

Principal Component Analysis :-

MNIST \rightarrow 784 dim \rightarrow 2 dimensional & visualize -

d dims \rightarrow d' dimensions & $d' < d$

\rightarrow PCA is one of the simplest dimensionality reduction techniques.

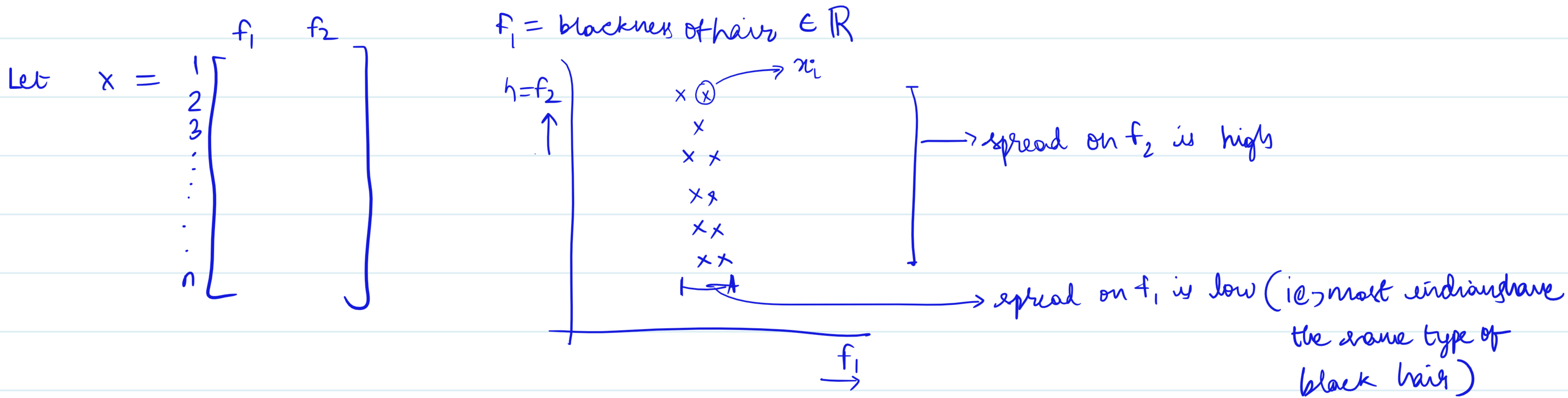
* Feature Engineering \rightarrow Reducing dimensionality by creating new features from given features (PCA, t-SNE)

* Feature Selection \rightarrow Reducing dimensionality by selecting subset of important features from given features.
(Forward Feature Selection, Backward Elimination)

\rightarrow Categorical features :- (contain finite number of categories/groups). They might not have a logical order (gender, country etc.)

Continuous features :- Numerical features that have an infinite number of values b/w two values (time, price etc.)

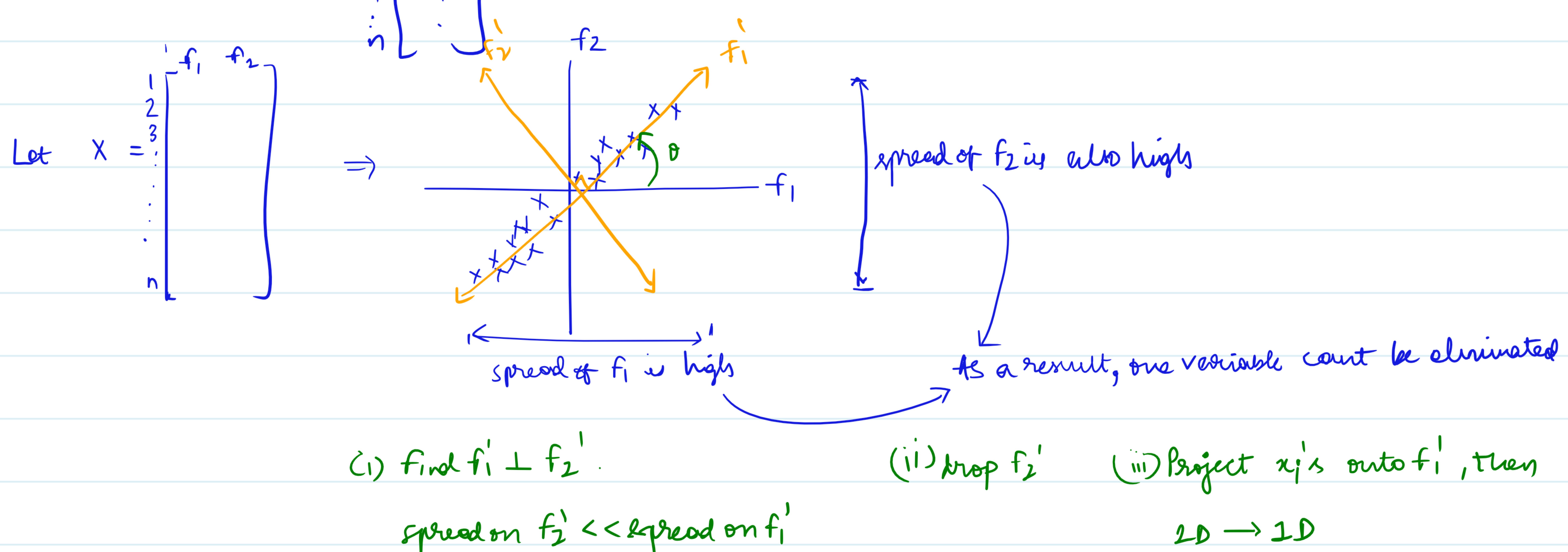
Geometric Intuition :-



If we want to convert this 2dimensional data to 1dimension

we can skip f_2 as it has low variance & keep only f_1

$\Rightarrow x' = \begin{bmatrix} f_1 \\ \vdots \\ f_1 \end{bmatrix} \Rightarrow$ We are preserving the direction with maximal spread/variance/information



\rightarrow We want to find directions of f_1' such that the variance of x 's projected onto f_1' is max. Achieved by rotating axis by θ° to find f_1' with max variance.

\rightarrow Before going for PCA, column standardization needs to be performed. As it brings features into a standard metric.

\rightarrow PCA is useless for dimensionality reduction if features are completely uncorrelated.

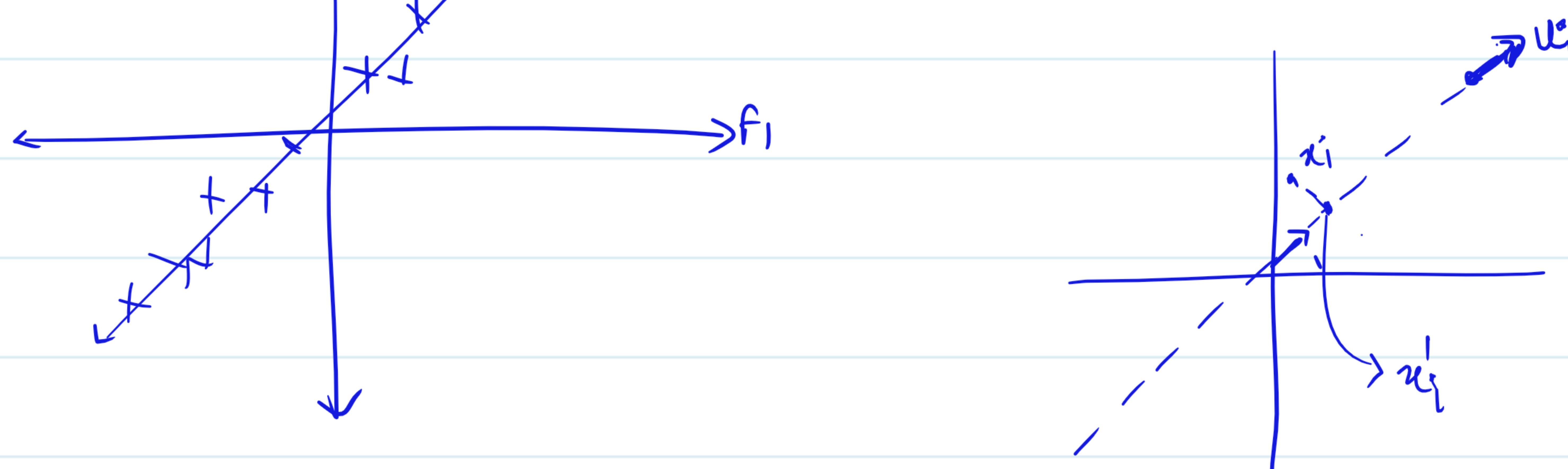
Variance Inflation :- When variance of dataset increases with addition of new columns that are completely unrelated to already existing features of dataset.
(Completely correlated = repeated features).

First principal component :- The direction in space along which projections have the largest variance.

Second principal component :- Direction which maximizes variance among all directions orthogonal to the first.

Mathematical objective function of PCA :-

$f_1 = u_1$ we only need the direction u_1 is the unit vector/direction $\Rightarrow \|u\| = 1$



x_i^1 = projection of x_i on u_1

$$D = \{x_i\}_{i=1}^n$$

$$D' = \{x_i^1\}_{i=1}^n$$

$$x_i^1 = \text{proj}_{u_1} x_i = \frac{u_1 \cdot x_i}{\|u_1\|^2} = u_1^T x_i$$

$$\begin{aligned} x_i^1 &= u_1^T x_i \\ \Rightarrow \bar{x}_i^1 &= u_1^T \bar{x}_i \quad \text{mean } \{x_i^1\}_{i=1}^n \\ &\text{mean } \{x_i\}_{i=1}^n \end{aligned}$$

Objective :- Find u_1 such that $\text{var}\{\text{proj}_{u_1} x_i\}_{i=1}^n$ is maximal.

$$\text{var}\{\text{proj}_{u_1} x_i\}_{i=1}^n = \frac{1}{n} \sum_{i=1}^n (u_1^T x_i - \bar{u}_1^T \bar{x})^2$$

$$u_1^T x_i = (u_1^T)_{1 \times n} x_i_{n \times 1} \rightarrow \text{scalar}$$

If X is column standardized, then mean becomes 0,
⇒ this value = 0

$$\Rightarrow \text{var}\{x_i^1\}_{i=1}^n = \frac{1}{n} \sum_{i=1}^n (u_1^T x_i)^2$$

We want to find u_1 such that this is maximized while keeping $u_1^T u_1 = 1 = \|u\|^2$

Objective of an optimization problem

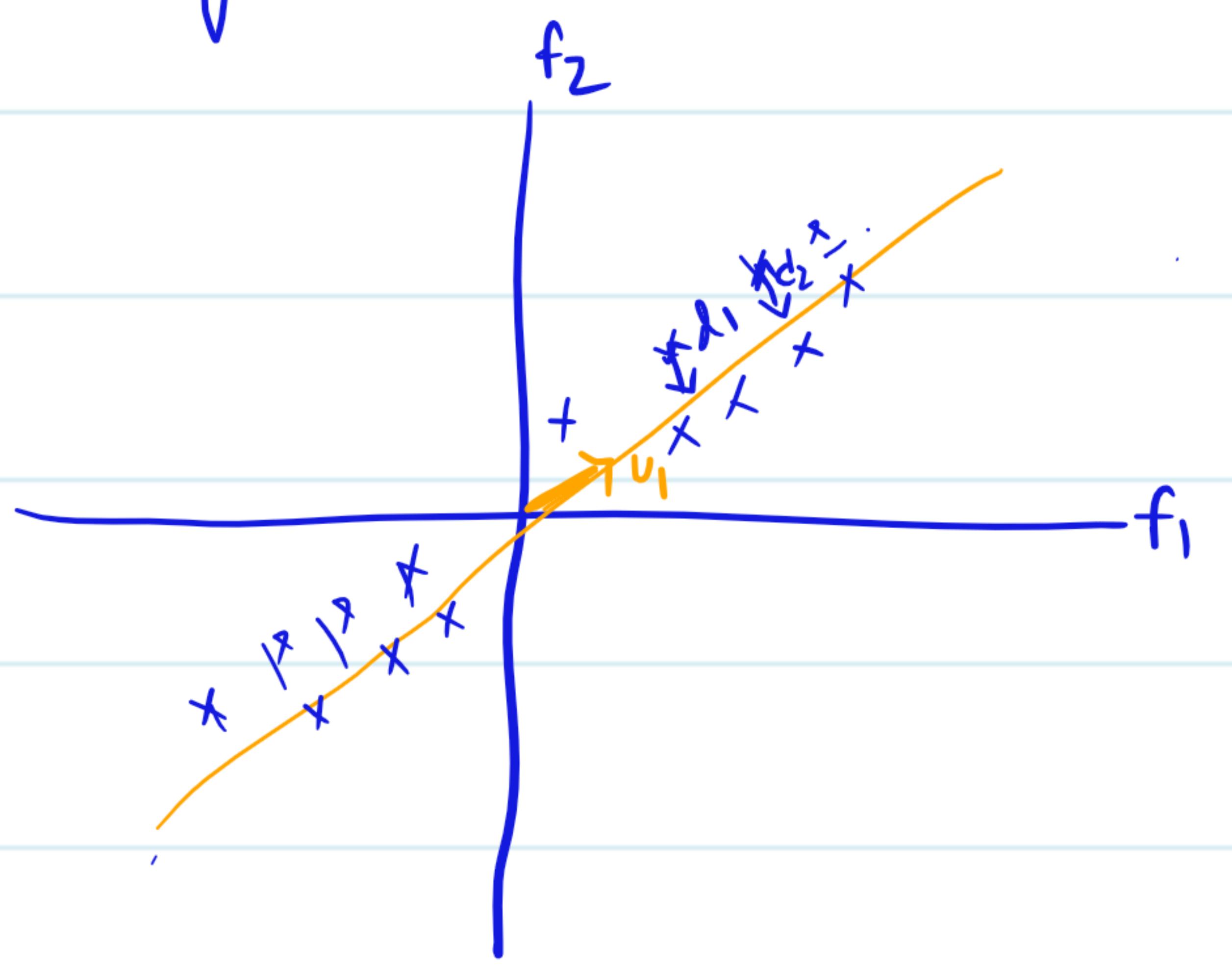
constraint

Optimization problem.

Data → column standardization → After standardization, all features have same spread, so we don't know which one to choose → Axis transformation → Find the unit vector in order for axis to line up with features that have the highest spread

Alternative Formulation of PCA :- Distance Minimization

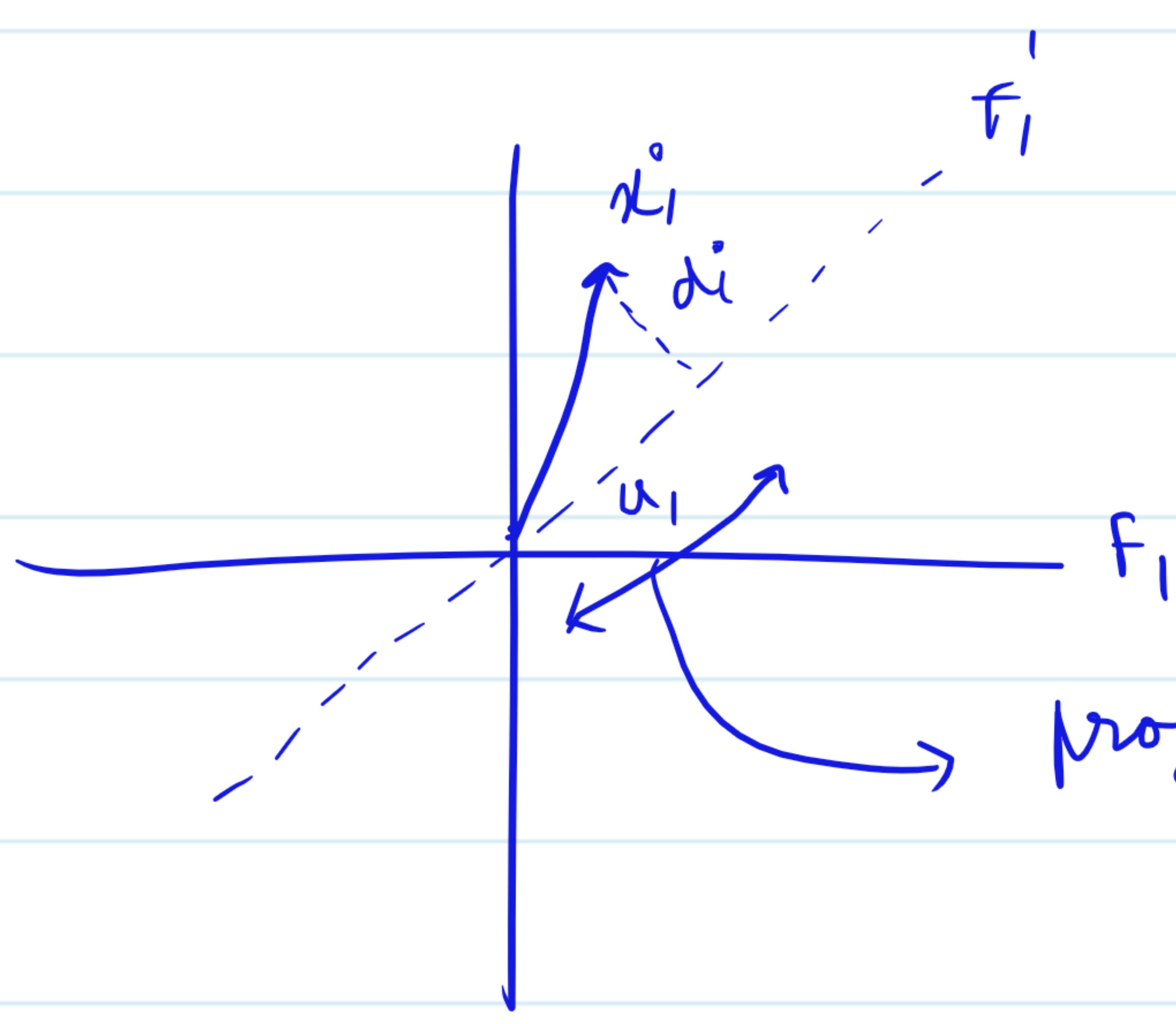
Previously - Find v_1 , that maximizes variance -



Alternatively we want to find out v_1 , that minimizes

$$\sum_{i=1}^n d_i^2$$

$$\min_{v_1} \sum_{i=1}^n d_i^2$$



v_1 : unit vector $v_1^T v_1 = 1 = \|v_1\|^2$

projection of x_i on v_1 = $\text{proj}_{v_1} x_i = v_1^T x_i$

$$\begin{aligned} \text{Pythagorean theorem: } d_i^2 &= \|x_i\|^2 - (v_1^T x_i)^2 \\ &= x_i^T x_i - (v_1^T x_i)^2 \end{aligned}$$

$$\Rightarrow \min_{v_1} \sum_{i=1}^n (x_i^T x_i - (v_1^T x_i)^2) \text{ such that } v_1^T v_1 = 1 \rightarrow \begin{array}{l} \text{distance minimization} \\ \text{PCA} \end{array}$$

Previously, we had $\max_{v_1} \frac{1}{n} \sum_{i=1}^n (v_1^T x_i)^2$ such that $v_1^T v_1 = 1 \rightarrow \text{Variance Maximization PCA}$

The v_1 that achieves these two is the same

Eigen Values & Eigen Vectors in PCA :-

Solution to optimization problems :-

$$\text{Let } X = \begin{bmatrix} 1 & 2 & 3 & \dots & d \\ 2 & 3 & \vdots & & \\ 3 & & & & \\ \vdots & & & & \\ n & & & & \end{bmatrix}_{n \times d}$$

→ column standardized $\Rightarrow S_{d \times d} = X^T X$

square symmetric

Eigen values :- $\lambda_1, \lambda_2, \lambda_3, \dots$
Eigen vectors :- v_1, v_2, v_3, \dots

Eigen Values of $S = \lambda_1, \lambda_2, \dots, \lambda_d$

d values since S is a $d \times d$ matrix

Eigen Vectors = $v_1, v_2, v_3, \dots, v_d$

one eigen vector for corresponding
Eigen value.

Assume $\lambda_1 > \lambda_2 > \lambda_3 > \lambda_4 \dots > \lambda_d$

Maximal Eigen value.

Eigen Value def :- $\lambda, v_i = S_{d \times d} v_i$

↓
Scalar
↓
Eigen Value

↓
Vector
↓
Eigen Vector

If $\lambda_1 > \lambda_2 > \lambda_3 \dots > \lambda_d$

$v_1 \downarrow v_2 \downarrow v_3 \dots \downarrow v_d$

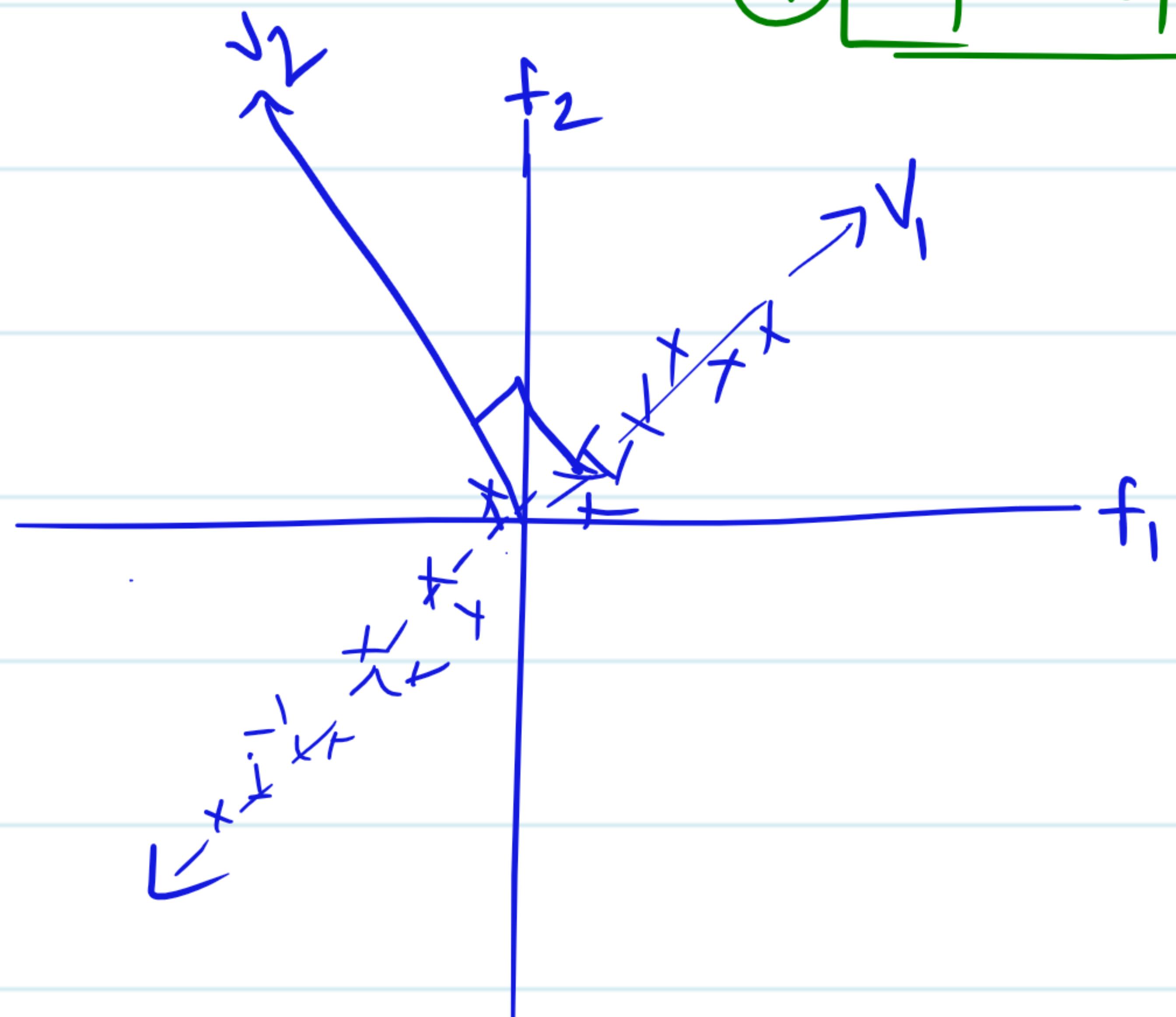
Any two $v_i \times v_j \rightarrow v_i \perp v_j \Rightarrow v_i^T v_j = 0 = v_i \cdot v_j$

$(v_1) = v_1$ = Eigen vector of $S (= X^T X)$ corresponding to the largest eigen value ($= \lambda_1$)

From Optimization problem.
Maximal Variance direction

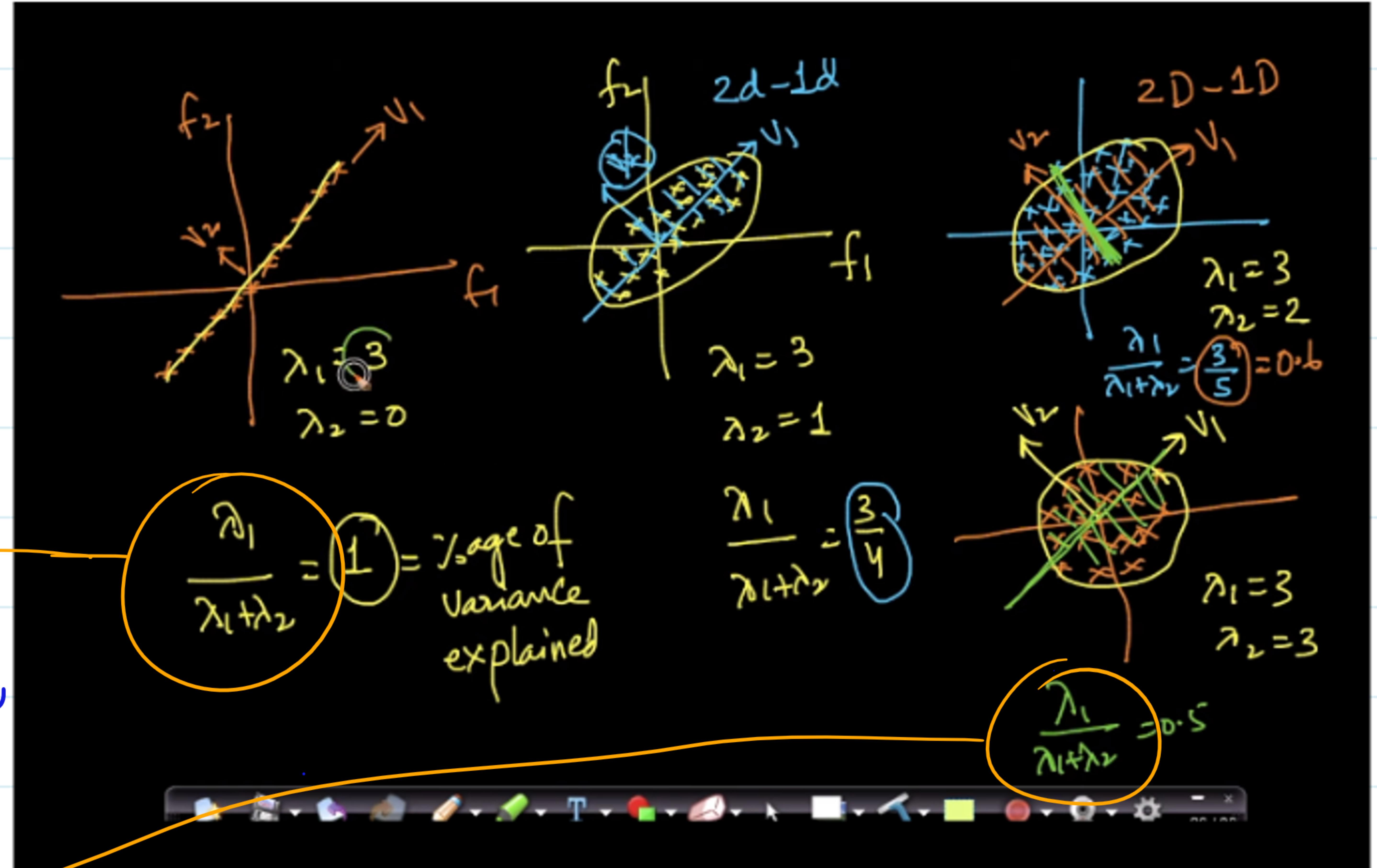
Steps to find v_1 :-

- ① Column Standardize X
- ② Find $S_{d \times d} = X^T X$
- ③ Compute 'd' Eigen Values & vectors
- ④ $v_1 = v_1$



v_1 → First eigen vector with maximal variance.
 v_2 → Second Eigen Vectors with maximal variance.
 v_d → Least Variance last eigen vector.

Significance of λ 's (Eigen values) :-



Since data is circular, 50% of it is on v_1 & 50% is on v_2 .

Eigen value :- How much data is present on each Eigen vectors.

Eigen Vector :- Which direction the first maximal is present & the second maximal is present.

PCA for Dimensionality Reduction & Normalization:-

$$X = \begin{bmatrix} f_1 & f_2 \\ \vdots & \vdots \\ 1 & 2 \\ 2 & 3 \\ \vdots & \vdots \\ n & n \end{bmatrix} \xrightarrow{\text{PCA}} X' = \begin{bmatrix} v_1 \\ \vdots \\ v_i \\ \vdots \\ n \end{bmatrix}$$

$$x'_i = x_i^T v_i \quad (\text{from max variance method PCA})$$

Alternative method :-

$$X' = \begin{bmatrix} v_1 & v_2 \\ \vdots & \vdots \\ 1 & 2 \\ 2 & 3 \\ \vdots & \vdots \\ n & n \end{bmatrix}$$

$\lambda_1 > \lambda_2 > \lambda_3 > \dots > \lambda_d$

$x'_i = [x_i^T v_1, x_i^T v_2]$

Eigen vals Eigen vectors (top 2)

→ PCA can not only be used for converting to 2D.

ex:- $X = \begin{bmatrix} 1 & 2 & 3 & \dots & 100 \\ 2 & 3 & 4 & \dots & 101 \\ 3 & 4 & 5 & \dots & 102 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ n & n+1 & n+2 & \dots & n+99 \end{bmatrix}$ → 100 features

$$\lambda_1 > \lambda_2 > \lambda_3 \dots > \lambda_d$$
$$v_1 > v_2 > v_3 \dots > v_d$$

↓ PCA → 100 to 50d.

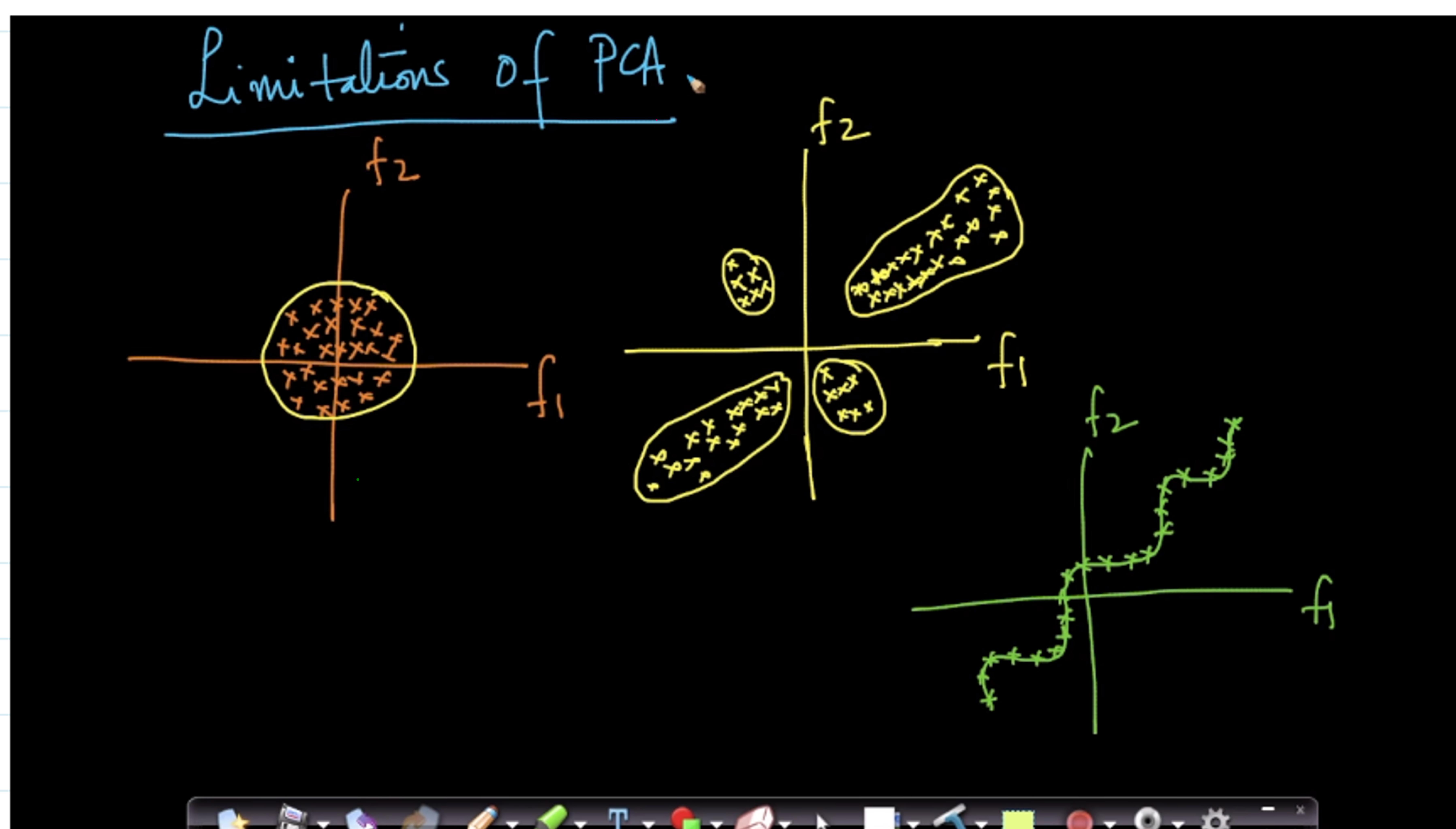
$$X' = \begin{bmatrix} 1 & 2 & 3 & \dots & 50 \\ 2 & 3 & 4 & \dots & 51 \\ 3 & 4 & 5 & \dots & 52 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ n & n+1 & n+2 & \dots & n+49 \end{bmatrix}$$
$$x_{ij}' = x_i^T v_j$$

→ d features to d' features where d' < d

→ Sometimes it's asked such that PCA needs to preserve 99% of the variance & the original number of features = 100.

Then we look for a value like let $\frac{\lambda_1 + \lambda_2 + \dots + \lambda_{51}}{\sum_{i=1}^{100} \lambda_i} = 0.99 \Rightarrow d' = 51$. It could be anything

→ sklearn inverse_transform can be used to convert d' to d dimensions.



→ In all these cases some information is lost.

Advantages :-

- ① Removes correlated features by reducing dimensions
- ② Improves algorithm performance
- ③ Reduces overfitting which mainly occurs when there are too many variables
- ④ Improves visualization

Disadvantages :-

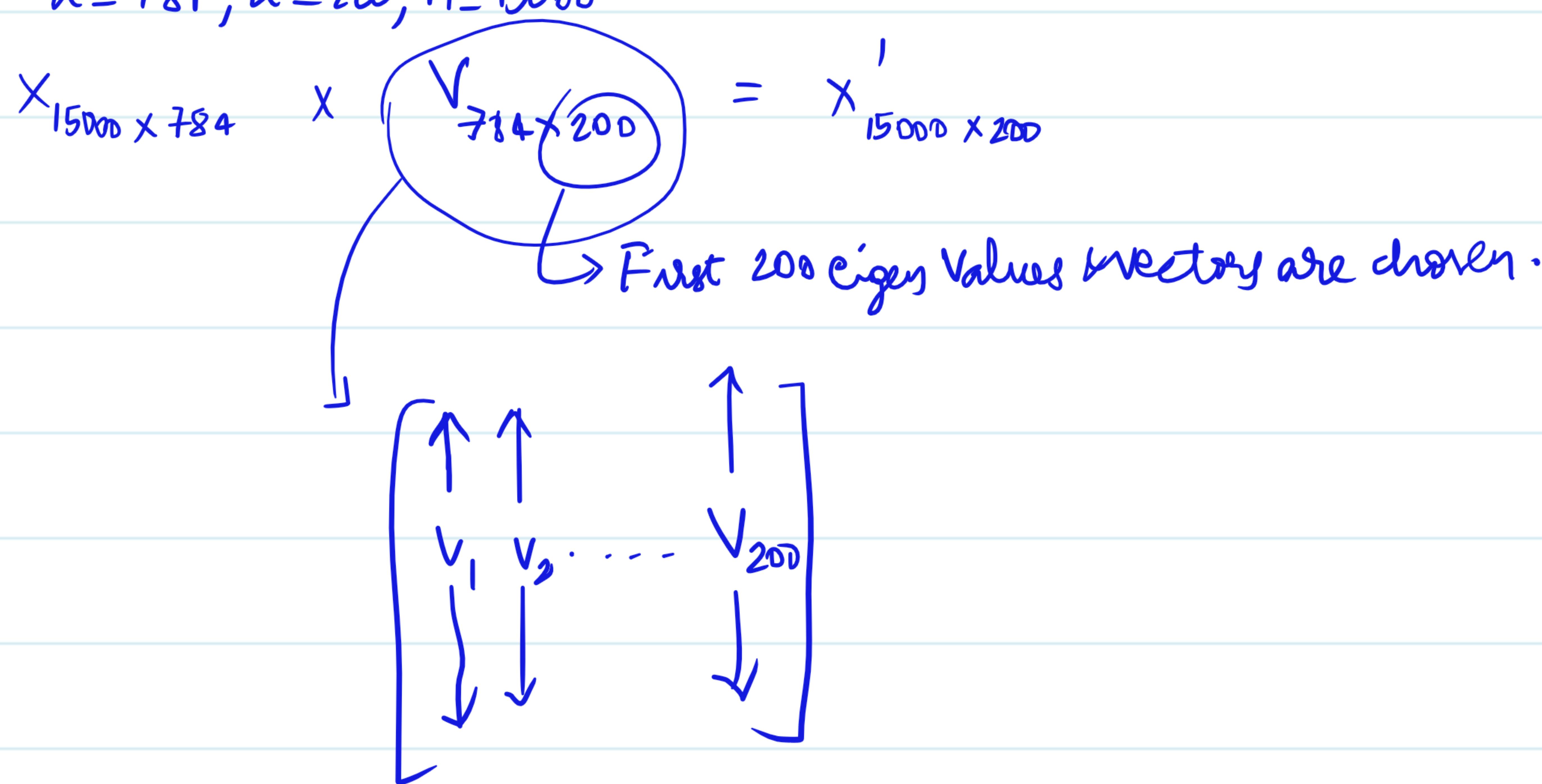
- ① Independent variables become less interpretable. They are converted to PCs.
- ② Data standardization is must before performing PCA.
- ③ Information loss- If number of PCs are not selected with care.

PCA for dimensionality Reduction (Not visualization) :-

$$d \xrightarrow{\text{PCA}} d'$$

↪ d when it's 2/3... it can be used for visualization sometimes, we just need to reduce the dimensions & visualization is not the main purpose.

ex:- $d=784, d'=200, n=15000$



→ $\frac{\lambda_1 + \lambda_2 + \dots + \lambda_{d'}}{\sum_{i=1}^d \lambda_i}$ gives the percentage of variance explained in d' dimensions. This helps us in picking the correct value of d' .

Generally we want to find 90% retaining d' .

→ It's not always guaranteed that applying ML model after PCA shows better results than just model. So performing both & using performance metrics helps pick which one is better.