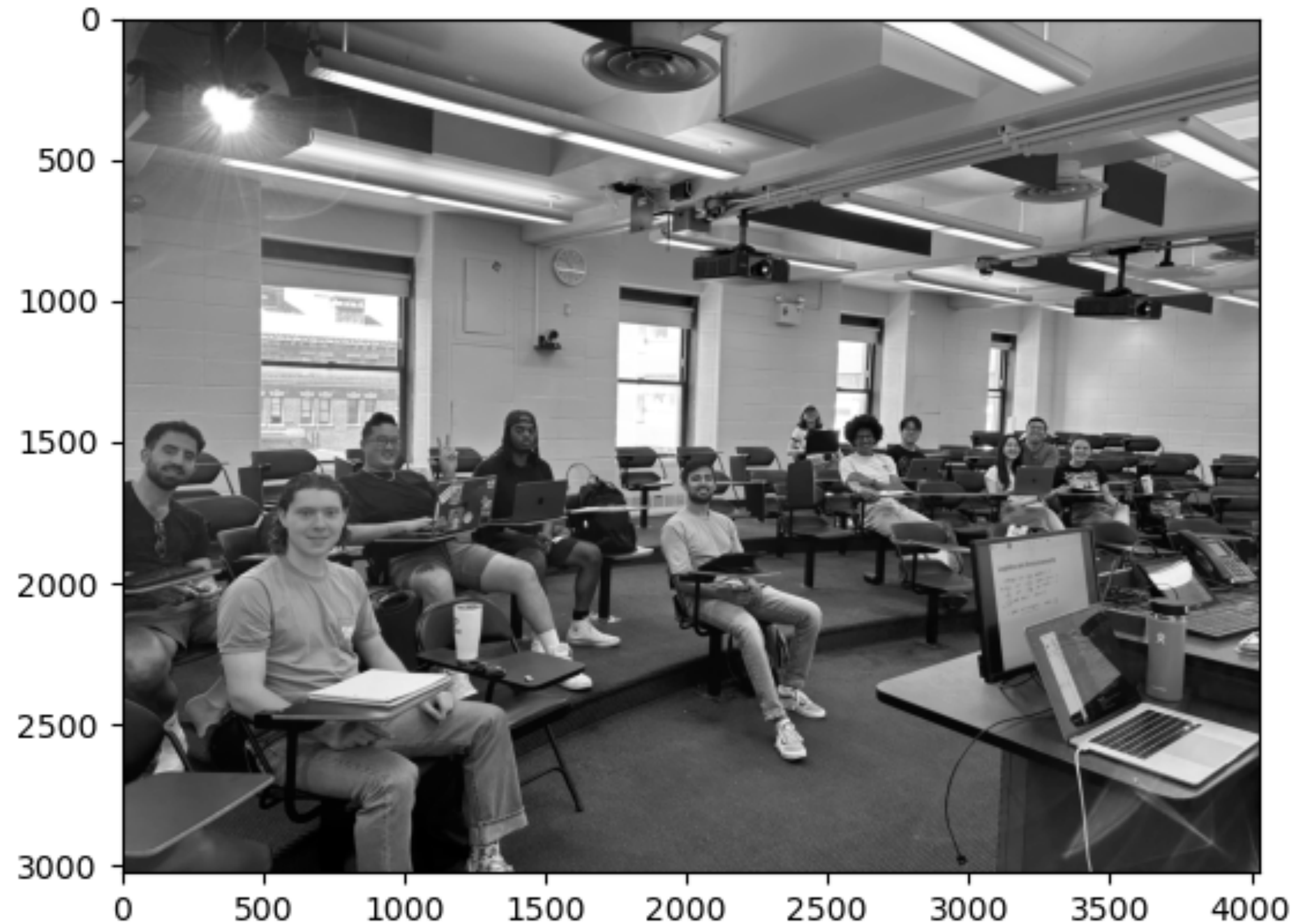
$$X = \begin{bmatrix} 100 & 0 & 0 & 0 \\ 0 & 90 & 0 & 0 \\ 0 & 0 & 8 & 0 \\ 0 & 0 & 0 & 2 \end{bmatrix}$$

$$X = \underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}}_V \underbrace{\begin{bmatrix} 100 & 0 & 0 & 0 \\ 0 & 90 & 0 & 0 \\ 0 & 0 & 8 & 0 \\ 0 & 0 & 0 & 2 \end{bmatrix}}_S \underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}}_U$$

Rank- k Approximation

Application in Image Processing

$$\begin{aligned} n &= 3000 \\ d &= 4000 \end{aligned}$$

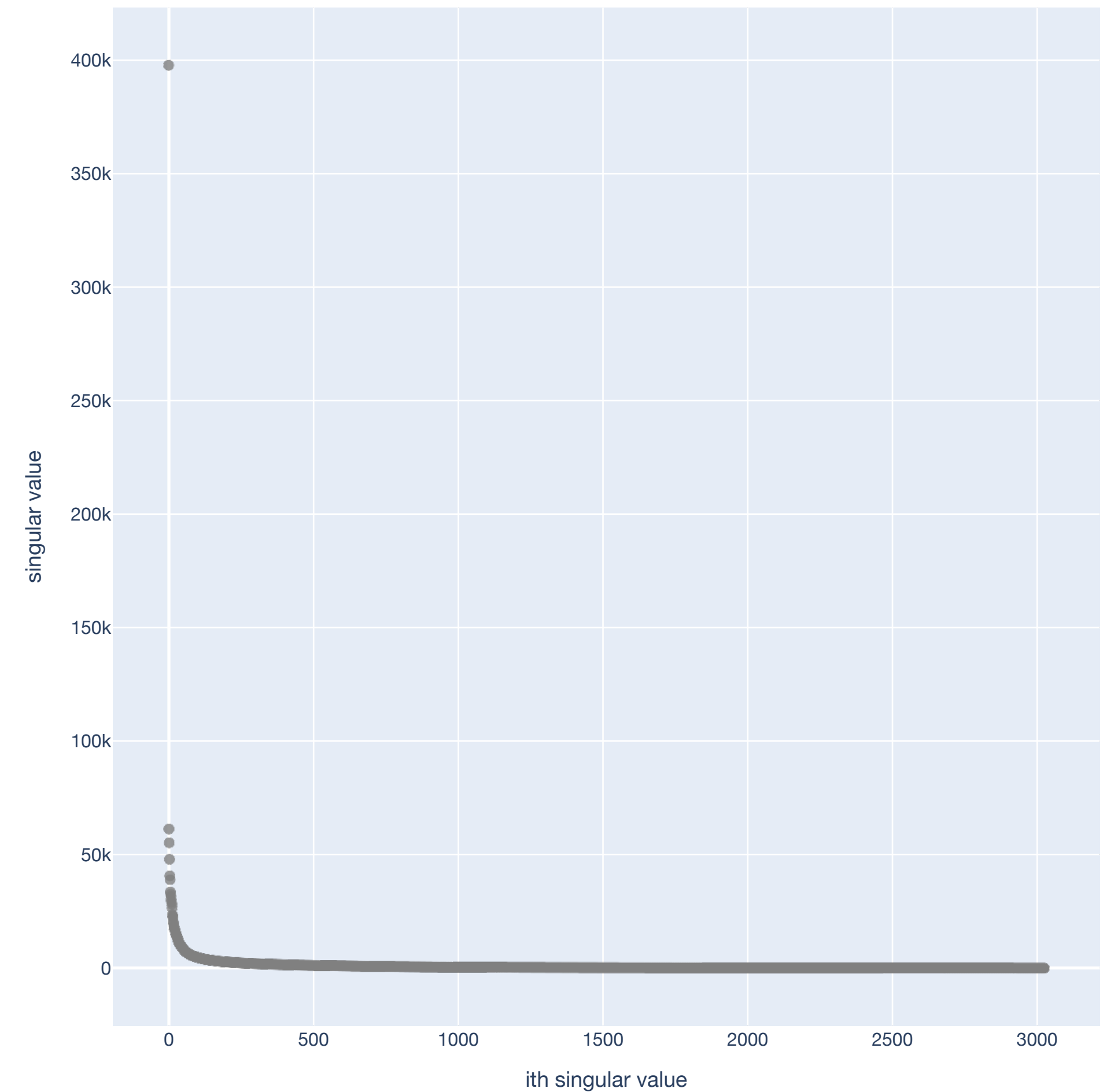
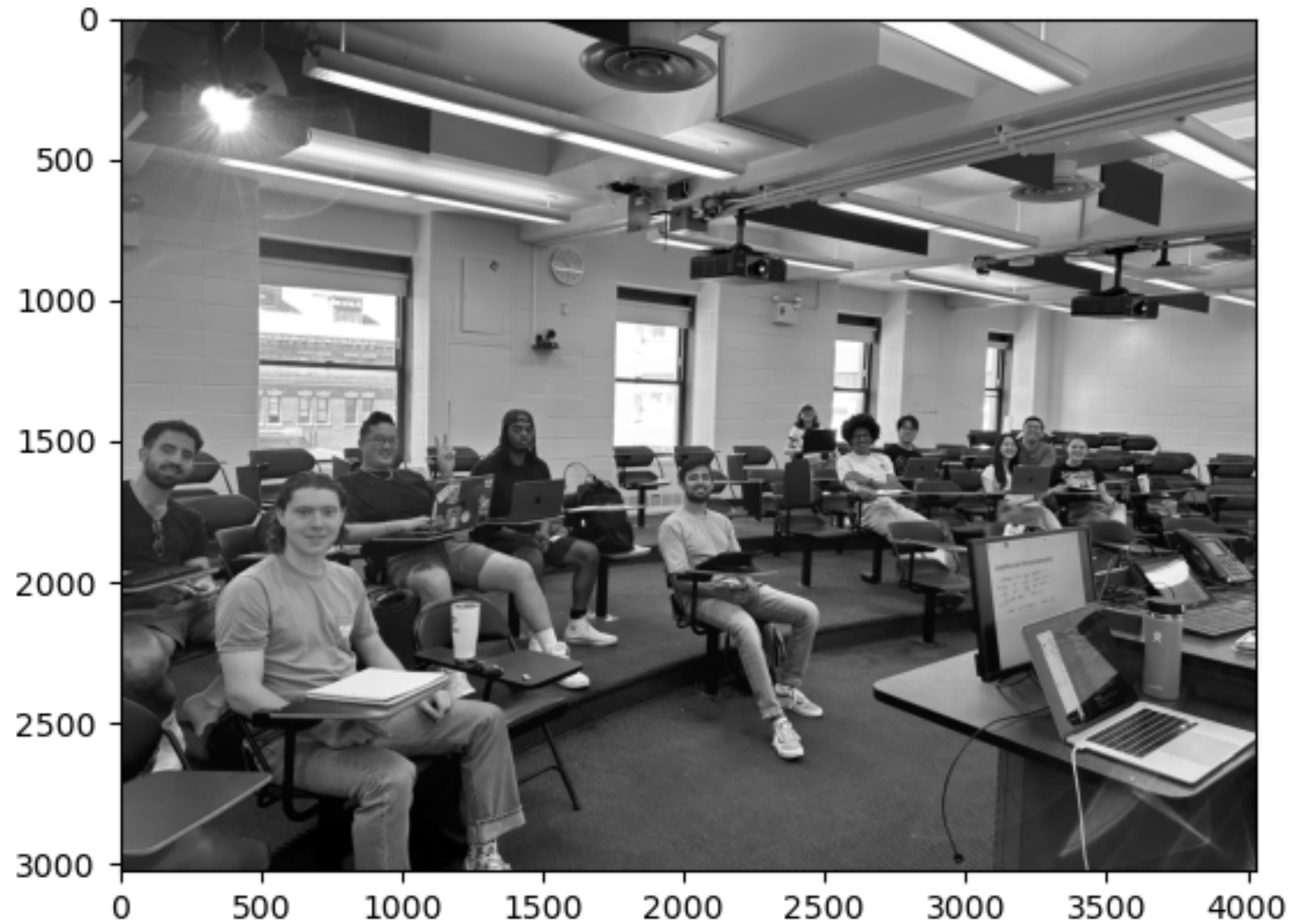


```
print(X)
print("Shape: {}".format(X.shape))
✓ 0.0s

[[ 78  78  78 ... 124 122 129]
 [ 82  81  79 ... 124 121 126]
 [ 81  80  78 ... 120 123 127]
 ...
 [ 40  42  40 ... 116  99 118]
 [ 40  41  40 ... 111 114 119]
 [ 41  39  40 ... 120 122  96]]
Shape: (3024, 4032)
```

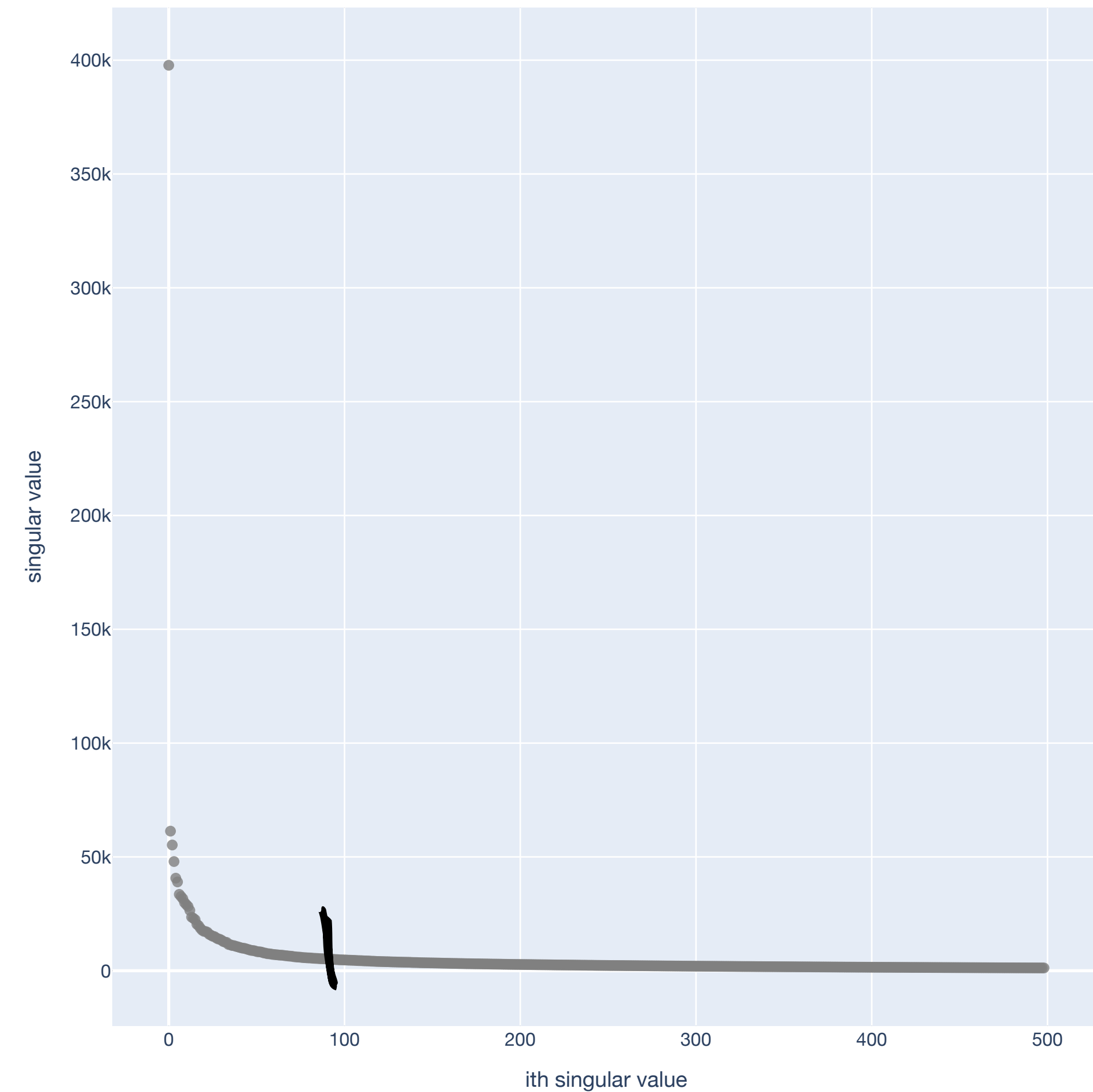
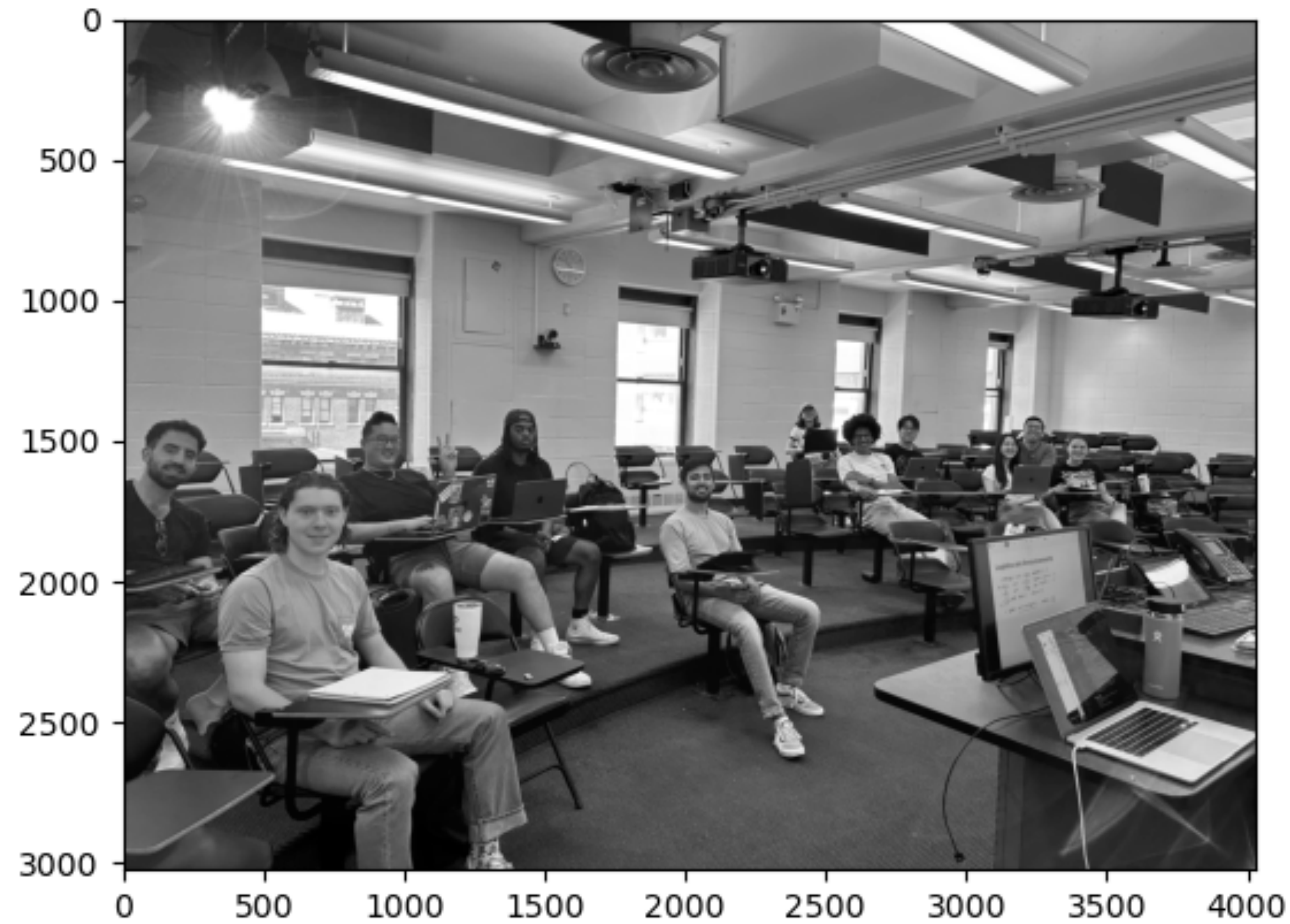

Rank- k Approximation

Application in Image Processing



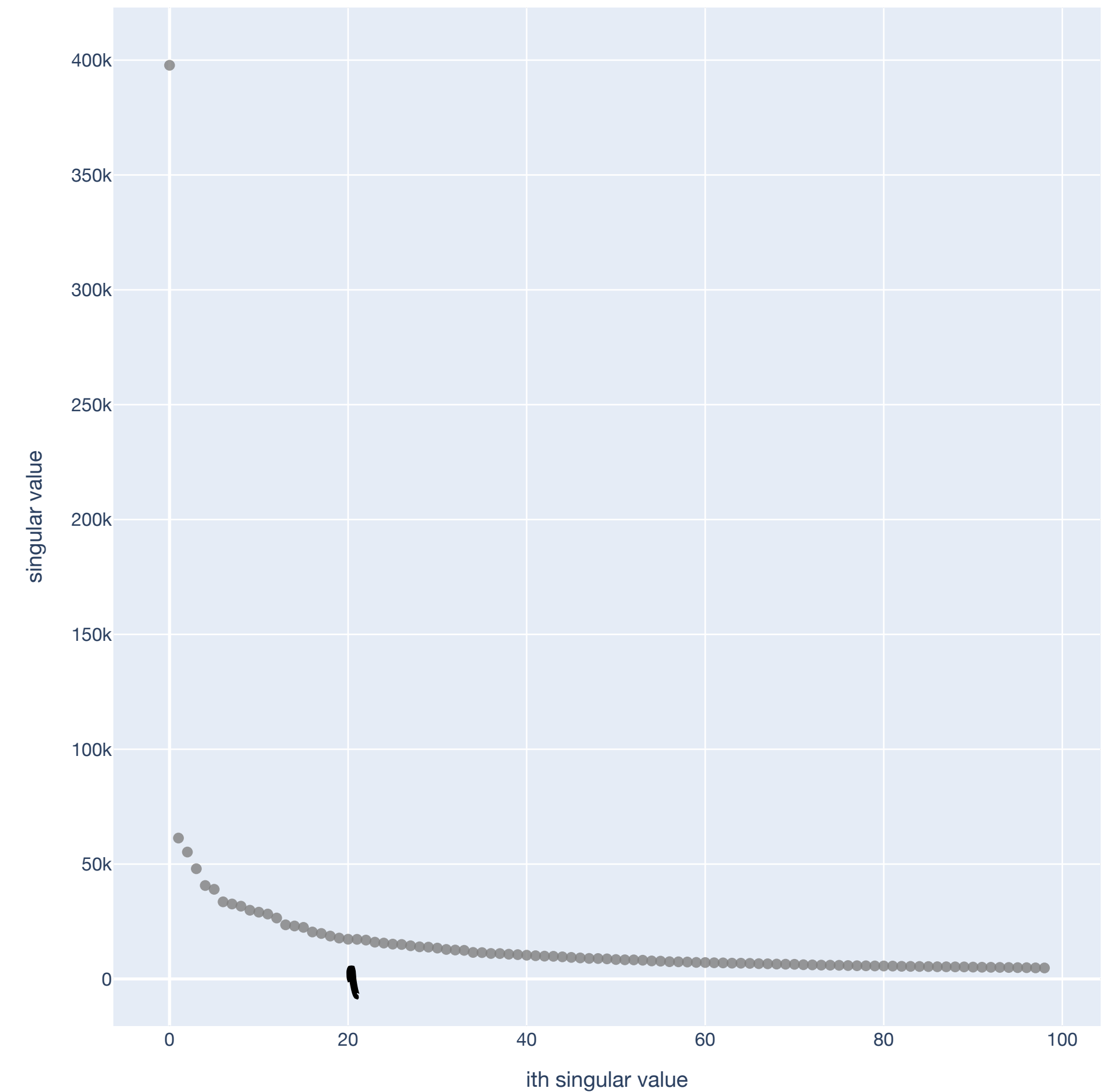
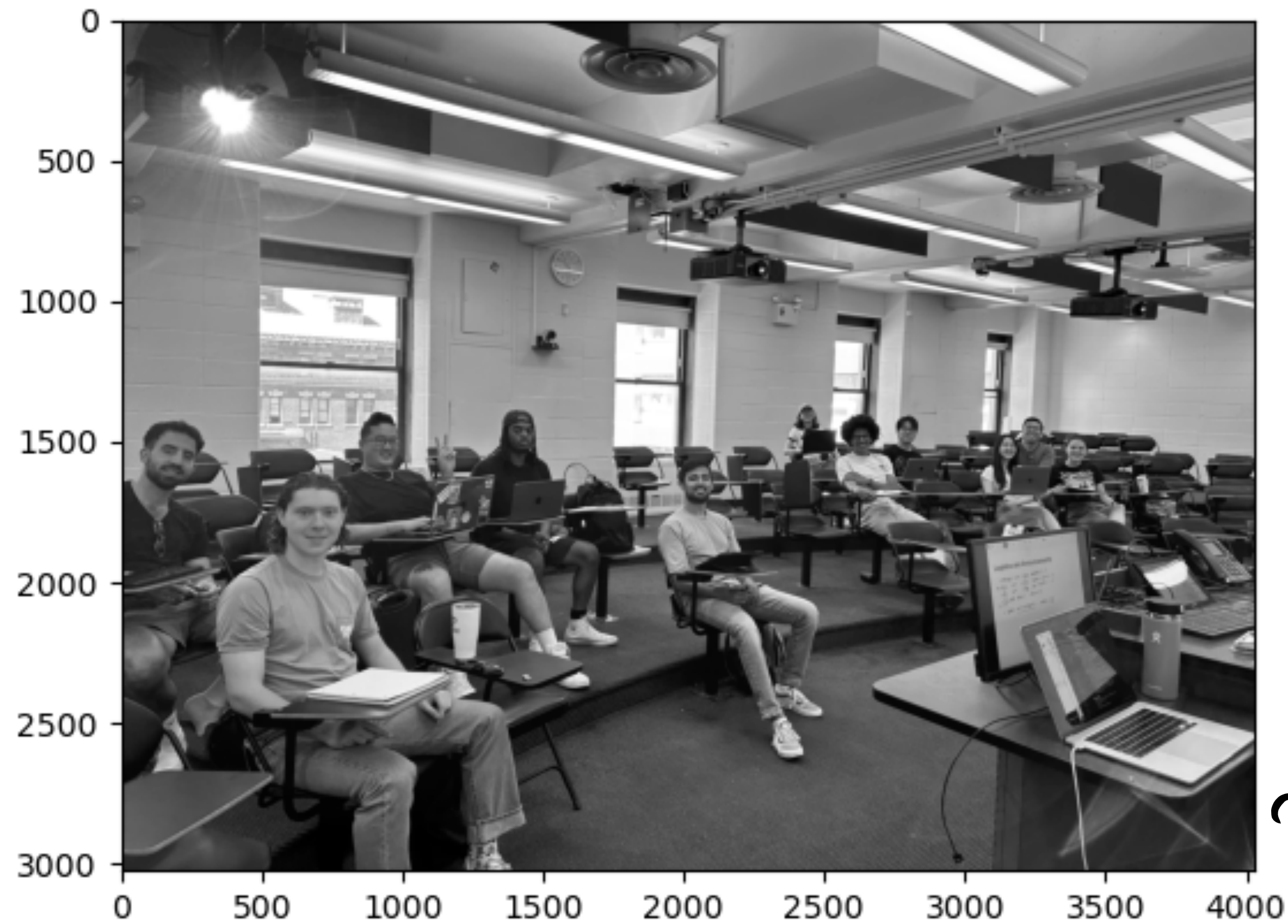
Rank- k Approximation

Application in Image Processing ($k = 500$)



Rank- k Approximation

Application in Image Processing ($k = 100$)



Rank- k Approximation

Application in Image Processing ($k = 20$)



$$X \in \mathbb{R}^{3000 \times 4000}$$

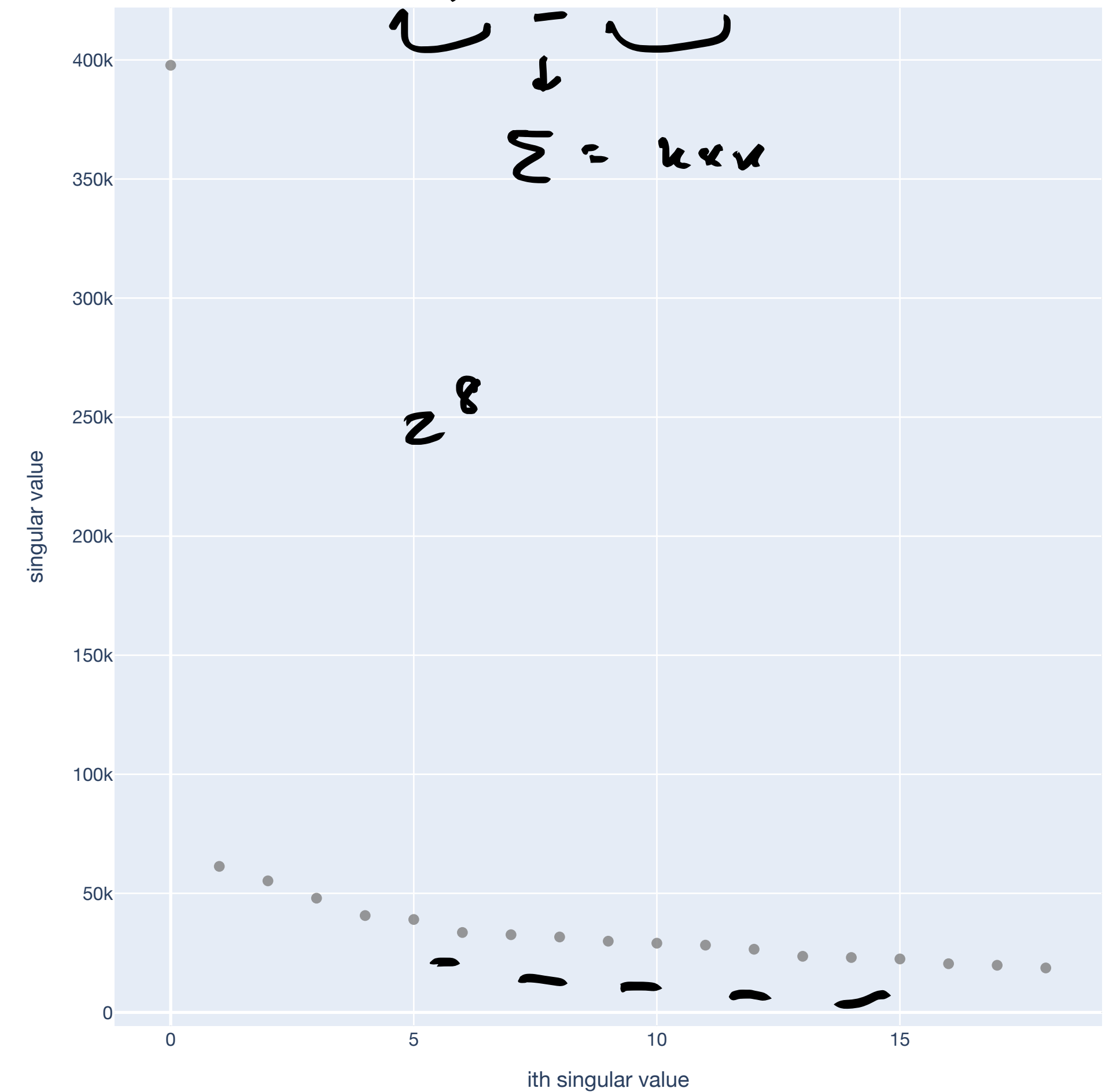
np.linalg.svd(X)

$$X = U \Sigma V^T \quad \kappa \ll \ll r$$

$n \times r$ $r \times r$ $r \times d$

$$\Sigma = k \times k$$

2^8



Least Squares

SVD and the Pseudoinverse

Regression Setup

$$\begin{array}{c} X \mathbf{w} = \mathbf{y} \\ \begin{array}{ccc} d \times d & d & d \end{array} \end{array} \quad \begin{array}{l} \textcircled{1} X \in \mathbb{R}^{d \times d} \\ \textcircled{2} \text{rank}(X) = d \end{array}$$
$$\mathbf{w} = X^{-1} \mathbf{y}$$

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*.

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

Least Squares

Main Theorem

Theorem (Ordinary Least Squares). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^n$. Let $\hat{\mathbf{w}} \in \mathbb{R}^d$ be the least squares minimizer:

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

If $n \geq d$ and $\text{rank}(\mathbf{X}) = d$, then:

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

To get predictions $\hat{\mathbf{y}} \in \mathbb{R}^n$:

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

Least Squares: SVD Perspective

Plugging in the SVD

$n \times d$
↗

$$\textcircled{1} (AB)^T = B^T A^T$$
$$\textcircled{2} (A^{-1})^T = (A^T)^{-1}$$

By the full SVD, we can represent $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. How can we interpret the least squares solution now that we know the SVD?

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

Least Squares: SVD Perspective

Plugging in the SVD

$$(AB)^T = B^T A^T$$

$$\underline{(ABC)^T = C^T B^T A^T}$$

By the full SVD, we can represent $\mathbf{X} = \underline{\mathbf{U}\Sigma\mathbf{V}^T}$. How can we interpret the least squares solution now that we know the SVD?

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

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

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

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

Least Squares: SVD Perspective

Plugging in the SVD

By the full SVD, we can represent $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. How can we interpret the least squares solution now that we know the SVD?

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\underbrace{\mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T})^{-1} \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{y} && (\mathbf{X}^T = \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T) \\ &= (\mathbf{V}\mathbf{\Sigma}^T \mathbf{\Sigma}\mathbf{V}^T)^{-1} \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \mathbf{y} && (\underbrace{\mathbf{U}^T \mathbf{U} = \mathbf{I}})\end{aligned}$$

Least Squares: SVD Perspective

Plugging in the SVD

By the full SVD, we can represent $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. How can we interpret the least squares solution now that we know the SVD?

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{y} && (\mathbf{X}^T = \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T) \\ &= (\mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{y} && (\mathbf{U}^T \mathbf{U} = \mathbf{I}) \\ &= \underbrace{(\mathbf{V}^T)^{-1}}_{\mathbf{B}} \underbrace{(\mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma})^{-1}}_{\mathbf{A}} \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{y} && \underline{((\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1})}\end{aligned}$$

Least Squares: SVD Perspective

Plugging in the SVD

By the full SVD, we can represent $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. How can we interpret the least squares solution now that we know the SVD?

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{y} && (\mathbf{X}^T = \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T) \\ &= (\mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{y} && (\mathbf{U}^T \mathbf{U} = \mathbf{I}) \\ &= (\mathbf{V}^T)^{-1} (\mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma})^{-1} \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{y} && ((\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}) \\ &= \mathbf{V} (\mathbf{\Sigma}^T \mathbf{\Sigma})^{-1} \mathbf{V}^T \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{y} && (\mathbf{V}^{-1} = \mathbf{V}^T)\end{aligned}$$

$$\mathbf{V}^T \mathbf{V} = \mathbf{V} \mathbf{V}^T = \mathbf{I}$$

Least Squares: SVD Perspective

Plugging in the SVD

By the full SVD, we can represent $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. How can we interpret the least squares solution now that we know the SVD?

$$\begin{aligned}\hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{y} && (\mathbf{X}^T = \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T) \\ &= (\mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{y} && (\mathbf{U}^T \mathbf{U} = \mathbf{I}) \\ &= (\mathbf{V}^T)^{-1} (\mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma})^{-1} \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{y} && ((\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}) \\ &= \mathbf{V} (\mathbf{\Sigma}^T \mathbf{\Sigma})^{-1} \mathbf{V}^T \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{y} && (\mathbf{V}^{-1} = \mathbf{V}^T) \\ &= \mathbf{V} (\mathbf{\Sigma}^T \mathbf{\Sigma})^{-1} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{y} && (\mathbf{V}^T \mathbf{V} = \mathbf{I})\end{aligned}$$

Least Squares: SVD Perspective

$\frac{n \times d}{1}$

Plugging in the SVD

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

By the full SVD, we can represent $\mathbf{X} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$. How can we interpret the least squares solution now that we know the SVD?

$$\begin{aligned} \hat{\mathbf{w}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{y} \\ &= (\mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V}^T)^{-1} \mathbf{V} \mathbf{\Sigma} \mathbf{U}^T \mathbf{y} \\ &= (\mathbf{V}^T)^{-1} (\mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma})^{-1} \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V} (\mathbf{\Sigma}^T \mathbf{\Sigma})^{-1} \mathbf{V}^T \mathbf{V} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{y} \\ &= \mathbf{V} (\mathbf{\Sigma}^T \mathbf{\Sigma})^{-1} \mathbf{\Sigma}^T \mathbf{U}^T \mathbf{y} \end{aligned}$$

$\underbrace{\mathbf{V} (\mathbf{\Sigma}^T \mathbf{\Sigma})^{-1} \mathbf{\Sigma}^T}_{\mathbf{\Sigma}^+} \mathbf{U}^T \mathbf{y}$

$\mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^T \mathbf{y}$

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

$$(\mathbf{U}^T \mathbf{U} = \mathbf{I})$$

$$((\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1})$$

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

$$(\mathbf{V}^T \mathbf{V} = \mathbf{I})$$

Pseudoinverse

Idea

$$\Sigma \in \mathbb{R}^{n \times d}$$

$$\begin{bmatrix} \sigma_1 & & & 0 \\ & \sigma_2 & & \\ & & \ddots & \\ 0 & & & \sigma_d \end{bmatrix}$$

$$\Sigma^T \Sigma \in \mathbb{R}^{d \times d}$$

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

Therefore, we derived:

$$\hat{w} = \mathbf{V} (\Sigma^T \Sigma)^{-1} \Sigma^T \mathbf{U}^T y \quad (\text{when } n \geq d \text{ and } \text{rank}(\mathbf{X}) = d).$$

d pos. singular values $\sigma_1, \dots, \sigma_d$

Taking a closer look at the matrix $(\Sigma^T \Sigma)^{-1} \Sigma^T \in \mathbb{R}^{d \times n}$, we have:

$$\underbrace{(\Sigma^T \Sigma)^{-1}}_{d \times d} \underbrace{\Sigma^T}_{d \times n} \underbrace{\Sigma}_{n \times d} = \mathbf{I}_{d \times d}$$

$$\underbrace{(\Sigma^T \Sigma)^{-1}}_{d \times d} \underbrace{(\Sigma^T \Sigma)}_{d \times d} = \mathbf{I}$$

$$\Sigma^T \Sigma \in \mathbb{R}^{d \times d}$$

In this way, $(\Sigma^T \Sigma)^{-1} \Sigma^T$ acts "like an inverse" to Σ , though Σ may not be square.

Pseudoinverse

Definition

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



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.

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

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

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

$$\begin{aligned} \mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T &\Rightarrow \mathbf{X}^+ = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^{-1} \\ &= (\mathbf{V}^T)^{-1} \mathbf{\Sigma}^+ \mathbf{U}^{-1} \\ &= \mathbf{V}\mathbf{\Sigma}^+\mathbf{U}^T \end{aligned}$$

Note: If using the notation of the compact SVD, this is written differently.

Pseudoinverse

Main Property

Prop (Pseudoinverse as left/right inverse). For any matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ with full SVD $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ and $\text{rank}(\mathbf{A}) = \min\{n, d\}$, the pseudo inverse

$$\mathbf{A}^+ = \mathbf{V}\mathbf{\Sigma}^+\mathbf{U}^T = \mathbf{V}(\mathbf{\Sigma}^T\mathbf{\Sigma})^{-1}\mathbf{\Sigma}^T\mathbf{U}^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$
$$= \mathbf{V}(\mathbf{\Sigma}^T\mathbf{\Sigma})^{-1}\mathbf{\Sigma}^T\mathbf{\Sigma}\mathbf{V}^T$$
$$= \mathbf{V}\mathbf{V}^T = \mathbf{I}_d$$

has the following properties:

- If $n = d$, then \mathbf{A}^+ is the *inverse*: $\mathbf{A}^+ = \mathbf{A}^{-1}$ and $\mathbf{A}^+\mathbf{A} = \mathbf{A}\mathbf{A}^+ = \mathbf{I}$.
- If $n > d$, then \mathbf{A}^+ is a *left inverse*: $\mathbf{A}^+\mathbf{A} = \mathbf{I}_{d \times d}$.
- If $d > n$, then \mathbf{A}^+ is a *right inverse*: $\mathbf{A}\mathbf{A}^+ = \mathbf{I}_{n \times n}$.

Pseudoinverse

Shape of Σ^+

What does $\Sigma^+ = (\Sigma^T \Sigma)^{-1} \Sigma^T$ look like?

$\Sigma \in \mathbb{R}^{n \times d}$ is a diagonal matrix with [singular values](#) $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$, with $r \leq \min\{n, d\}$.

$$\underbrace{\Sigma = \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ 0 & 0 & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_d \end{bmatrix}}_{n=d} \text{ or } \underbrace{\Sigma = \begin{bmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ 0 & 0 & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_d \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}}_{n>d} \text{ or } \underbrace{\Sigma = \begin{bmatrix} \sigma_1 & 0 & \dots & 0 & 0 & 0 & \dots \\ 0 & \sigma_2 & \dots & 0 & 0 & 0 & \dots \\ 0 & 0 & \ddots & \vdots & \vdots & \vdots & \dots \\ 0 & 0 & \dots & \sigma_n & 0 & 0 & \dots \end{bmatrix}}_{d>n}$$

Pseudoinverse

Shape of Σ^+

What does $\Sigma^+ = (\Sigma^T \Sigma)^{-1} \Sigma^T$ look like?

$\Sigma \in \mathbb{R}^{n \times d}$ is a diagonal matrix with [singular values](#) $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$, with $r \leq \min\{n, d\}$.

$$\Sigma^+ = \underbrace{\begin{bmatrix} 1/\sigma_1 & 0 & \dots & 0 \\ 0 & 1/\sigma_2 & \dots & 0 \\ 0 & 0 & \ddots & \vdots \\ 0 & 0 & \dots & 1/\sigma_d \end{bmatrix}}_{n=d} \text{ or } \Sigma^+ = \underbrace{\begin{bmatrix} 1/\sigma_1 & 0 & \dots & 0 & 0 & 0 & \dots \\ 0 & 1/\sigma_2 & \dots & 0 & 0 & 0 & \dots \\ 0 & 0 & \ddots & \vdots & \vdots & \vdots & \dots \\ 0 & 0 & \dots & 1/\sigma_d & 0 & 0 & \dots \end{bmatrix}}_{n>d} \text{ or } \Sigma^+ = \underbrace{\begin{bmatrix} 1/\sigma_1 & 0 & \dots & 0 \\ 0 & 1/\sigma_2 & \dots & 0 \\ 0 & 0 & \ddots & \vdots \\ 0 & 0 & \dots & 1/\sigma_n \\ 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}}_{d>n}$$

Least Squares: SVD Perspective

Using the pseudoinverse

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^n$. Let $\hat{\mathbf{w}} \in \mathbb{R}^d$ be the least squares minimizer:

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

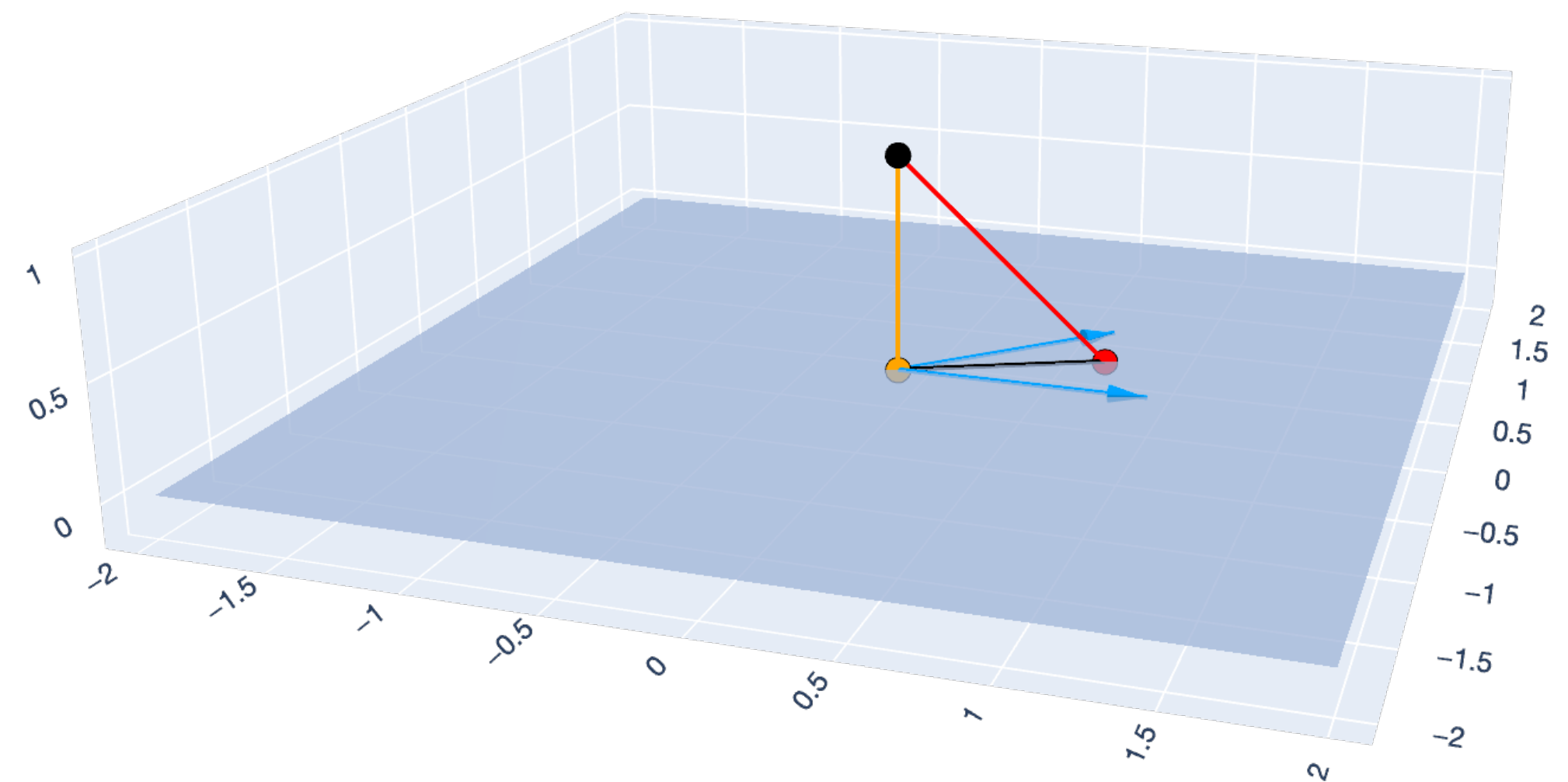
Theorem (Ordinary Least Squares).

If $n \geq d$ and $\text{rank}(\mathbf{X}) = d$, then:

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

To get predictions $\hat{\mathbf{y}} \in \mathbb{R}^n$:

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



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

Click to

Least Squares: SVD Perspective

Using the pseudoinverse

$$n \geq d \quad \text{rank}(X) = d$$

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^n$. Let $\hat{\mathbf{w}} \in \mathbb{R}^d$ be the least squares minimizer:

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

If $n = d$ and $\text{rank}(\mathbf{X}) = d$, then we are just solving the system $\mathbf{X}\mathbf{w} = \mathbf{y}$, and:

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

We solved this by the principle of least squares because, when $n > d$, we don't have an inverse. We are solving for an *approximation*:

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

$$\rightarrow \min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = 0.$$

Least Squares: SVD Perspective

Using the pseudoinverse

We solved this by the principle of least squares because, when $n > d$, we don't have an inverse. We are solving for an *approximation*:

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

We don't have an inverse — but now we have a *pseudoinverse*:

$$\underline{\mathbf{X}^+ \mathbf{X}} \underline{\mathbf{w}} \approx \underline{\mathbf{X}^+ \mathbf{y}} \implies \underbrace{\hat{\mathbf{w}}}_{\downarrow} = \underline{\mathbf{X}^+ \mathbf{y}} = \underbrace{\mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^T \mathbf{y}}_{\downarrow}.$$

Least Squares: SVD Perspective

Main Theorem (with pseudoinverse)

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\mathbf{y} \in \mathbb{R}^n$. Let $\hat{\mathbf{w}} \in \mathbb{R}^d$ be the least squares minimizer:

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

Theorem (OLS with pseudoinverse).

If $n \geq d$ and $\text{rank}(\mathbf{X}) = d$, then:

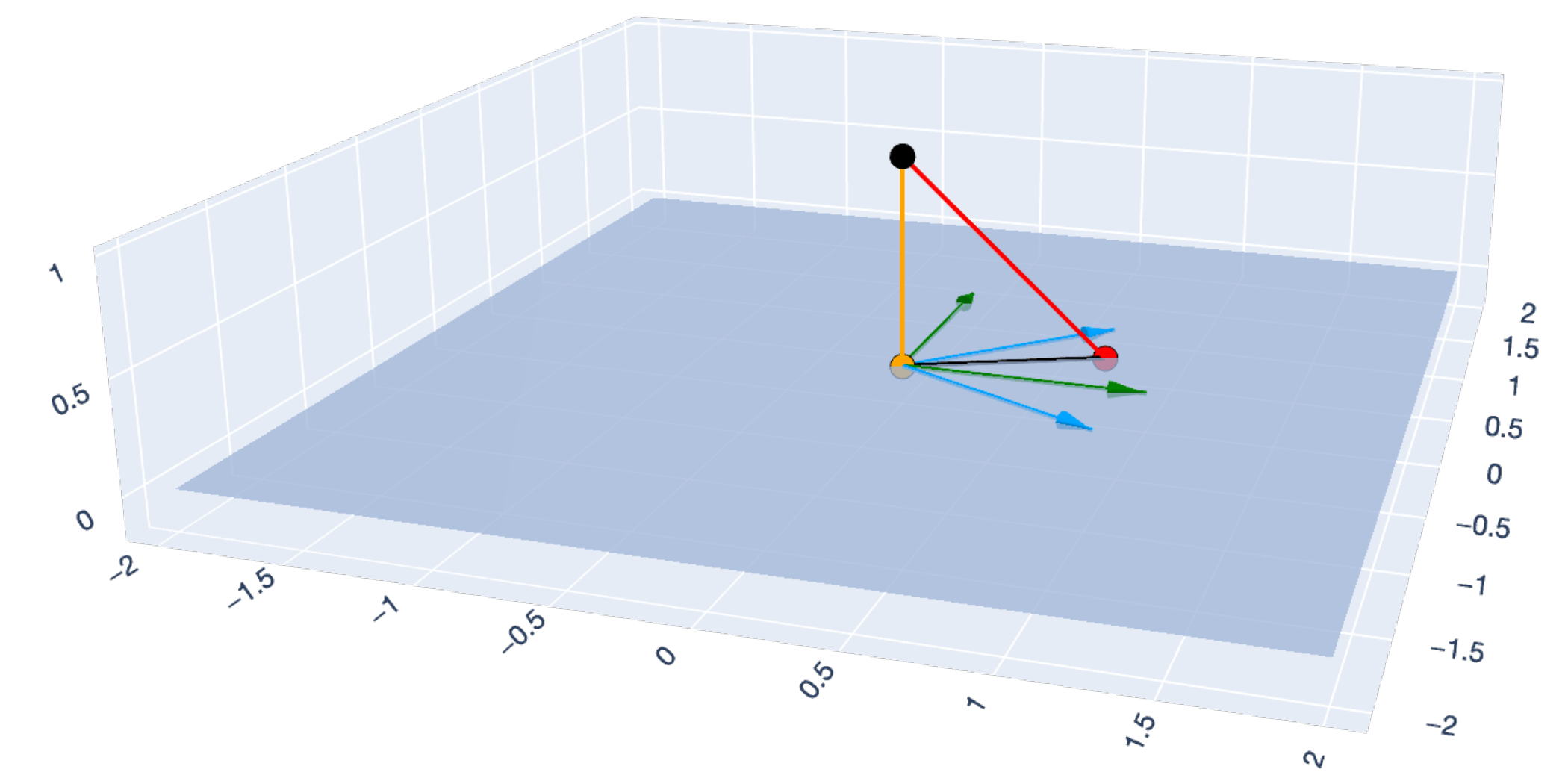
$$\hat{\mathbf{w}} = \mathbf{X}^+ \mathbf{y} = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^T \mathbf{y}.$$

To get predictions $\hat{\mathbf{y}} \in \mathbb{R}^n$:

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

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

$Xw \approx y$
 $\hat{w} \approx X^+y$



$d > n \approx$

Least Squares with $d \geq n$

Review: Systems of Linear Equations

So far, we've considered the case where $\mathbf{X} \in \mathbb{R}^{n \times d}$, $n \geq d$, and $\text{rank}(\mathbf{X}) = d$.

In general, our goal is to solve the system of linear equations:

$$\boxed{\mathbf{X}\mathbf{w} = \mathbf{y}.}$$
$$\rightarrow \begin{bmatrix} \text{---} x_1 \text{---} \\ \text{---} x_n \text{---} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$$

We know that there are three scenarios, if \mathbf{X} is full rank (i.e., $\text{rank}(\mathbf{X}) = \min\{n, d\}$)...

If $n = d$, then number of equations = number of unknowns. One unique solution: $\hat{\mathbf{w}} = \mathbf{X}^{-1}\mathbf{y}$.

If $n > d$, then number of equations > number of unknowns. One unique solution: $\hat{\mathbf{w}} = \mathbf{X}^+\mathbf{y}$.

If $d > n$, then number of unknowns > number of equations. *Infinitely many solutions!*

Systems of Linear Equations

Example: no solutions

In general, our goal is to solve the system of linear equations:

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

Consider the system:

$$\begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

Systems of Linear Equations

Example: one unique solution, $n = d$

In general, our goal is to solve the system of linear equations:

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

Consider the system:

$$\begin{bmatrix} 2 & 1 \\ 2 & -1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \end{bmatrix}$$

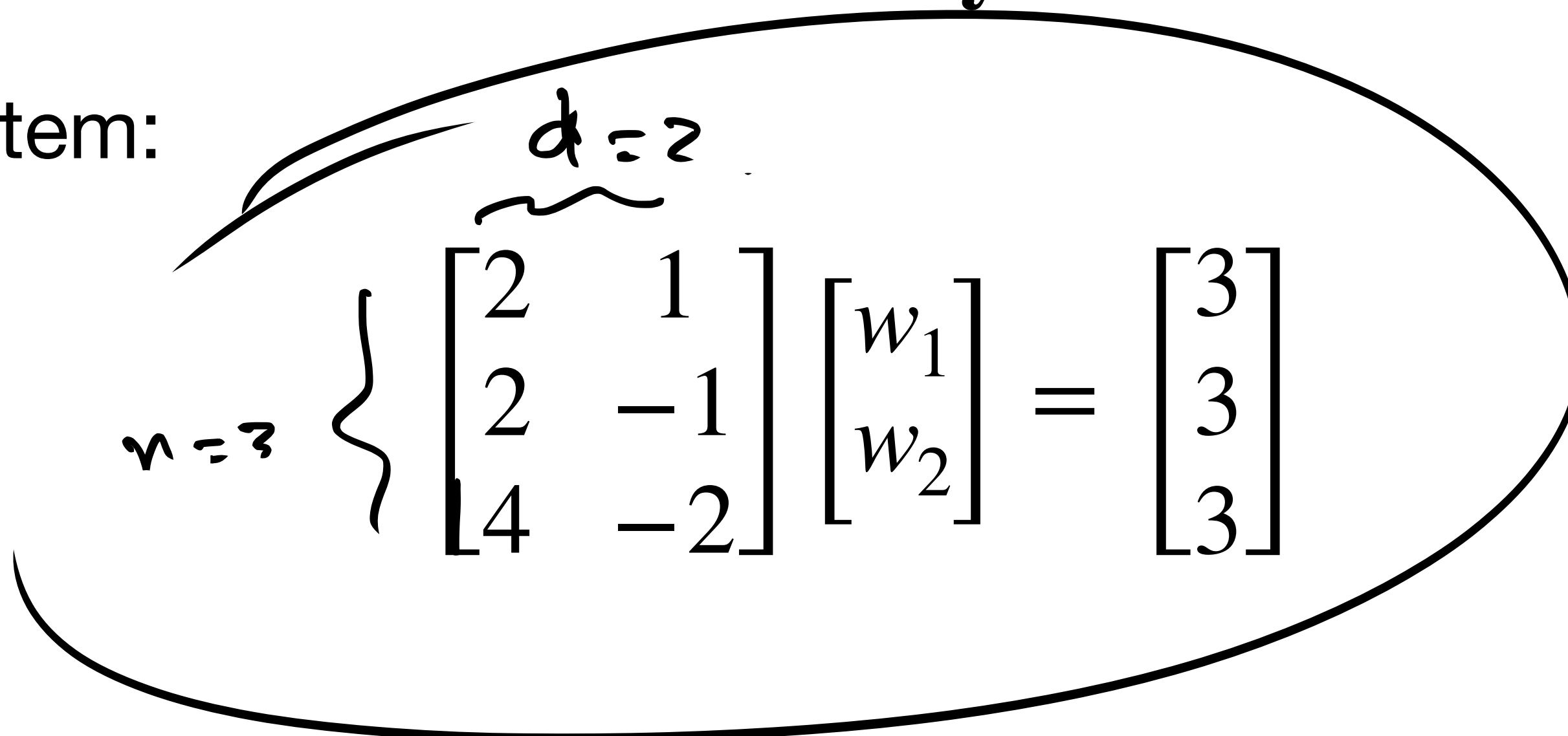
Systems of Linear Equations

Example: one unique solution, $n > d$

In general, our goal is to solve the system of linear equations:

$$Xw = y.$$

Consider the system:


$$n=3 \left\{ \begin{array}{cc} \underbrace{d=2} \\ \begin{bmatrix} 2 & 1 \\ 2 & -1 \\ 4 & -2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \\ 3 \end{bmatrix} \end{array} \right.$$

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

Systems of Linear Equations

Example: infinitely many solutions, $d > n$

In general, our goal is to solve the system of linear equations:

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

Consider the system:

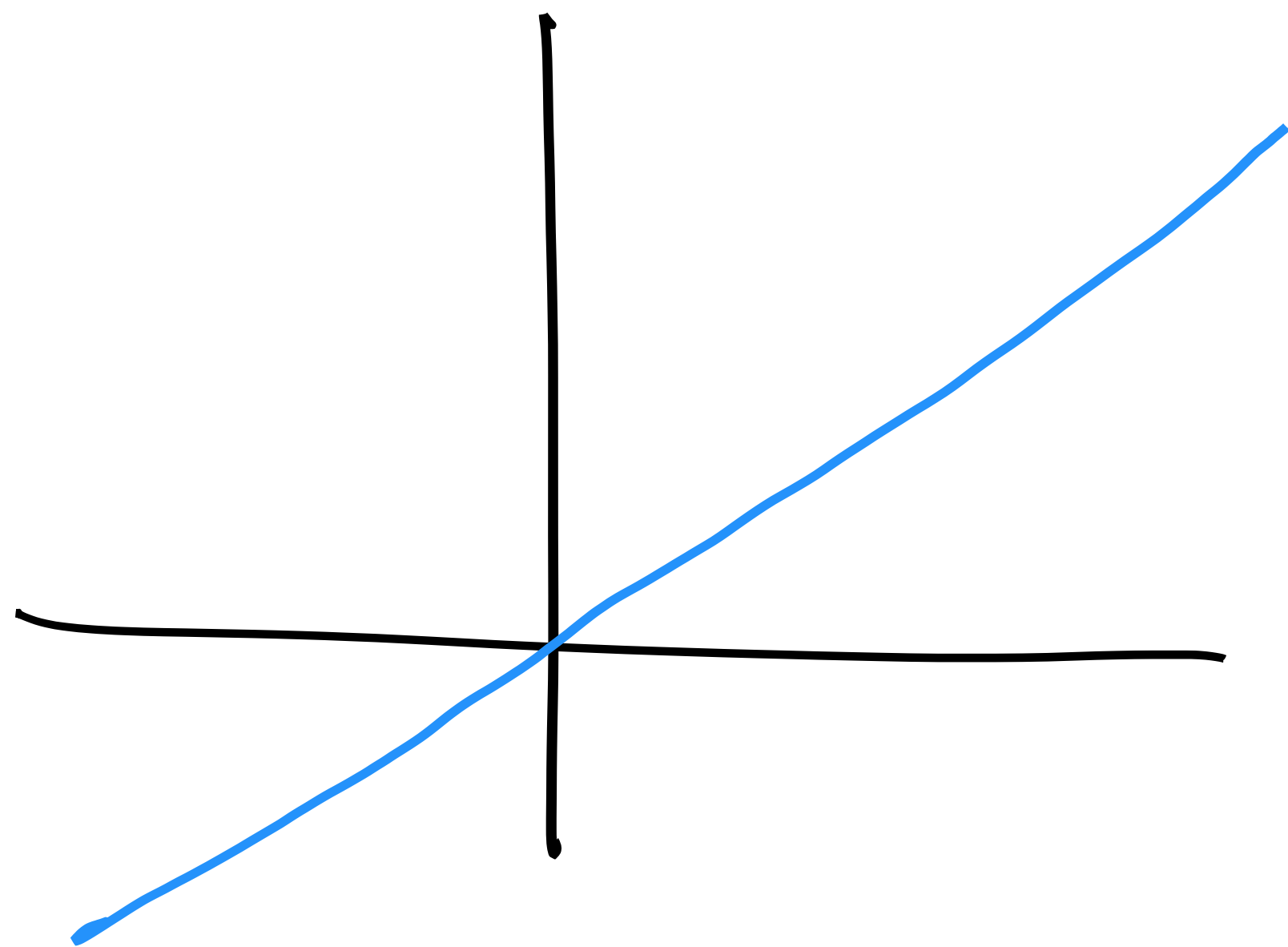
$$\begin{bmatrix} 2 & 1 & 1 \\ 2 & -1 & 0 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \end{bmatrix}$$

Least Squares with $d \geq n$

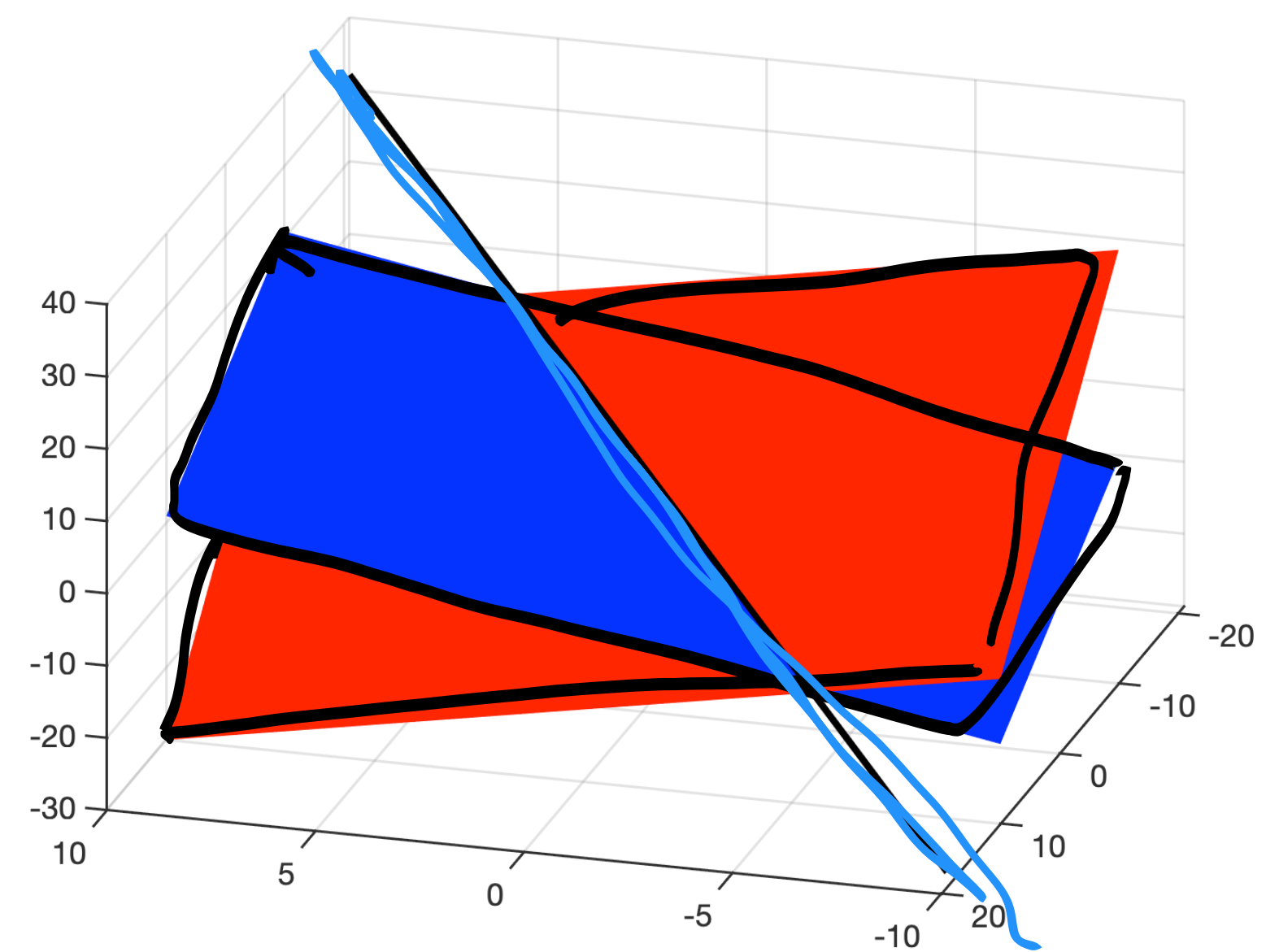
Review: Systems of Linear Equations

When the number of equations $<$ number of unknowns...

$n = 2, \mathbb{R}^2$



$n = 3, \mathbb{R}^3$



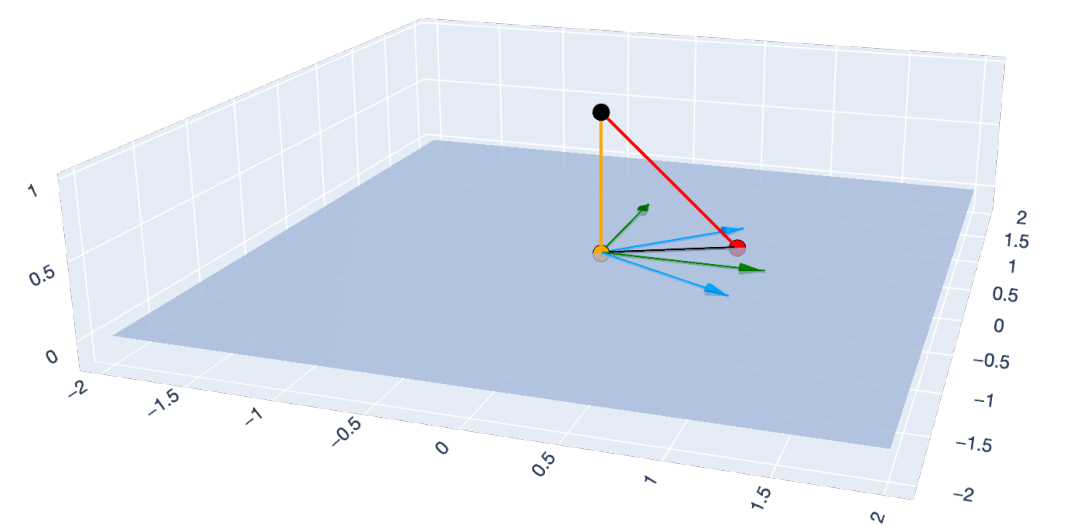
Least Squares with $d \geq n$

Problem Statement

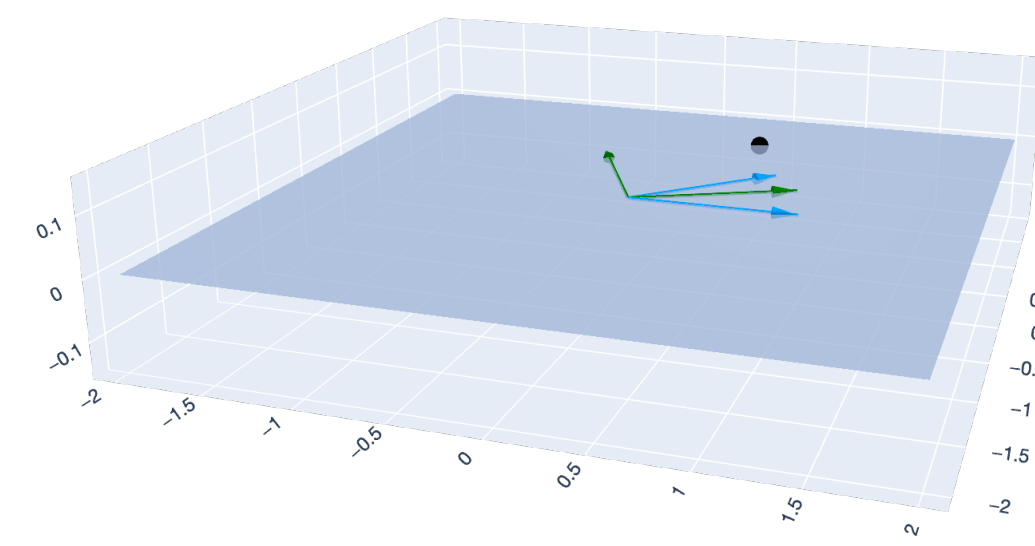
Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, let $d \geq n$, and let $\text{rank}(\mathbf{X}) = n$. We want to solve the system of linear equations:

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

Because $\text{rank}(\mathbf{X}) = n$, infinitely many *exact* solutions exist. Which to choose?



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



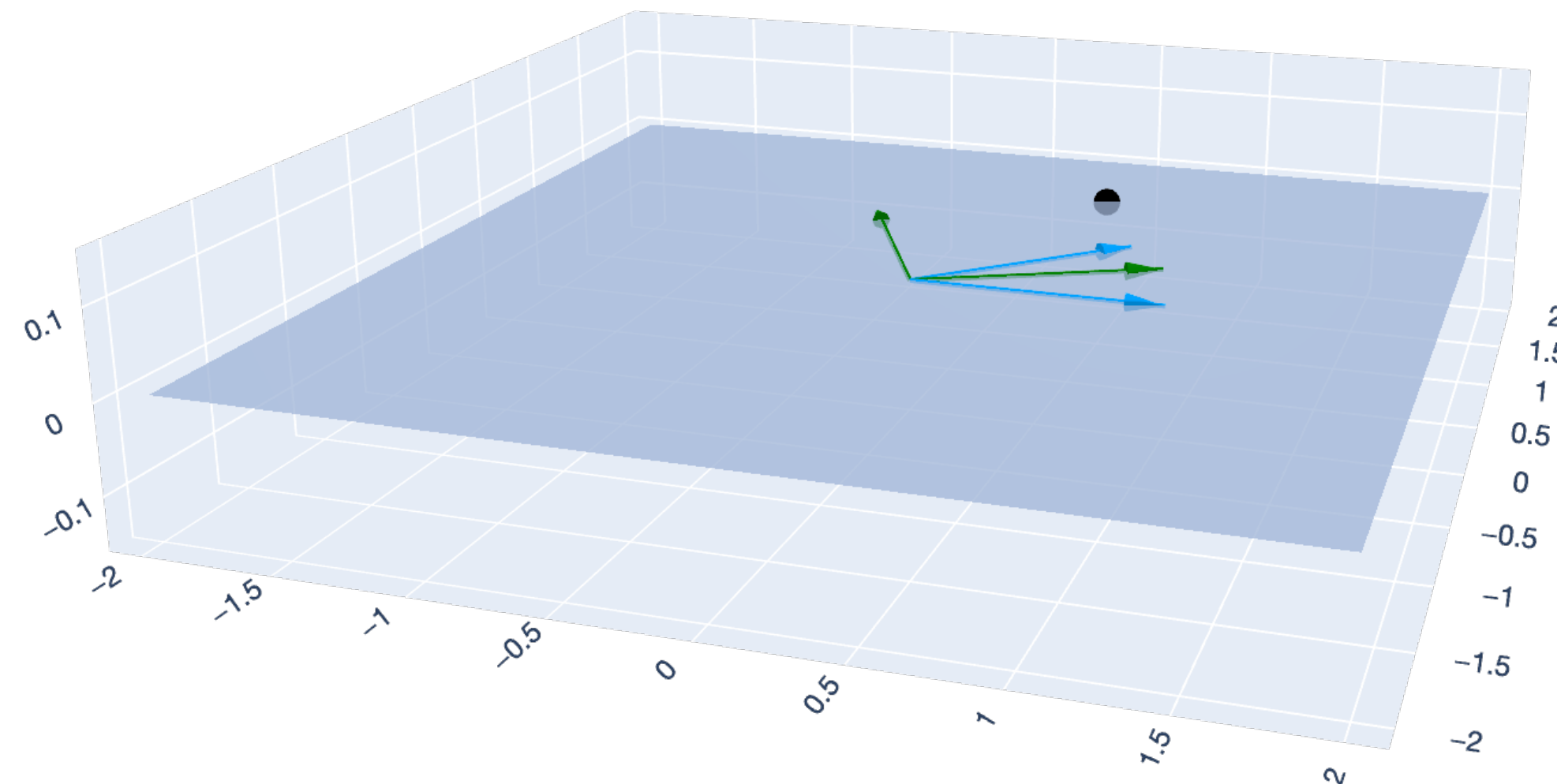
— x1 — x2 — u1 — u2 — y

Least Squares with $d \geq n$

Using the Pseudoinverse

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

There are now infinitely many $\hat{\mathbf{w}} \in \mathbb{R}^d$ such that $\mathbf{X}\hat{\mathbf{w}} = \mathbf{y}$. Which $\hat{\mathbf{w}}$ to pick?



— x_1 — x_2 — u_1 — u_2 • y

Pseudoinverse

Main Property

Prop (Pseudoinverse as left/right inverse). For any matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ with full SVD $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ and $\text{rank}(\mathbf{A}) = \min\{n, d\}$, the pseudo inverse

$$\mathbf{A}^+ = \mathbf{V}\mathbf{\Sigma}^+\mathbf{U}^\top = \mathbf{V}(\mathbf{\Sigma}^\top\mathbf{\Sigma})^{-1}\mathbf{\Sigma}^\top\mathbf{U}^\top$$

has the following properties:

- If $n = d$, then \mathbf{A}^+ is the *inverse*: $\mathbf{A}^+ = \mathbf{A}^{-1}$ and $\mathbf{A}^+\mathbf{A} = \mathbf{A}\mathbf{A}^+ = \mathbf{I}$.
- If $n > d$, then \mathbf{A}^+ is a *left inverse*: $\mathbf{A}^+\mathbf{A} = \mathbf{I}_{d \times d}$.
- If $d > n$, then \mathbf{A}^+ is a *right inverse*: $\mathbf{A}\mathbf{A}^+ = \mathbf{I}_{n \times n}$.

Least Squares with $d \geq n$

Using the Pseudoinverse

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ have the full SVD $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$.

Choose $\hat{\mathbf{w}} = \mathbf{X}^+ \mathbf{y} = \mathbf{V}\mathbf{\Sigma}^+ \mathbf{U}^T \mathbf{y}$ to use the pseudoinverse.

Least Squares with $d \geq n$

Using the Pseudoinverse

Let $\mathbf{X} \in \mathbb{R}^{n \times d}$ have the full SVD $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$.

Choose $\hat{\mathbf{w}} = \mathbf{X}^+\mathbf{y} = \mathbf{V}\mathbf{\Sigma}^+\mathbf{U}^\top\mathbf{y}$ to use the pseudoinverse.

Then, $\hat{\mathbf{w}} \in \mathbb{R}^d$ is a solution:

$$\mathbf{X}\hat{\mathbf{w}} = \mathbf{X}\mathbf{X}^+\mathbf{y} = \mathbf{I}_{n \times n}\mathbf{y} = \mathbf{y},$$

where $\mathbf{X}^+ \in \mathbb{R}^{d \times n}$ is a right inverse by the previous property.

Least Squares with $d \geq n$

Theorem: Minimum norm solution

Theorem (Minimum norm least squares solution). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, let $d \geq n$, and let $\text{rank}(\mathbf{X}) = n$. Then, $\hat{\mathbf{w}} = \mathbf{X}^+ \mathbf{y} = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^T \mathbf{y}$ is the exact solution $\mathbf{X} \hat{\mathbf{w}} = \mathbf{y}$ with smallest Euclidean norm:

$$\|\mathbf{w}\|_2^2 \geq \|\hat{\mathbf{w}}\|_2^2 \text{ for all } \mathbf{w} \in \mathbb{R}^d.$$

Least Squares with $d \geq n$

Theorem: Minimum norm solution

Theorem (Minimum norm least squares solution). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, let $d \geq n$, and let $\text{rank}(\mathbf{X}) = n$. Then, $\hat{\mathbf{w}} = \mathbf{X}^+ \mathbf{y} = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^T \mathbf{y}$ is the exact solution $\mathbf{X} \hat{\mathbf{w}} = \mathbf{y}$ with smallest Euclidean norm:

$$\|\mathbf{w}\|_2^2 \geq \|\hat{\mathbf{w}}\|_2^2 \text{ for all } \mathbf{w} \in \mathbb{R}^d.$$

Proof. Consider any arbitrary $\mathbf{w} \in \mathbb{R}^d$.

Least Squares with $d \geq n$

Theorem: Minimum norm solution

Theorem (Minimum norm least squares solution). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, let $d \geq n$, and let $\text{rank}(\mathbf{X}) = n$. Then, $\hat{\mathbf{w}} = \mathbf{X}^+ \mathbf{y} = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^\top \mathbf{y}$ is the exact solution $\mathbf{X} \hat{\mathbf{w}} = \mathbf{y}$ with smallest Euclidean norm:

$$\|\mathbf{w}\|_2^2 \geq \|\hat{\mathbf{w}}\|_2^2 \text{ for all } \mathbf{w} \in \mathbb{R}^d.$$

Proof. Consider any arbitrary $\mathbf{w} \in \mathbb{R}^d$. We can write \mathbf{w} 's Euclidean norm as:

$$\|\mathbf{w}\|^2 = \|(\mathbf{w} - \hat{\mathbf{w}}) + \hat{\mathbf{w}}\|^2 = \|\mathbf{w} - \hat{\mathbf{w}}\|^2 - 2(\mathbf{w} - \hat{\mathbf{w}})^\top \hat{\mathbf{w}} + \|\hat{\mathbf{w}}\|^2$$

Least Squares with $d \geq n$

Theorem: Minimum norm solution

Theorem (Minimum norm least squares solution). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, let $d \geq n$, and let $\text{rank}(\mathbf{X}) = n$. Then, $\hat{\mathbf{w}} = \mathbf{X}^+ \mathbf{y} = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^T \mathbf{y}$ is the exact solution $\mathbf{X} \hat{\mathbf{w}} = \mathbf{y}$ with smallest Euclidean norm:

$$\|\mathbf{w}\|_2^2 \geq \|\hat{\mathbf{w}}\|_2^2 \text{ for all } \mathbf{w} \in \mathbb{R}^d.$$

Proof. Consider any arbitrary $\mathbf{w} \in \mathbb{R}^d$. We can write \mathbf{w} 's Euclidean norm as:

$$\|\mathbf{w}\|^2 = \|(\mathbf{w} - \hat{\mathbf{w}}) + \hat{\mathbf{w}}\|^2 = \|\mathbf{w} - \hat{\mathbf{w}}\|^2 - 2(\mathbf{w} - \hat{\mathbf{w}})^T \hat{\mathbf{w}} + \|\hat{\mathbf{w}}\|^2.$$

Consider the term $(\mathbf{w} - \hat{\mathbf{w}})^T \hat{\mathbf{w}}$:

$$(\mathbf{w} - \hat{\mathbf{w}})^T \hat{\mathbf{w}} = (\mathbf{w} - \hat{\mathbf{w}})^T \underbrace{\mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}}_{\mathbf{X}^+ \text{ if } d > n}$$

Least Squares with $d \geq n$

Theorem: Minimum norm solution

Theorem (Minimum norm least squares solution). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, let $d \geq n$, and let $\text{rank}(\mathbf{X}) = n$. Then, $\hat{\mathbf{w}} = \mathbf{X}^+ \mathbf{y} = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^T \mathbf{y}$ is the exact solution $\mathbf{X} \hat{\mathbf{w}} = \mathbf{y}$ with smallest Euclidean norm:

$$\|\mathbf{w}\|_2^2 \geq \|\hat{\mathbf{w}}\|_2^2 \text{ for all } \mathbf{w} \in \mathbb{R}^d.$$

Proof. Consider any arbitrary $\mathbf{w} \in \mathbb{R}^d$. We can write \mathbf{w} 's Euclidean norm as:

$$\|\mathbf{w}\|^2 = \|(\mathbf{w} - \hat{\mathbf{w}}) + \hat{\mathbf{w}}\|^2 = \|\mathbf{w} - \hat{\mathbf{w}}\|^2 - 2(\mathbf{w} - \hat{\mathbf{w}})^T \hat{\mathbf{w}} + \|\hat{\mathbf{w}}\|^2.$$

Consider the term $(\mathbf{w} - \hat{\mathbf{w}})^T \hat{\mathbf{w}}$:

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

Least Squares with $d \geq n$

Theorem: Minimum norm solution

Theorem (Minimum norm least squares solution). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, let $d \geq n$, and let $\text{rank}(\mathbf{X}) = n$. Then, $\hat{\mathbf{w}} = \mathbf{X}^+ \mathbf{y} = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^T \mathbf{y}$ is the exact solution (i.e., $\mathbf{X} \hat{\mathbf{w}} = \mathbf{y}$) with smallest Euclidean norm:

$$\|\mathbf{w}\|_2^2 \geq \|\hat{\mathbf{w}}\|_2^2 \text{ for all } \mathbf{w} \in \mathbb{R}^d.$$

Proof. Consider any solution $\mathbf{w} \in \mathbb{R}^d$, such that $\mathbf{X} \mathbf{w} = \mathbf{y}$. We can write \mathbf{w} 's Euclidean norm as:

$$\|\mathbf{w}\|^2 = \|(\mathbf{w} - \hat{\mathbf{w}}) + \hat{\mathbf{w}}\|^2 = \|\mathbf{w} - \hat{\mathbf{w}}\|^2 - 2(\mathbf{w} - \hat{\mathbf{w}})^T \hat{\mathbf{w}} + \|\hat{\mathbf{w}}\|^2.$$

Consider the term $(\mathbf{w} - \hat{\mathbf{w}})^T \hat{\mathbf{w}}$:

$$(\mathbf{w} - \hat{\mathbf{w}})^T \hat{\mathbf{w}} = (\mathbf{w} - \hat{\mathbf{w}})^T \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y} = (\mathbf{X} \mathbf{w} - \mathbf{X} \hat{\mathbf{w}})^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y} = 0,$$

because \mathbf{w} and $\hat{\mathbf{w}}$ are both exact solutions.

Least Squares with $d \geq n$

Theorem: Minimum norm solution

Theorem (Minimum norm least squares solution). Let $\mathbf{X} \in \mathbb{R}^{n \times d}$, let $d \geq n$, and let $\text{rank}(\mathbf{X}) = n$. Then, $\hat{\mathbf{w}} = \mathbf{X}^+ \mathbf{y} = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^\top \mathbf{y}$ is the exact solution (i.e., $\mathbf{X} \hat{\mathbf{w}} = \mathbf{y}$) with smallest Euclidean norm:

$$\|\mathbf{w}\|_2^2 \geq \|\hat{\mathbf{w}}\|_2^2 \text{ for all } \mathbf{w} \in \mathbb{R}^d.$$

Proof. Consider any solution $\mathbf{w} \in \mathbb{R}^d$, such that $\mathbf{X} \mathbf{w} = \mathbf{y}$. We can write \mathbf{w} 's Euclidean norm as:

$$\|\mathbf{w}\|^2 = \|(\mathbf{w} - \hat{\mathbf{w}}) + \hat{\mathbf{w}}\|^2 = \|\mathbf{w} - \hat{\mathbf{w}}\|^2 - 2(\mathbf{w} - \hat{\mathbf{w}})^\top \hat{\mathbf{w}} + \|\hat{\mathbf{w}}\|^2.$$

Consider the term $(\mathbf{w} - \hat{\mathbf{w}})^\top \hat{\mathbf{w}}$:

$$(\mathbf{w} - \hat{\mathbf{w}})^\top \hat{\mathbf{w}} = (\mathbf{w} - \hat{\mathbf{w}})^\top \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y} = (\mathbf{X} \mathbf{w} - \mathbf{X} \hat{\mathbf{w}})^\top (\mathbf{X} \mathbf{X}^\top)^{-1} \mathbf{y} = 0,$$

because \mathbf{w} and $\hat{\mathbf{w}}$ are both exact solutions. Therefore,

$$\|\mathbf{w}\|^2 = \|\mathbf{w} - \hat{\mathbf{w}}\|^2 + \|\hat{\mathbf{w}}\|^2 \implies \|\mathbf{w}\|^2 \geq \|\hat{\mathbf{w}}\|^2.$$

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.

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

We approximate by least squares:

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

Choose

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

the best approximate solution:

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

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

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

Choose

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

the minimum norm (exact) solution:

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

Least Squares: SVD Perspective

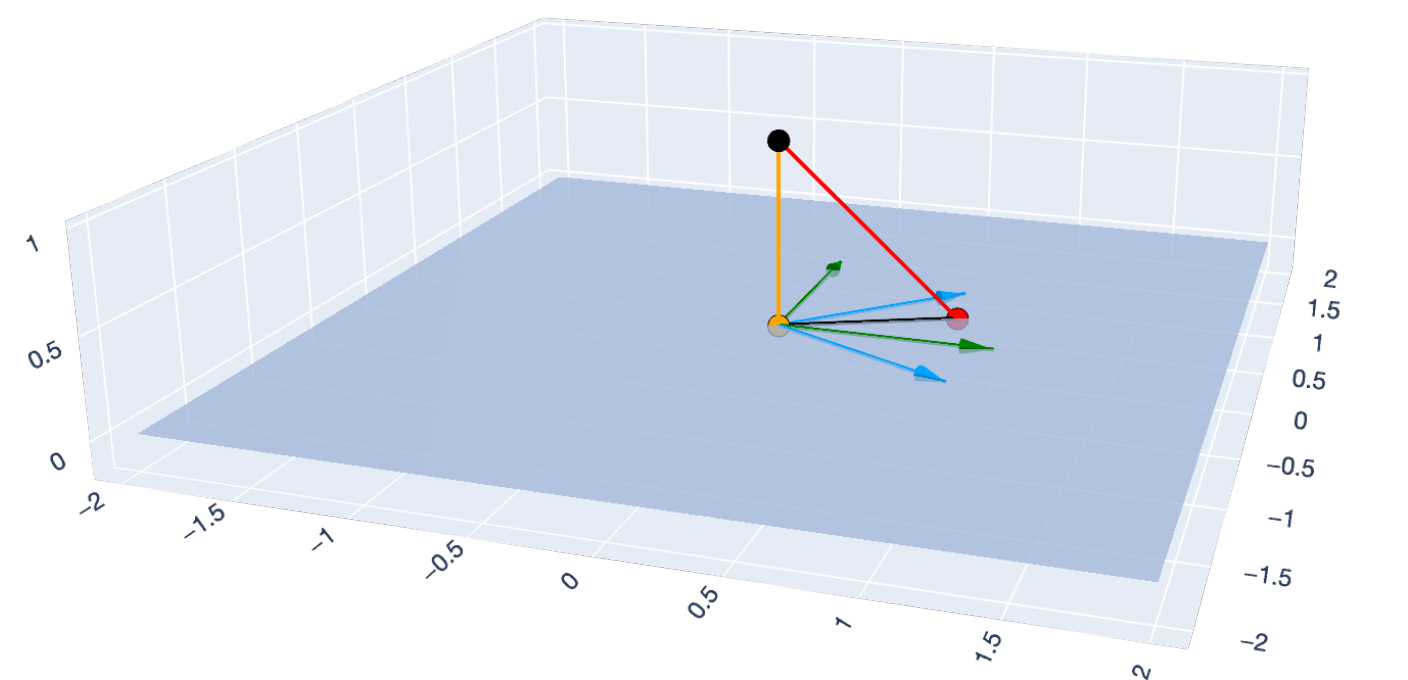
Unified Picture

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

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

We approximate by least squares:

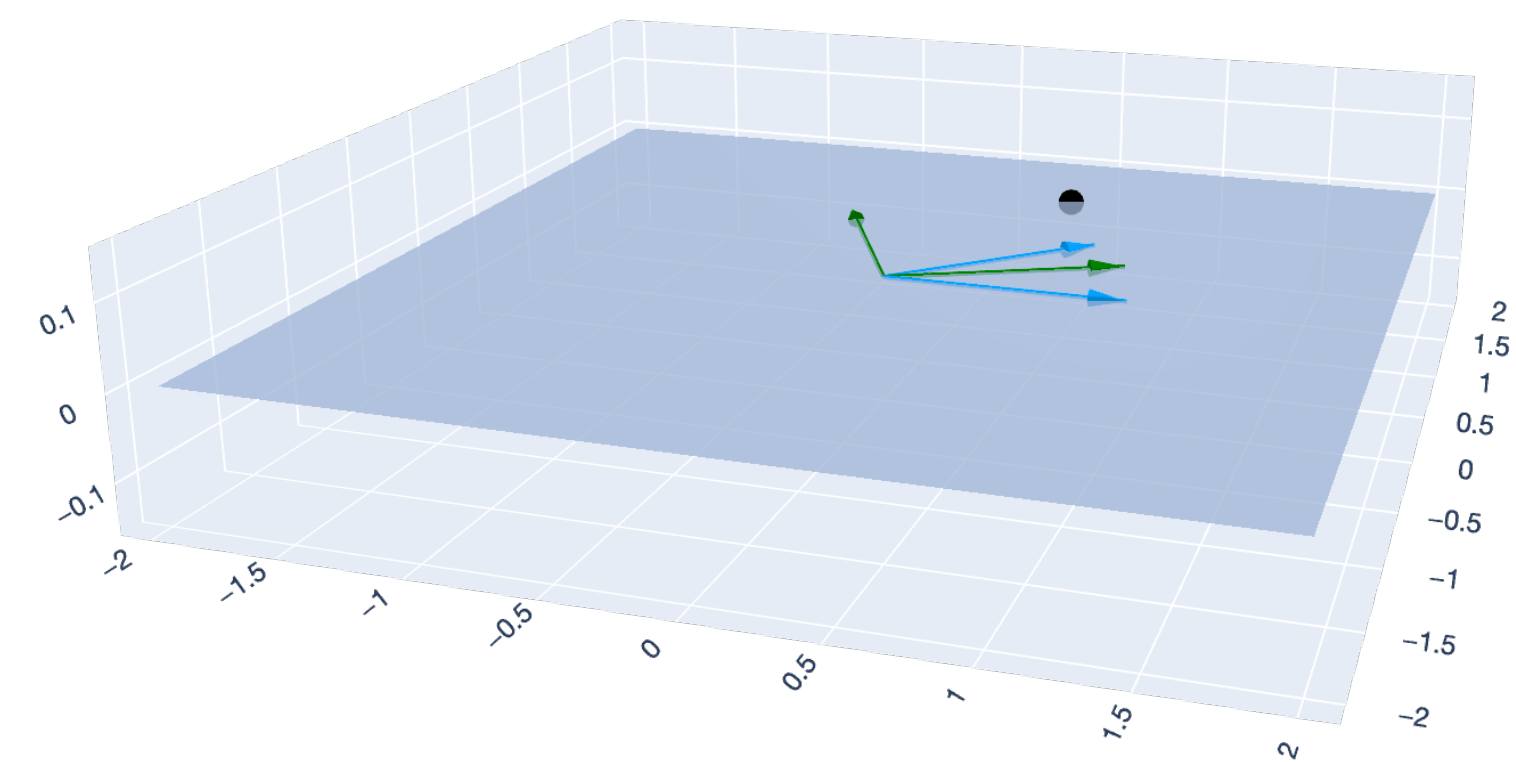
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2.$$



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

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

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

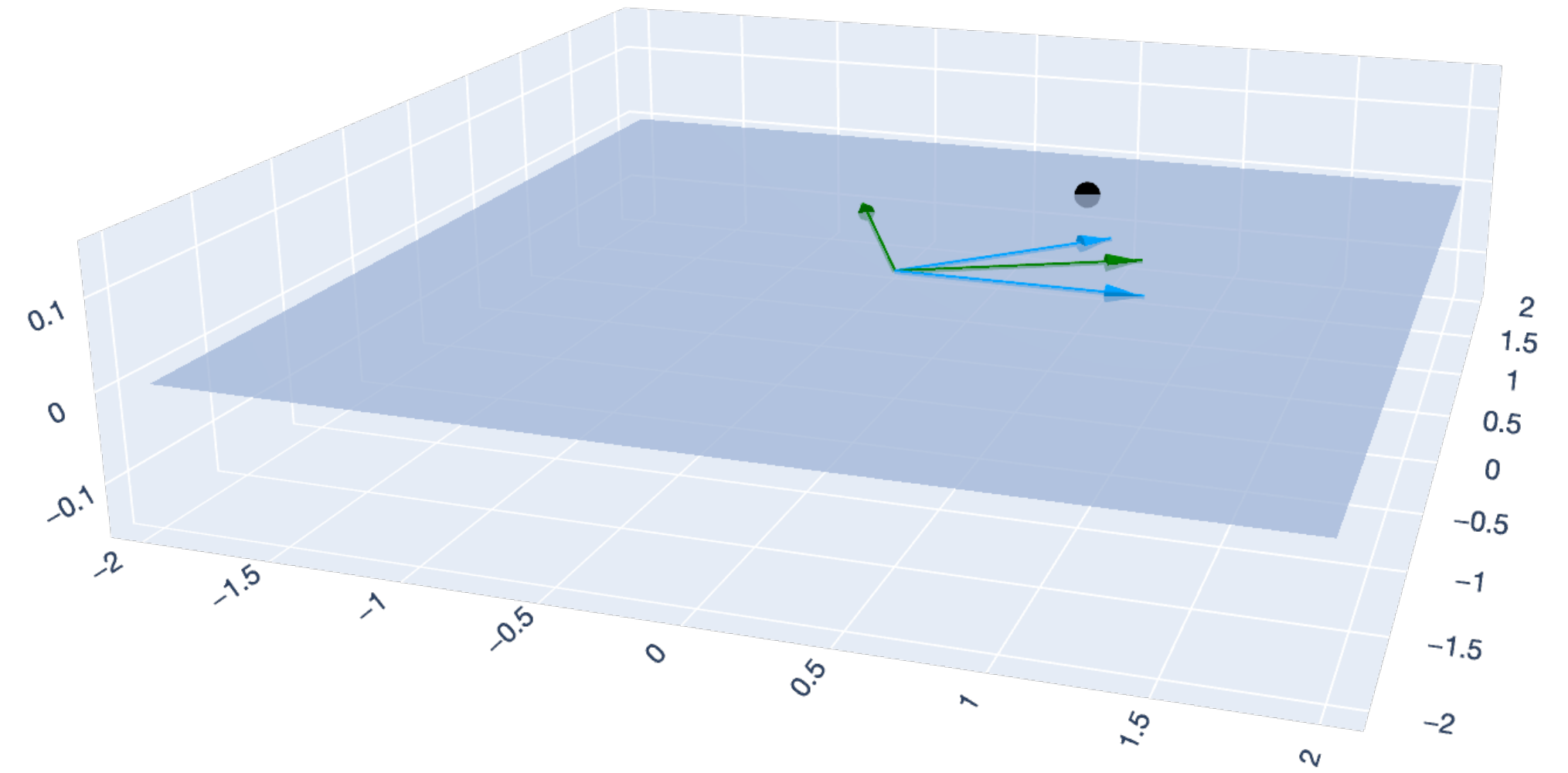
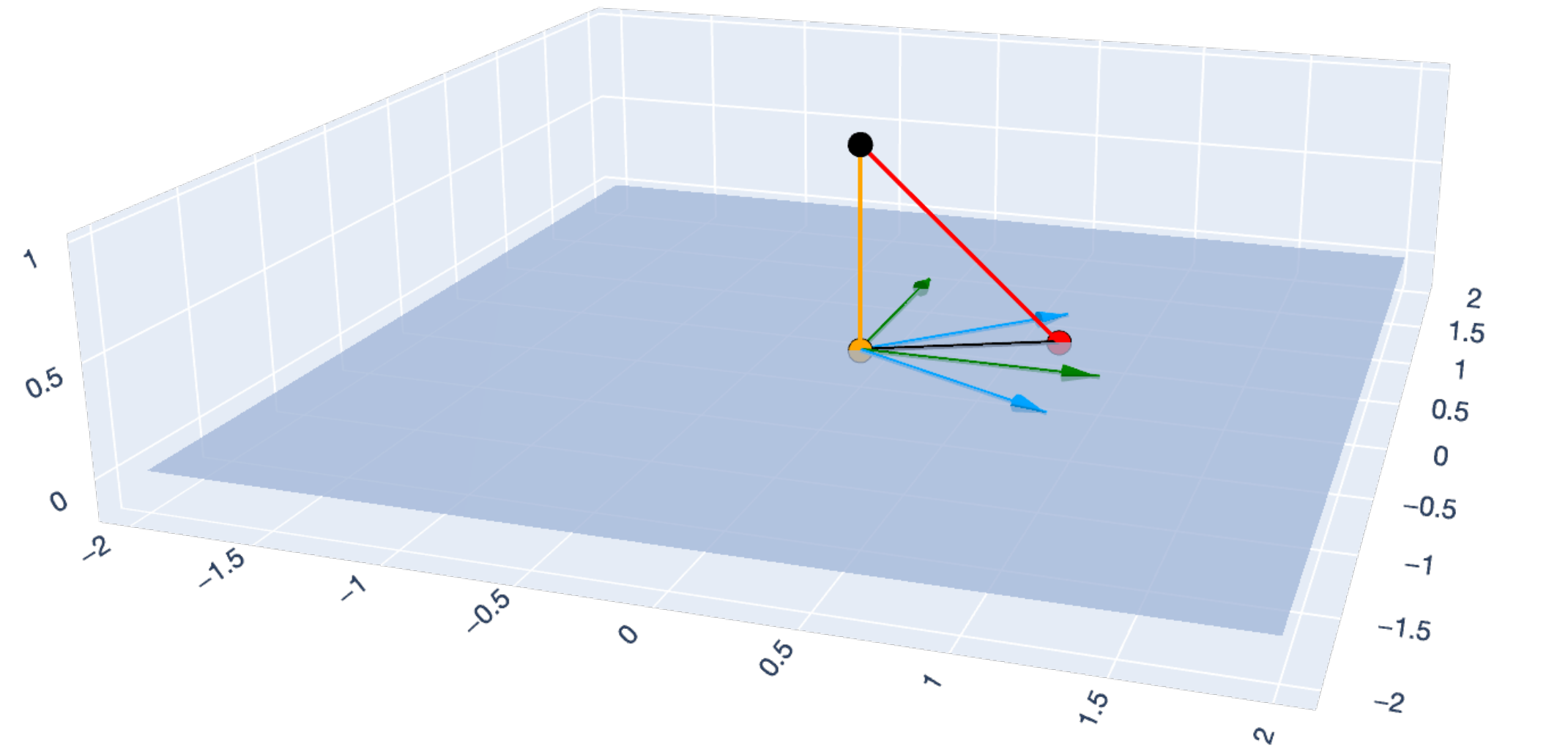


— x1 — x2 — u1 — u2 • y

Recap

Lesson Overview

Big Picture: Least Squares

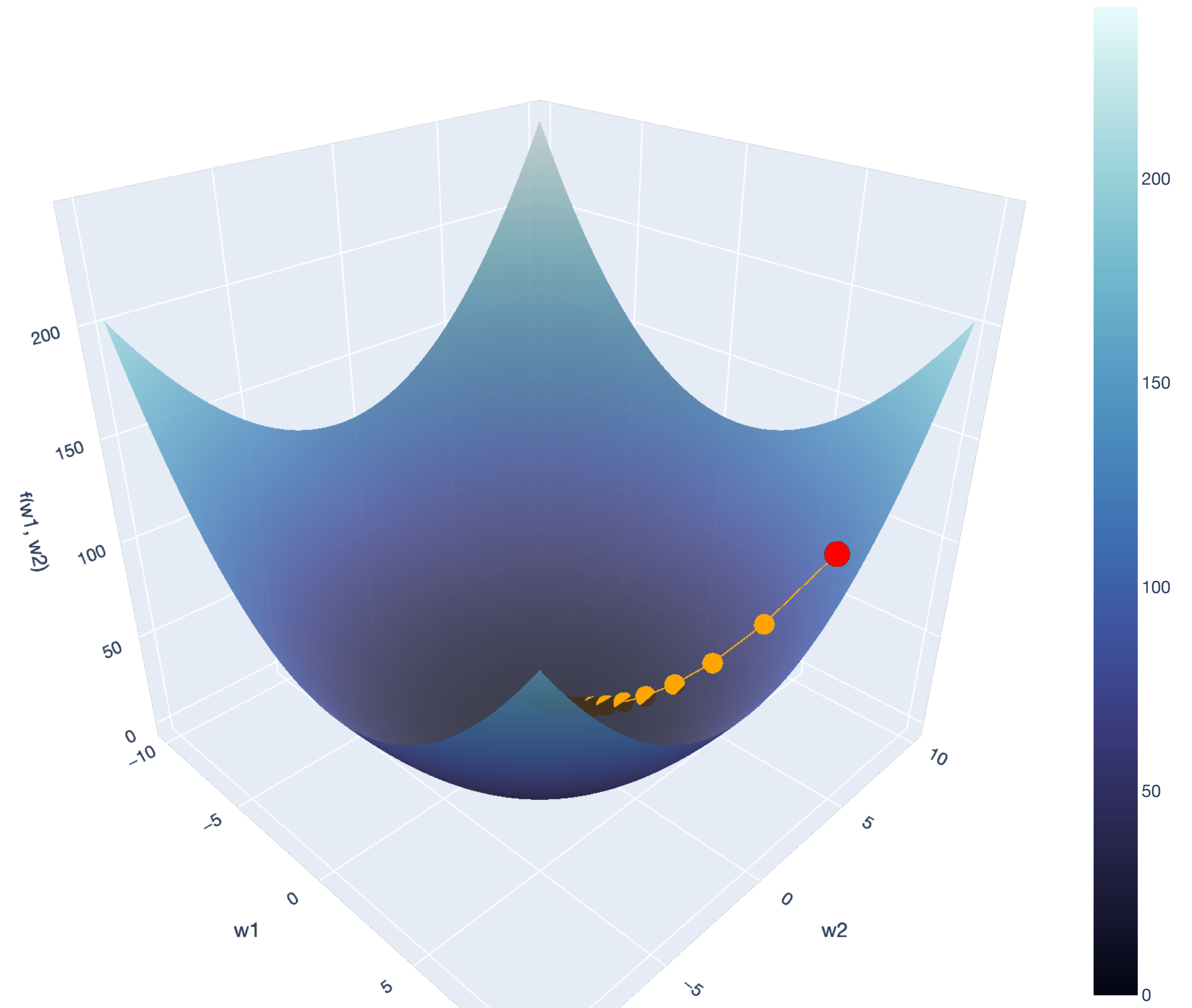
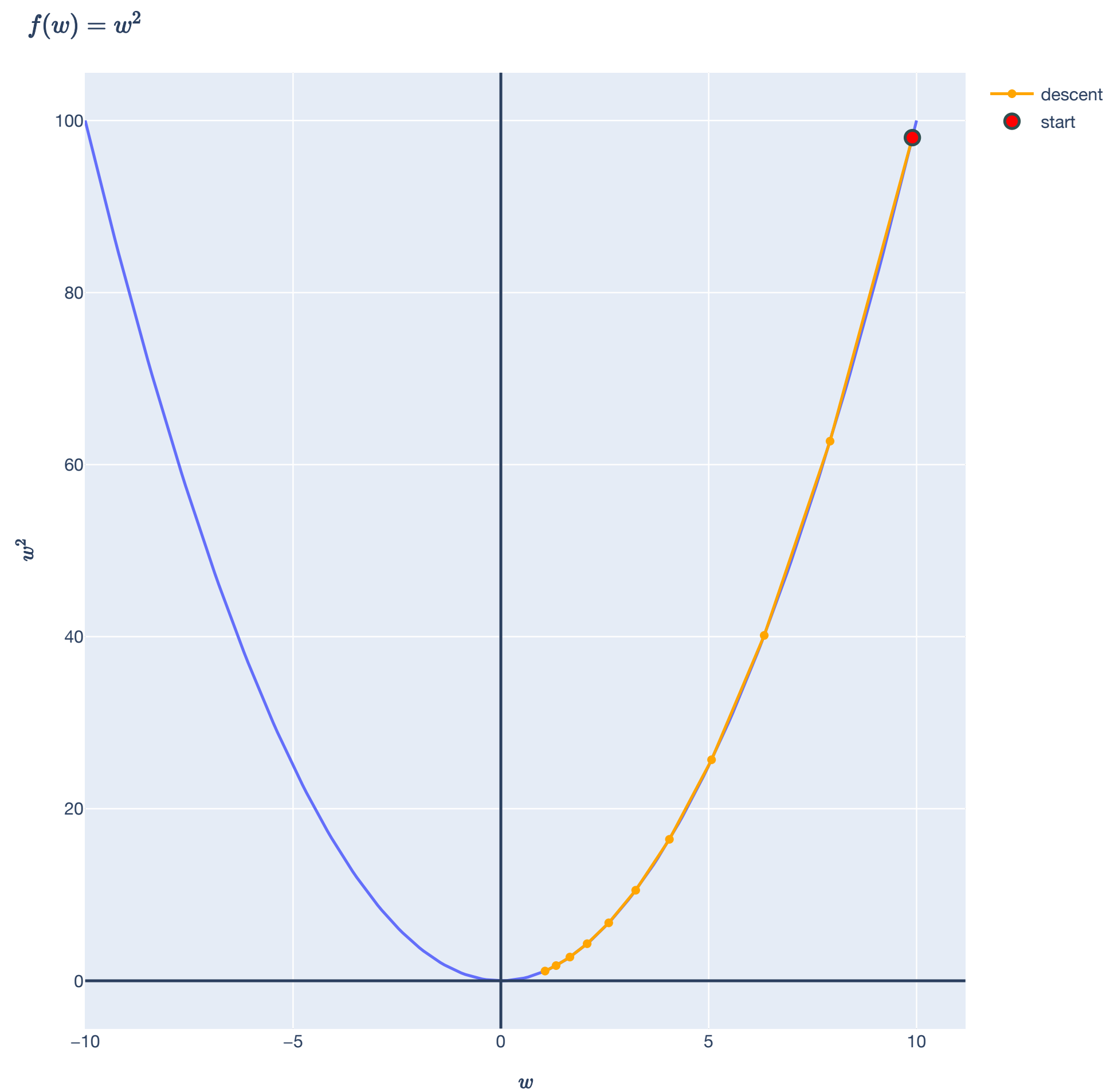


— x1
 — x2
 — u1
 — u2
 — $y - \hat{y}$
 — $\tilde{y} - \hat{y}$
 — $\tilde{y} - y$
 ● y
 ● \hat{y}
 ● \tilde{y}

— x1
 — x2
 — u1
 — u2
 ● y

Lesson Overview

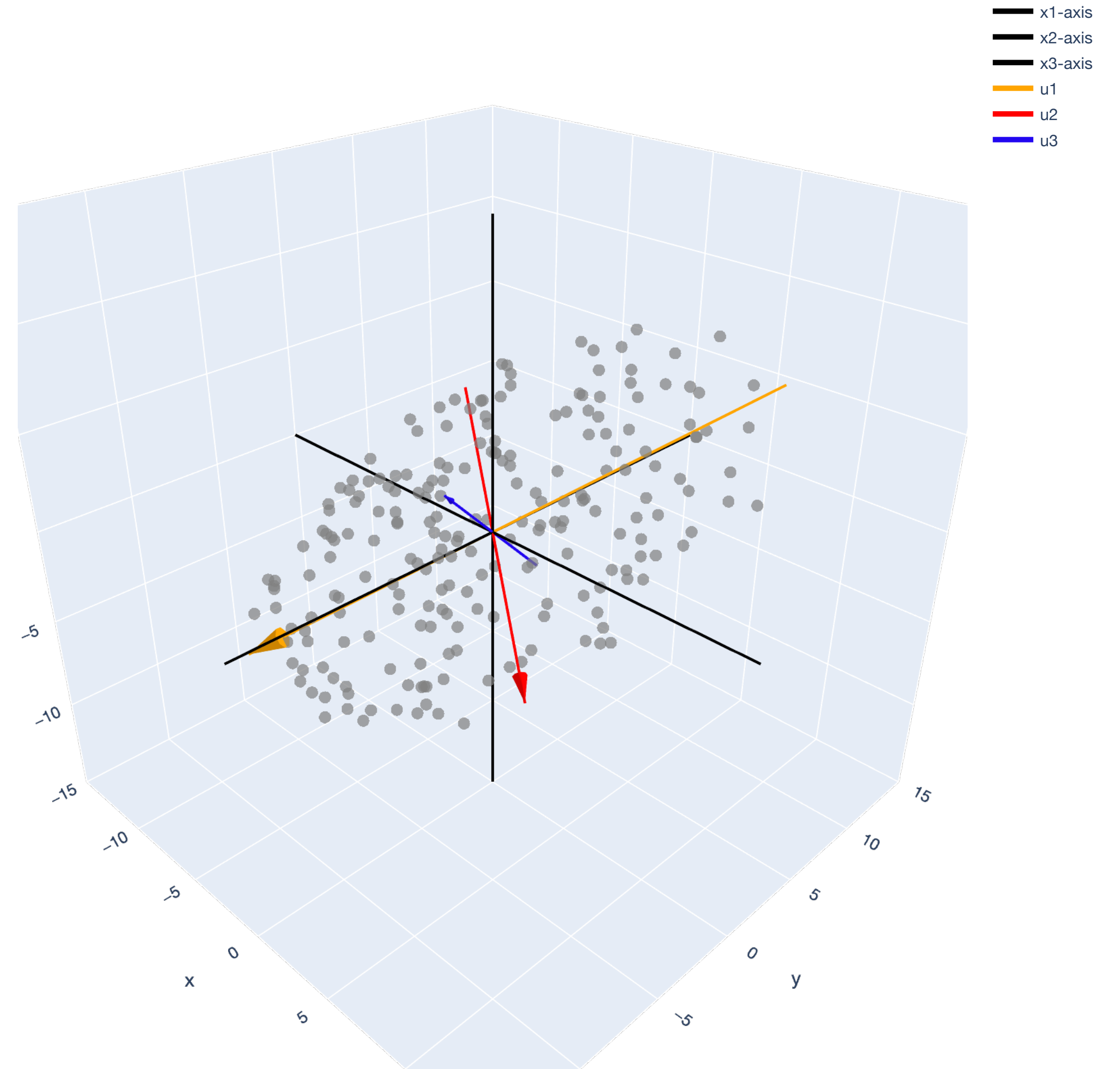
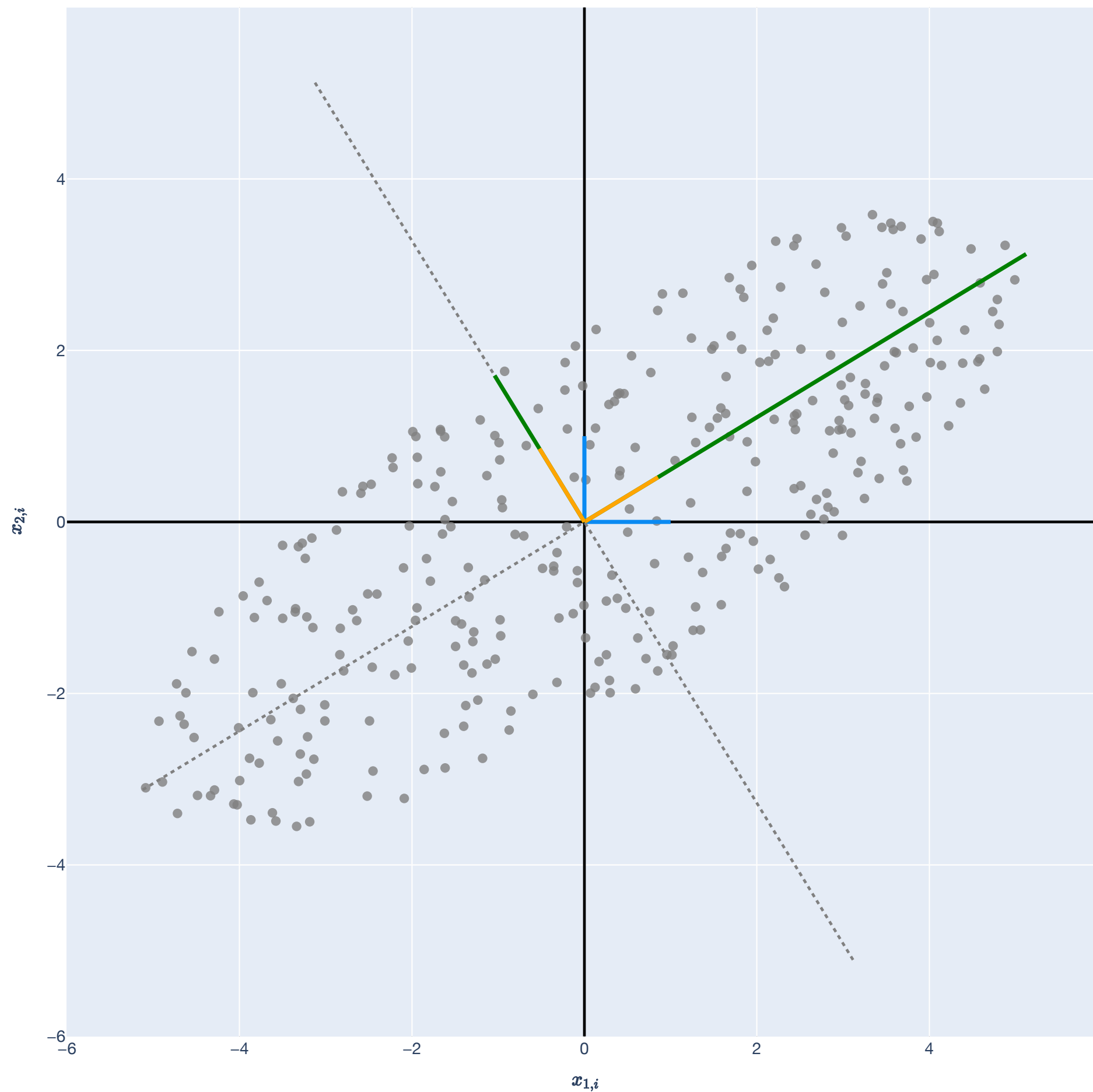
Big Picture: Gradient Descent



[Click to interact](#)

Lesson Overview

Big Picture: Singular Value Decomposition (SVD)



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: Singular Value Decomposition. Daniel Hsu.

Mathematical Foundations for Machine Learning. Rebecca Willett.

Elements of Statistical Learning: Data Mining, Inference, and Prediction. Trevor Hastie, Robert Tibshirani, and Jerome Friedman.