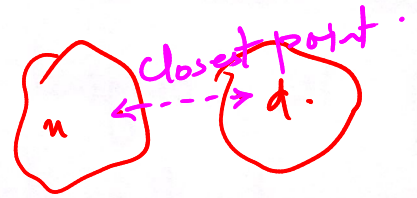


# Dimensionality Reductions (PCA)

Consider examples  $\underline{x}_i$  in a  $d$ -dimensional space  $\underline{x}_i \in \mathbb{R}^d$

A matrix  $\underline{W} \in \mathbb{R}^{n \times d}$  can project the examples  $\underline{x}_i$  to a  $n$ -dimensional space.



$$\underline{y}_i = \underline{W} \underline{x}_i$$

$n \times 1$        $n \times d$        $d \times 1$

The vector  $\underline{y}_i$  is the projection of  $\underline{x}_i$  onto an  $n$ -dimensional linear subspace.

We can also have a linear transformation defined by matrix  $\underline{U} \in \mathbb{R}^{d \times n}$  to project  $\underline{y}_i$  back to the  $d$ -dim subspace. The mapped point in the  $d$ -dimensional space is denoted as  $\tilde{\underline{x}}_i$ . It is the reconstruction of the original vector  $\underline{x}_i$ .

The Principal Component Analysis (PCA) problem can be defined as that of finding matrices  $\underline{W}_{n \times d}$  and  $\underline{U}_{d \times n}$  such that the total reconstruction error for all the examples is minimized.

$$\underline{U}, \underline{W} = \arg \min_{\underline{U}, \underline{W}} \sum_{i=1}^m \|\underline{x}_i - \tilde{\underline{x}}_i\|^2$$

$$= \arg \min_{\underline{U}, \underline{W}} \sum_{i=1}^m \|\underline{x}_i - \underline{U} \underline{W} \underline{x}_i\|^2$$

$$\underline{y}_i = \underline{W} \underline{x}_i$$

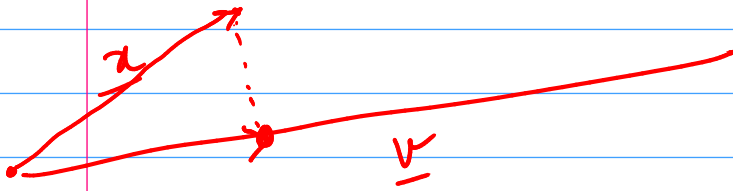
$$\tilde{\underline{x}}_i = \underline{U} \underline{y}_i$$

$$\tilde{\underline{x}}_i = \underline{U} \underline{W} \underline{x}_i$$

This is the first interpretation of the PCA problem.

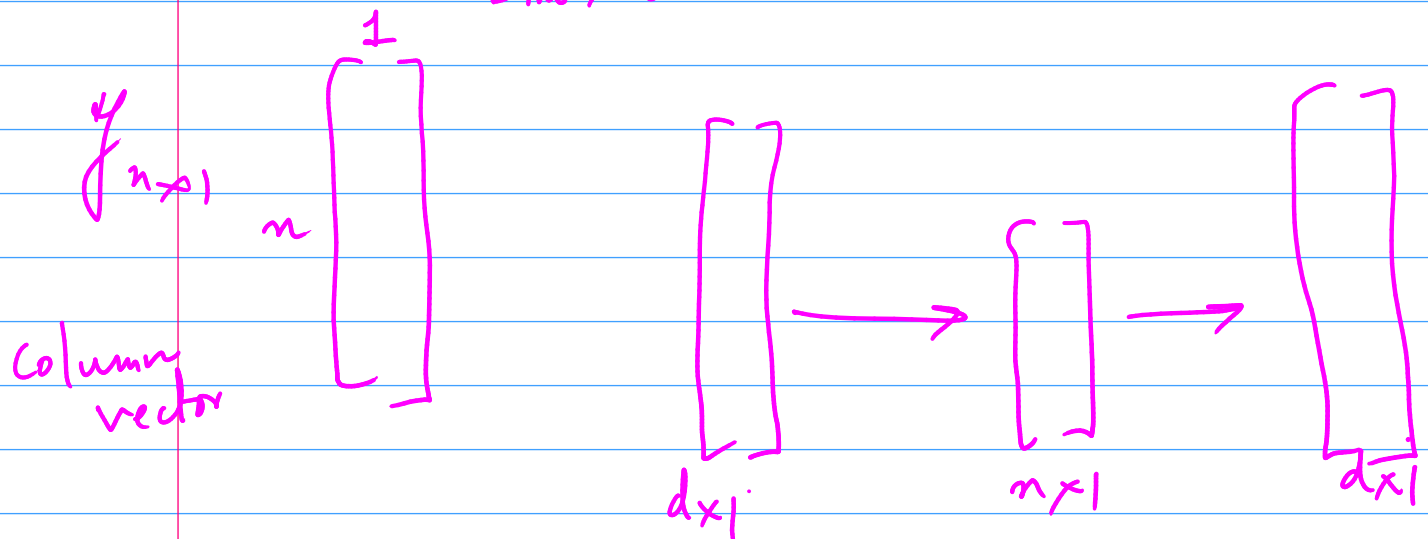
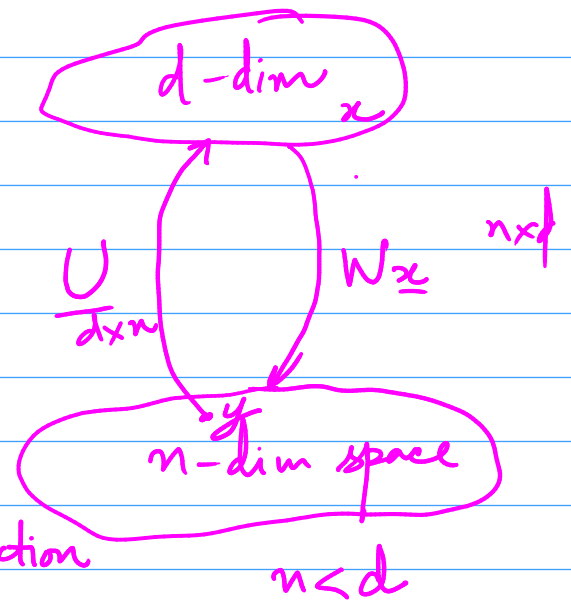
$$y_{n \times 1} = W_{n \times d} x_{d \times 1}$$

$W^{n \times d}$  has mapped a  $d$ -dim vector  
linear mapping to an  $n$ -dim vector



$$\underline{x} \xrightarrow{W} y \xrightarrow{U} \hat{x}$$

↑  
Dim reduction  
reconstruction



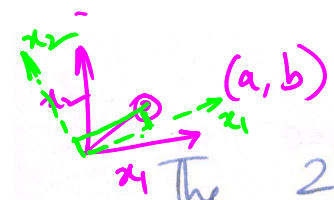
Feature Selection  
Dim Reduction

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \xrightarrow{FS} \begin{bmatrix} x_1 \\ x_4 \\ x_7 \\ \vdots \end{bmatrix}_{n < d}$$

Dim Reduction

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \xrightarrow{W} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \xleftarrow{U}$$





The 2<sup>nd</sup> interpretation of the PCA problem is to find new axis directions which are rotated versions of the original axes directions such that if we represent the data ( $\underline{x}_i \ i \in [1, m]$ ) in terms of the new features axes then the correlations and redundancies in the transformed feature values  $y_i$  are removed.

PCA can be used to automatically determine the correlation removing axis system.

$\equiv$  Variance maximizing direction

To apply PCA, we first need to create a data matrix

$X_{m \times d}$  data matrix

$$\underline{X} = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_m^T \end{bmatrix} \quad \begin{matrix} \text{ex } 1 \rightarrow \\ \vdots \\ \text{ex } m \rightarrow \end{matrix} \begin{bmatrix} x_1^1 & x_1^2 & x_1^3 & \dots & x_1^d \\ x_2^1 & x_2^2 & x_2^3 & \dots & x_2^d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_m^1 & x_m^2 & x_m^3 & \dots & x_m^d \end{bmatrix}$$

$\begin{matrix} \text{ex } 1 \rightarrow \\ \vdots \\ \text{ex } m \rightarrow \end{matrix}$

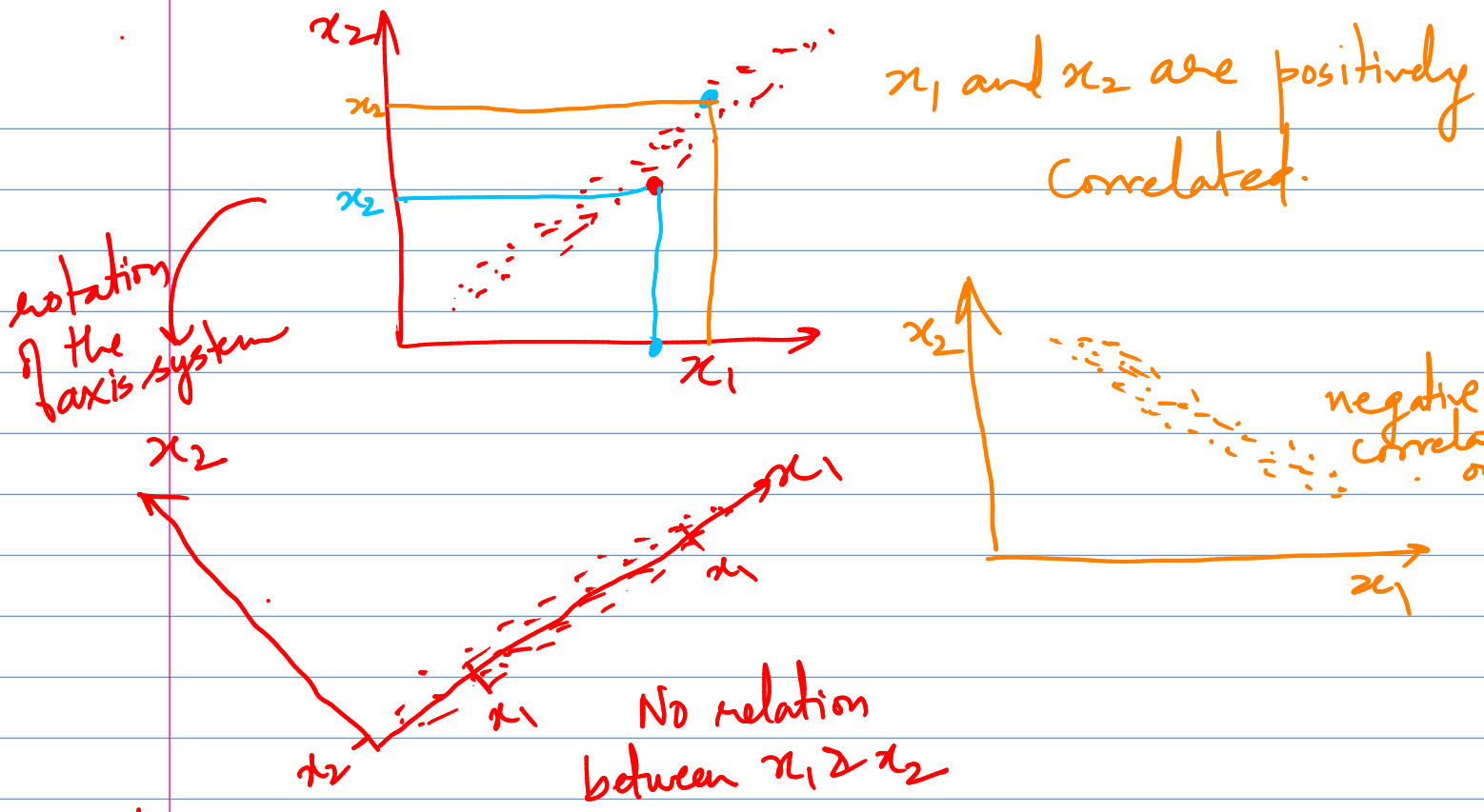
If the vectors  $\underline{x}_i$  are mean centered, then we can compute the covariance matrix as  $\frac{1}{m} \sum_{i=1}^m \underline{x}_i \underline{x}_i^T$   $d \times 1$   $1 \times d$   $d \times d$  outer product matrix

If the data is not mean centered, then the covariance matrix is given as

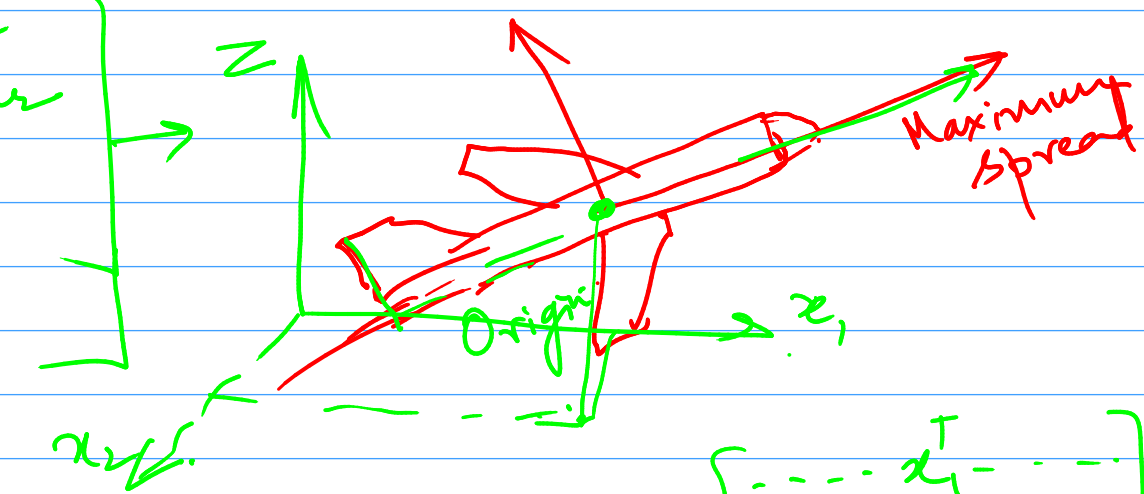
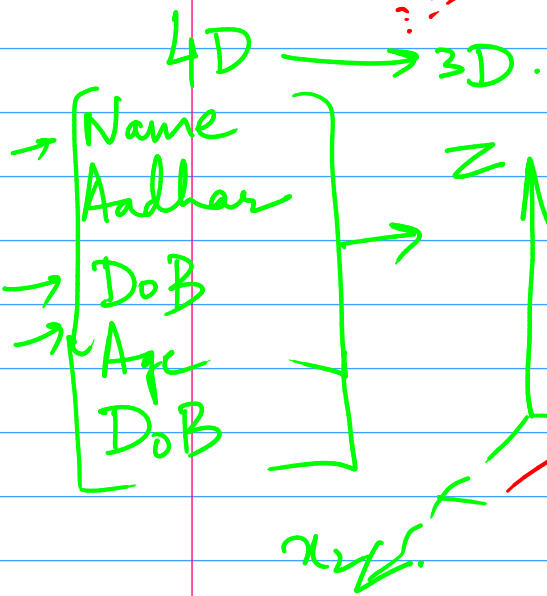
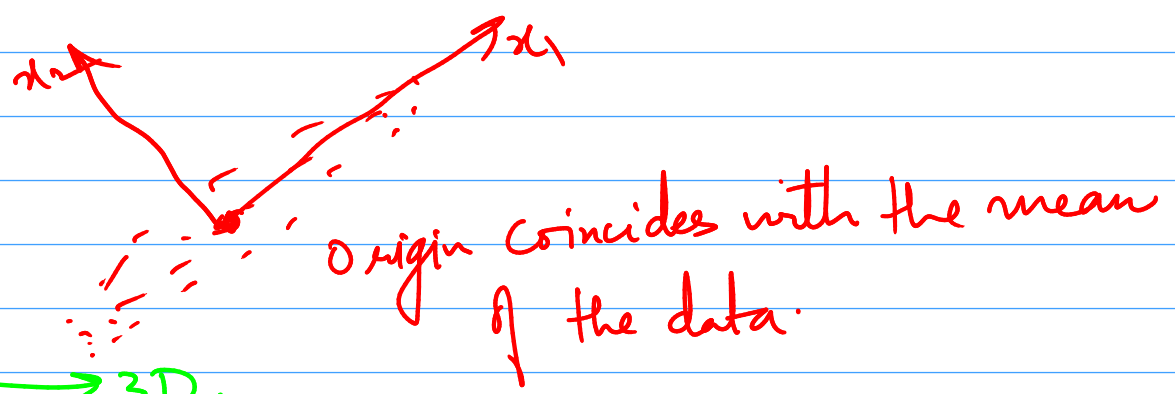
$$\frac{1}{m} \sum_{i=1}^m (\underline{x}_i - \underline{\mu})(\underline{x}_i - \underline{\mu})^T$$

$$= \frac{1}{m} \sum_{i=1}^m (\underline{x}_i - \underline{\mu})(\underline{x}_i^T - \underline{\mu}^T)$$

where  $\underline{\mu}$  is a  $d \times 1$  mean vector of the data set.



Features  $x_1$  &  $x_2$  have been de correlated.



$x^T x =$

$$\underline{X} = \begin{bmatrix} \dots & \underline{x}_1^T & \dots \\ \dots & \underline{x}_2^T & \dots \\ \dots & \dots & \dots \end{bmatrix} x^T \begin{bmatrix} x_1 & x_2 & \dots \\ \vdots & \vdots & \vdots \end{bmatrix}$$



$$= \frac{1}{m} \sum_i \left( \underline{x}_i \underline{x}_i^T - \underline{x}_i \underline{\mu}^T - \underline{\mu} \underline{x}_i^T + \underline{\mu} \underline{\mu}^T \right)$$

$$= \frac{1}{m} \left( \sum_i \underline{x}_i \underline{x}_i^T - \left( \sum_i \underline{x}_i \right) \underline{\mu}^T - \underline{\mu} \left( \sum_i \underline{x}_i^T \right) + \sum_{i=1}^m \underline{\mu} \underline{\mu}^T \right)$$

$$= \frac{1}{m} \left( \sum_i \underline{x}_i \underline{x}_i^T - m \underline{\mu} \underline{\mu}^T - m \cancel{\underline{\mu} \underline{\mu}^T} + m \cancel{\underline{\mu} \underline{\mu}^T} \right)$$

$$= \boxed{\frac{1}{m} \sum_{i=1}^m \underline{x}_i \underline{x}_i^T - \underline{\mu} \underline{\mu}^T} \quad \underline{C} \equiv \frac{\underline{X}^T \underline{X}}{m} - \underline{\mu} \underline{\mu}^T$$

We denote the covariance matrix as  $\underline{C} \equiv \frac{1}{m} \sum_i \underline{x}_i \underline{x}_i^T - \underline{\mu} \underline{\mu}^T$

The covariance matrix is a square  $d \times d$  matrix.

The diagonal entries of  $\underline{C}$  give the variance values along the  $d$  feature axes.

The off diagonal entries give the covariance values for the pairs of feature components.

The covariance matrix  $\underline{C}$  is positive semi-definite. 1x1 scalar

We find that for any vector  $\underline{v}$ , the value of  $\underline{v}^T \underline{C} \underline{v} \geq 0$ . 1x1 scalar

To show this we write the covariance matrix as

$$\underline{C} = \frac{1}{m} \underline{X}^T \underline{X} - \underline{\mu} \underline{\mu}^T$$

where  $\underline{X}$  is  $m \times d$  data matrix

$$\underline{X} : m \begin{bmatrix} \underline{x}_1^T & \dots & \dots \\ \underline{x}_2^T & \dots & \dots \\ \vdots & \ddots & \vdots \\ \underline{x}_m^T & \dots & \dots \end{bmatrix}$$

$$\underline{v}^T \underline{C} \underline{v} = \frac{\underline{v}^T \underline{X}^T \underline{X} \underline{v}}{m} - \underline{v}^T \underline{\mu} \underline{\mu}^T \underline{v}$$



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

$$= \left[ \frac{(\underline{X} \underline{v})^T (\underline{X} \underline{v})}{m} - (\underline{1}^T \underline{v})^T (\underline{1}^T \underline{v}) \right] \geq 0$$

Variance =  $\frac{1}{m} \sum x_i^2 - \bar{x}^2$

$\Rightarrow C$  is psd.

$\underline{X} \underline{v}$

row 1  
row 2  
row 3

col vec

1D proj

$\underline{1}^T \underline{v}$

1D projection

variance of the projection of the row vectors in  $\underline{X}$  on the column vector  $\underline{v}$ .

$\geq 0$

$\underline{1}^T \underline{v}$  is the projection of the mean vector on the direction  $\underline{v}$ . This will also be the mean of the 1-D projections of the examples  $\underline{x}_i$  on  $\underline{v}$ .

The goal of PCA is to one-by-one determine the orthonormal vectors  $\underline{v}$  maximizing  $\underline{v}^T C \underline{v}$ .

variance along direction  $\underline{v}$

Because the covariance matrix is symmetric and positive semi-definite, we can obtain its eigen value decomposition

$$C = V \Lambda V^T$$

$V$ : orthonormal eigen vectors of  $C$ , as columns.

$\Lambda$ : diagonal matrix of eigen values

$\Lambda_{ii}$ : eigen value corresponding to  $i$ th eigenvector of  $V$ .

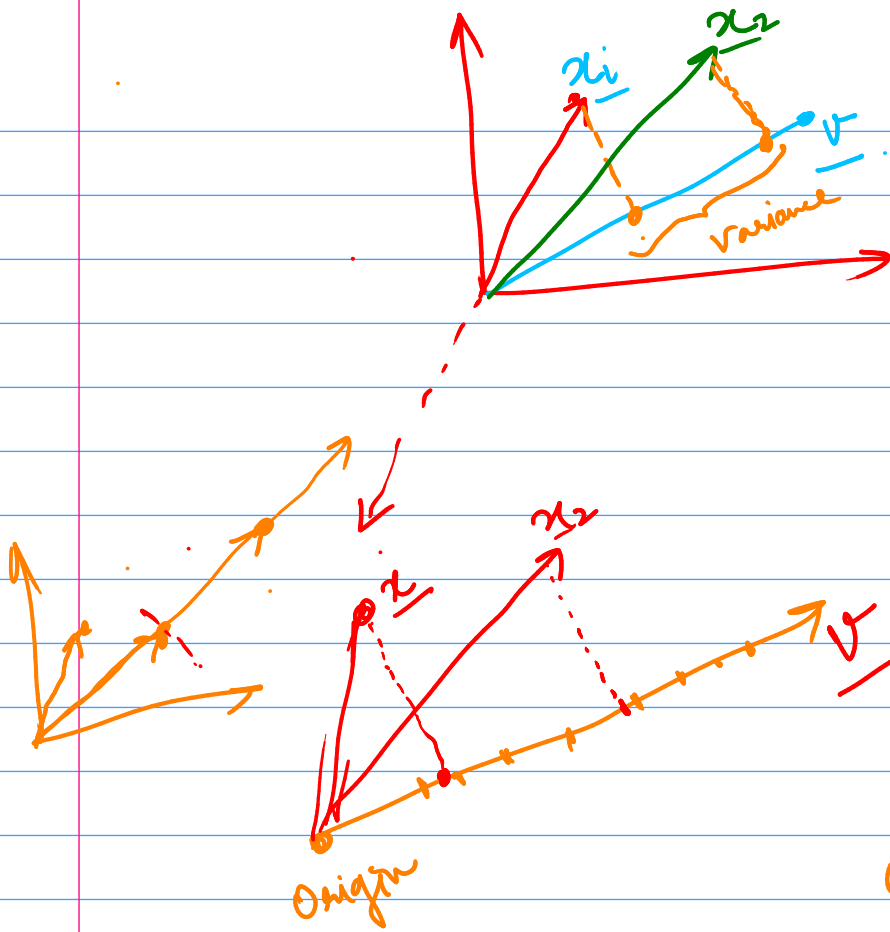
Objective is to maximize  $\underline{v}^T C \underline{v}$  subject to  $\underline{v}$  being a unit vector.

Formulating the Lagrangian  $\underline{v}^T C \underline{v} - \lambda (\|\underline{v}\|^2 - 1)$

Taking derivative w.r.t.  $\underline{v}$  and equating it to zero gives

1D projection

$\text{VC}$   
variance of all  
the points (examples)  
along direction  $\underline{v}$



$$\underline{v} = ( \dots )$$

direction

$\underline{v}$  defines a one-dim space

projecting any vector on  $\underline{v} \Rightarrow$  1-D projection