→ Pick eigen vectors corresponding to largest eigen value of the covariance matrix $\Sigma$.

→ reduce to one dimension (PCA)

→ To reduce to two dimensions → pick ~~second largest vat~~ eigen vector corresponding to second largest eigen value.

→ We implement PCA in order to find best set of directions such that variability of data is maximized in the reduced dimensional space.

i·e $\quad x' = Xp$

↓

new matrix after dimensional reduction. (PCA)

$$\max \| x' \|^2 = \max \| Xp \|^2 \qquad\qquad p^T p = 1$$

$$= \max \, tr\left((Xp)^T (Xp)\right) \qquad\quad \text{such that } p^T p = 1$$

$$= \max \left( p^T x^T x p \right) \qquad\qquad \text{such that } p^T p = 1$$

$$= \max \left( p^T S p \right) \qquad\qquad\quad \text{such that } p^T p = 1$$

where P is the direction.

S = Scatter matrix

(Un normalized Covariance matrix)

## Using Lagrangian multipliers:-

$$\max \left( p^T S p \right) \quad s.T \quad p^T p = 1$$

$$\rightarrow \max \left( p^T S p \right) - \lambda \left( p^T p - 1 \right) \qquad (\lambda > 0)$$

This is used to solve constrained optimization by introducing $\lambda$ to combine $p^T S p$ and the constraint

Then solve for stationary point:

$$\frac{\partial}{\partial \lambda} L(P, \lambda) = 0 \qquad \text{where } L(P, \lambda)$$

$$= P^T S P - \lambda (P^T P - 1)$$

$$\rightarrow P^T P = 1$$

$$\frac{\partial L}{\partial P} \rightarrow SP = \lambda P$$

$$\therefore L(P, \lambda) = P^T \lambda P - \lambda (1 - 1) = \lambda P^T P = \lambda$$

∴ optimization reduces to max $\lambda$ such that $\lambda P = SP$

This mean eigen vector of p of s is found corresponding to maximum eigen value $\lambda$.

To find k directions, find eigen vector corresponding to each of k eigenvalue sorted in descending order.

     i-e take $P_1$ & $P_2$ which $SP = \lambda P$

        $P_1$ is the eigen vector corresponding to highest $\lambda$.

∴ To maximize value, $P_2$ should correspond to second highest eigen vector.

$$\lambda = \begin{bmatrix} \lambda_1 & 0 & \cdots & \\ 0 & \lambda_2 & - - & - \\ 0 & 0 & \lambda_3 & - - - \\ 0 & 0 & 0 & - - \lambda \end{bmatrix}$$

②

Let $x_i \in R^n$ be a sample

& let $X$ be the matrix of all samples $x_i$.

$\Rightarrow X \in R^{m \times n}$

## Dimensionality reduction using PCA:-

On applying PCA to $X$, we obtain a matrix of principle components, $V \in R^{n \times d}$ (assuming $d = rank(x)$)

$\rightarrow$ Reconstructed matrix $\boxed{x' = XVV^T}$

Loss function $J = \frac{1}{m} \sum\limits_{i=1}^{m}$ (original matrix $-$ reconstructed matrix)$^2$ $\quad \forall x_i \in X$

$= \frac{1}{m} \sum\limits_{i=1}^{m} (x_i' - x_i)(x_i' - x_i)^T$

$= \frac{1}{m} \sum\limits_{i=1}^{m} (x_i VV^T - x_i)(x_i VV^T - x_i)^T$

## Minimising the objective function $J$ using gradient descent:-

$J$ is a function of $V$ $\qquad (\because x' = XVV^T)$

$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_j)$

Here $\theta_j \equiv V$

$V \leftarrow V - \alpha \frac{\partial}{\partial V} \left[ \frac{1}{m} \sum\limits_{i=1}^{m} (x_i VV^T - x_i)(x_i VV^T - x_i)^T \right]$

$\frac{dJ(V)}{dV} = \frac{d}{dV} \left( \frac{1}{m} \sum\limits_{i=1}^{m} (x_i VV^T - x_i)(x_i VV^T - x_i)^T \right)$

$= \frac{1}{m} \sum\limits_{j=1}^{m} \left[ 2 (x_i VV^T - x_i)^T (2 x_i V) \right]$

$= \frac{4}{m} \sum\limits_{q=1}^{m} \left[ x_i (VV^T - I) \right]^T (x_i V) = \frac{4}{m} \sum\limits_{i=1}^{m} (VV^T - I)(x_i^T x_i) V$

Hence, update function

$$v \leftarrow V - \frac{4\alpha}{m} \sum_{i=1}^{m} (Vv^T - I)(x_i^T x) v$$

Hence, update function

$$v \leftarrow v - \frac{4\alpha}{m} \sum_{i=1}^{m} (vv^T - I)(x_i^T x)v$$

3. In case of PCA

Loss function = reconstruction error.

each sample $x_i \in \mathbb{R}^n$

Given sample data $X \in \mathbb{R}^{m \times n}$

Reconstructed data $x' = xvv^T$

where $v \in \mathbb{R}^{T \times d}$ is matrix of principal components.

$$RSS = J = \frac{1}{m} \sum (error)^2 \quad \forall \text{ samples.}$$

$$= \frac{1}{m} \sum_{i=1}^{m} (x_i' - x_i)(x_i' - x_i)^T \quad \forall x_i \in X$$

$$\Rightarrow J(v) = \frac{1}{m} \sum_{i=1}^{m} (x_i vv^T - x_i)(x_i vv^T - x_i)^T \quad \forall x_i \in X$$

$$\frac{dJ(v)}{dv} = \frac{4}{m} \sum_{i=1}^{m} (vv^T - I)(x_i^T x)v$$

(i) For $L_1$ regularisation (Lasso regression) $\Rightarrow$ RSS + $L_1$-norm

$$L^{lasso}(\lambda) = (J(v) + \lambda \|v\|_1)$$

$$\Rightarrow \min \left( \frac{1}{m} \sum_{i=1}^{m} [x_i vv^T - x_i][x_i vv^T - x_i]^T + \lambda \|v\|_1 \right)$$

$$\Rightarrow \frac{\partial}{\partial v} L^{lasso}(\lambda) = 0$$

$$\Rightarrow \frac{4}{m} \sum_{i=1}^{m} (vv^T - I)(x_i x_i^T)v + \lambda (sign(v)) = 0$$

(ii) For $L_2$ regularisation (Ridge regression) $\rightarrow$ RSS + L2 $-$norm.

$$L^{Ridge}(\lambda) = (J(W) + \lambda \|V_{*}\|_2)$$

$$\rightarrow min \left( \frac{1}{m} \sum_{i=1}^{m} [x_i VV^T - x_i][x_i VV^T - x_i]^T + \lambda \|V\|_2 \right)$$

$$\rightarrow \frac{\partial}{\partial V} L^{Ridge}(\lambda) = 0$$

$$\rightarrow \frac{4}{m} \sum_{i=1}^{m} (VV^T - I)(x_i^T x_i) V + \lambda \frac{V}{\|V\|_2} = 0$$