1.

Continue ON PCA. Principal Component Analysis
a method for doing dimensionality Reduction.
(its an unsupervised ML technique).
We will see why it can be a very bad idea
to use it naively for supervised Learning).

( SVD, eigenvalue decompositions and PCA are
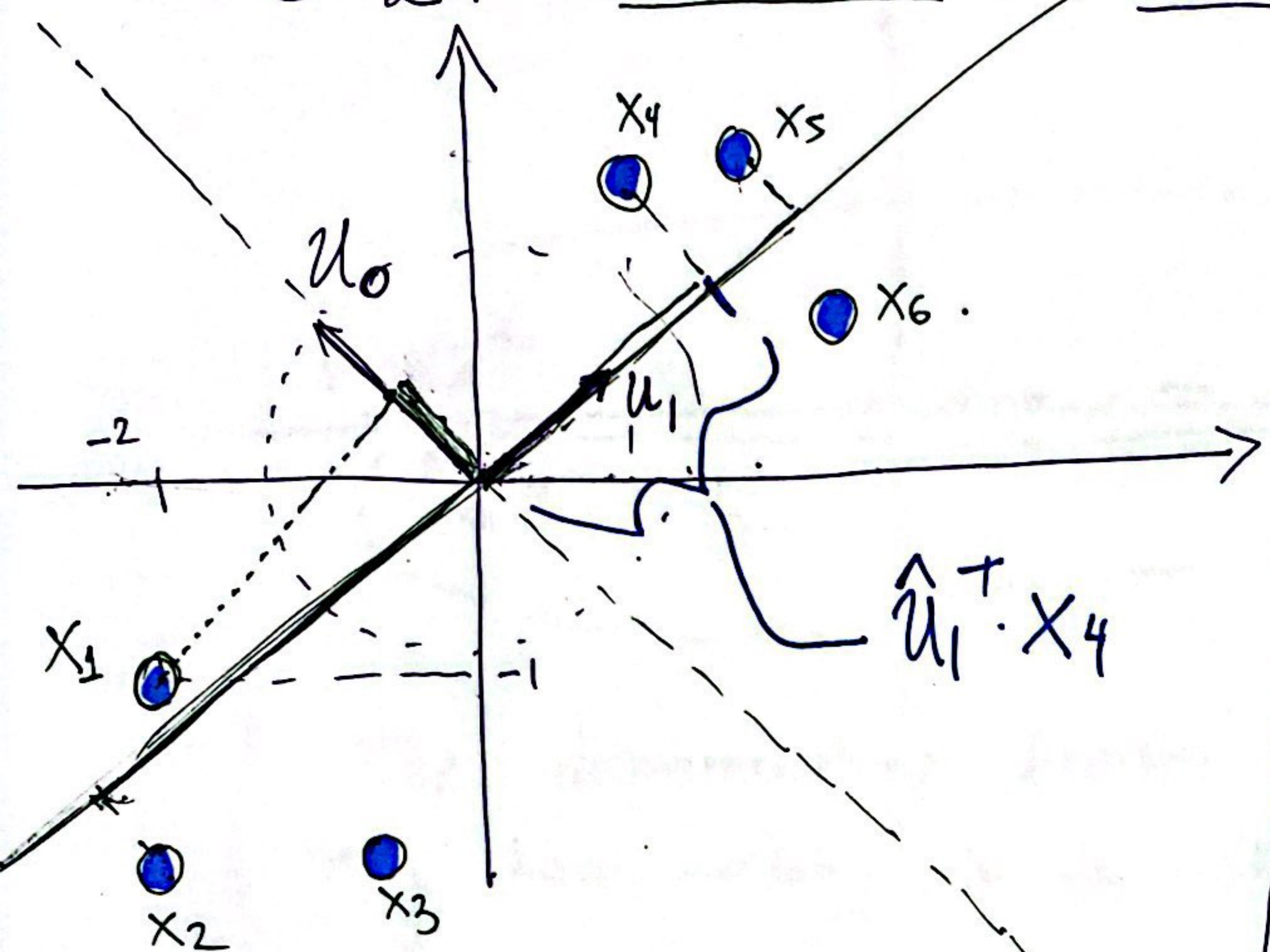closely connected tools. ).

Dimensionality Reduction : features are 10 dimensional
and you want to create new features that
are 1 dimensional OR 2D.

1. PCA is a projection of data points on a linear
subspace such that the variance of the
projected data is maximized.

2. This (turns out) to be equivalent to finding the
subspace that minimizes the error.

3. Also this turns out to be equivalent to
Finding new features that are uncorrelated.
- PCA whitens the data´ (white noise
symmetric in all directions).

[ This special subspace is called the principal subspace
and you can do it any dimension you want. ]

**2.**

Given a dataset of points $x_1, x_2 \ldots x_n \in \mathbb{R}^d$.

Lets see what the first principal component is.

$d = 2$. Centered data: $\underline{x_i^c = x_i - \bar{x}}$



$\hat{u_1}^T \cdot X_4$

since $x_1^T \cdot \hat{u_0} = 0.707$

the explained variance in the direction $\hat{u_0}$ is $0.707$.

---

Explained variance by $\hat{u_1}$

is $\left(\hat{u_1}^T \cdot x_1\right)^2 + \left(\hat{u_1}^T \cdot x_2\right)^2$

$+ \ldots \left(\hat{u_1}^T \cdot x_n\right)^2$

$= \sum_i \left(\hat{u_1}^T \cdot x_i\right)^2.$

if I stack $x_i$ as Rows of a Matrix $X = \frac{1}{n}\begin{bmatrix} -x_1- \\ -x_2- \\ \vdots \end{bmatrix}$

---

(Example)   $u_0 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$.

$X_1 = \begin{bmatrix} -2 \\ -1 \end{bmatrix}$.

· Project $X_1$ on $u_0$.

· First we normalize $u_0$.

$\hat{u_0} = u_0 \cdot \frac{1}{\|u_0\|} = u_0 \cdot \frac{1}{\sqrt{u_0^T u_0}}$

$\hat{u_0} = \frac{1}{\sqrt{2}} \cdot \begin{pmatrix} -1 \\ 1 \end{pmatrix}.$

Projection of $X_1$ on $u_0$.

is  $X_1^T \cdot \hat{u_0}$

$= \frac{1}{\sqrt{2}} \cdot (2 - 1) = \frac{1}{\sqrt{2}} = 0.707.$

$\hat{X_1} = (\text{Projected Length}) \cdot \overrightarrow{\text{Projected direction}}$
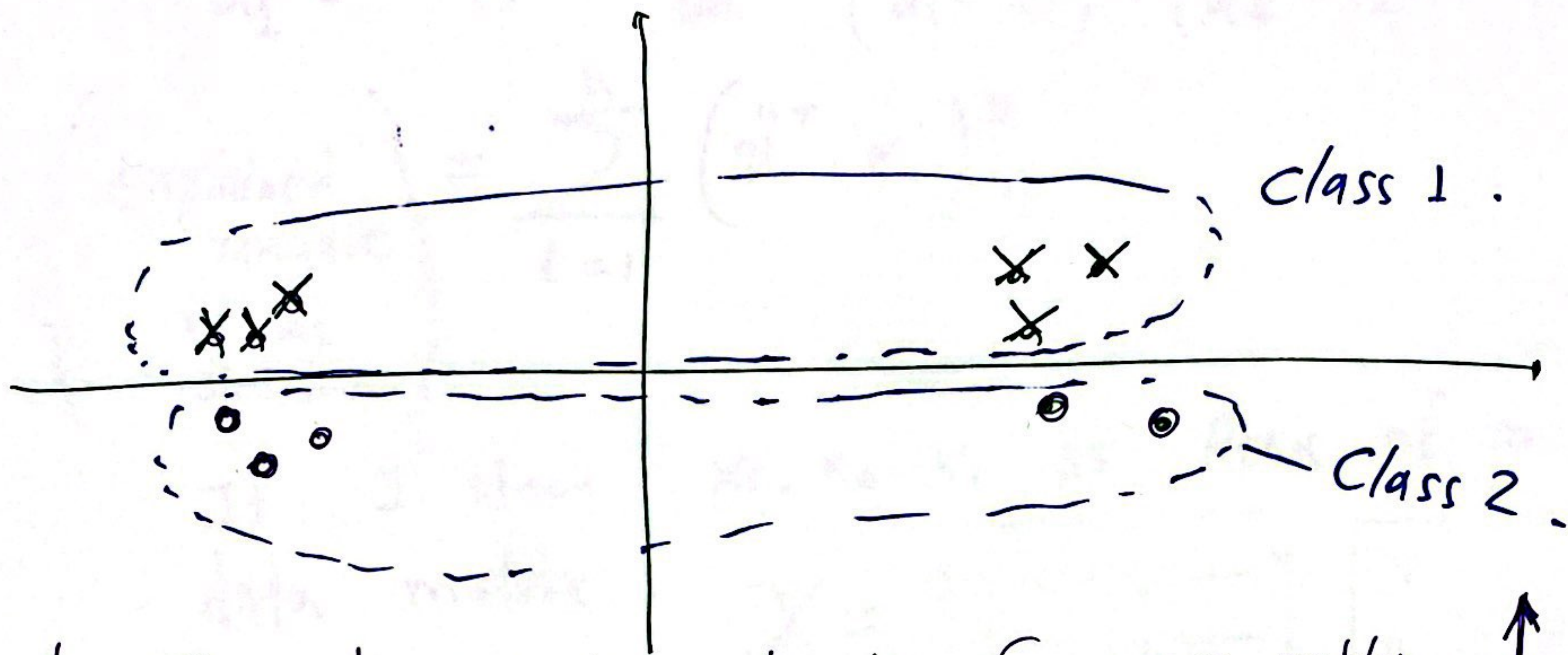
$= \frac{1}{\sqrt{2}} \cdot \hat{u_0}$

$\hat{X_1} = \left(X_1^T \cdot \hat{u_0}\right) \cdot \hat{u_0}.$

$= \frac{1}{\sqrt{u_0^T u_0}} \cdot \left(X_1^T u_0\right) \frac{1}{\sqrt{u_0^T u_0}} \cdot u_0$

$\hat{X_1} = \frac{1}{u_0^T u_0} \left(X_1^T u_0\right) \cdot u_0.$

When is PCA a terrible direction
1st component
for classification?



Class 1 .

Class 2 .

The discriminating direction for your problem
may have nothing to do with the direction of maximum
variance ($\longleftrightarrow$ .) .

( When this happens, PCA is a disaster for supervised
problems. )

Given data points $x_1, x_2 \ldots x_n \in \mathbb{R}^d$.
I am unit for a unit length vector $u_1$.
that maximizes

4

Given $X_1, X_2 .. X_n \in R^d$. (centered data!)

explained variance in a unit direction

$\hat{u}_1 \bullet$ is $\quad (\hat{u}_1^T \cdot X_1)^2 + (\hat{u}_1^T \cdot X_2)^2 + .. (\hat{u}_1^T \cdot X_n)^2.$

$\begin{pmatrix} \text{explained} \\ \text{variance} \\ \text{in } \hat{u}_1 \\ \text{direction.} \end{pmatrix} = \sum_{i=1}^{n} (\hat{u}_1^T \cdot X_i)^2.$

If I stack $X_1, X_2 .. X_n$ as Rows of a data matrix

$$X = \uparrow_n \begin{bmatrix} - X_1 - \\ - X_2 - \\ \vdots \\ - X_n - \end{bmatrix} \xleftarrow{\quad d \quad}$$

V.

$\|Y\|^2 = V^T \cdot V.$ (Trick 1.).

$X \cdot \hat{u}_1 = \uparrow_n \begin{bmatrix} - X_1 - \\ - X_2 - \\ \vdots \\ - X_n - \end{bmatrix} \cdot \begin{bmatrix} | \\ \hat{u}_1 \\ | \end{bmatrix} = \begin{bmatrix} X_1^T \cdot \hat{u}_1 \\ X_2^T \cdot \hat{u}_1 \\ \vdots \\ X_n^T \cdot \hat{u}_1 \end{bmatrix}$

$(Xu)^T = u^T X^T.$ (trick 2).

$(ABC)^T = C^T \cdot B^T \cdot A^T.$

Explained variance $= \|X \cdot \hat{u}_1\|_2^2 = \sum (X_i^T \cdot \hat{u}_1)^2.$

1st Principal component is the vector $\hat{u}_1$

that $\max_{\|u_1\| = 1} \sum (X_i^T \cdot u_1)^2 = \max_{\|u_1\| = 1} \|X \cdot u_1\|_2^2.$

Using trick 1: $\max_{\|u_1\| = 1} (Xu)^T (Xu).$

trick 2: $\max_{\|u\| = 1} u^T X^T X \cdot u.$

This is the covariance matrix for centered data.

$\max_{\|u\| = 1} u^T C \cdot u.$