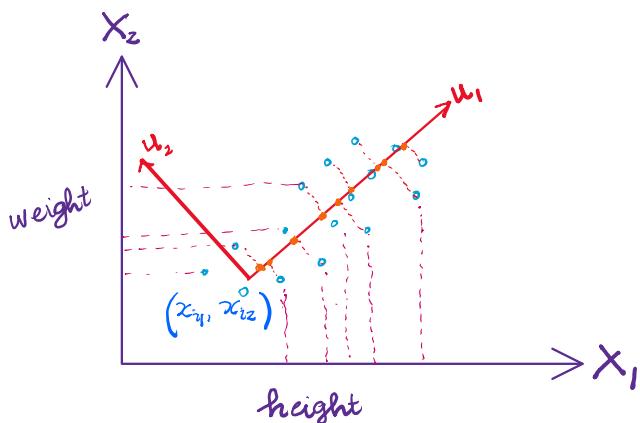


Q: Can we reduce the data without compromising the information?

Two methods for dimension reduction:

- Unsupervised:  $\{x_i\} \rightarrow$  Principal Component Analysis
- Supervised:  $\{(x_i, y_i)\} \rightarrow$  LDA.

## PCA:



Note:  $u_1$  is a better projection direction than  $u_2$  for this data.

PCA in action:

- Data visualization
- Feature extraction.

### Setup:

Consider a Dataset,  $D = \{x_1, x_2, \dots, x_n\}$ ,  $x_i \in \mathbb{R}^d$

Goal: To map these data on a dimension  $m < d$ .

The Question becomes: Which  $u_i$  to choose (see fig above)?

Objective: maximize the variance of projected data.

For a point  $x_i$ ,  $\|x_i\|$  is the length of  $\vec{x}_i = (x_{i1}, x_{i2})$  from origin.

$\therefore$  Projection of  $x$  on  $u_i$  direction:

$$\text{proj}_{u_i}(x_i) = (x_i^T u_i) u_i$$

Now to calculate the variance of projected data we first need mean:

$$\text{Mean of projected data} = \frac{1}{n} \sum_{i=1}^n \text{proj}_{u_i}(x_i) = (u_i^T \bar{x}) u_i.$$

(due to linearity)

$$\text{Mean of projected data} = \frac{1}{n} \sum_{i=1}^n \text{proj}_{u_i}(x_i) = (u_i^\top x_i) u_i$$

(due to linearity)

$$\begin{aligned}\text{Now, Variance of the projections} &= \frac{1}{n} \sum_{i=1}^n (\bar{u}_i^\top x_i - \bar{u}_i \bar{x})^2 \\ &= \frac{1}{n} \sum_{i=1}^n (u_i^\top (x_i - \bar{x}))^2 \\ &= \frac{1}{n} \sum_{i=1}^n u_i^\top (x_i - \bar{x})(x_i - \bar{x})^\top u_i \\ &= u_i^\top \left( \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top \right) u_i \\ &= u_i^\top S u_i\end{aligned}$$

covariance matrix

### Aside: (Covariance Matrix)

$$\text{Let } x_1 = \begin{pmatrix} x_{11} \\ x_{12} \end{pmatrix}, x_2 = \begin{pmatrix} x_{21} \\ x_{22} \end{pmatrix} \text{ and } \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$$

$$X = [x_1 - \mu \quad x_2 - \mu \quad \dots \quad x_n - \mu]$$

$$\frac{1}{n} X X^\top = \text{covariance matrix of data points} = S$$

### PCA as optimization Problem:

$$\text{maximize } u_i^\top S u_i, \text{ subject to constraint } u_i^\top u_i = 1$$

Notes:

1.  $\nabla^2 f(u_i)$  is Positive Semi-definite.
2.  $S$  is a real symmetric matrix
3.  $y^\top X X^\top y = \|X^\top y\|^2 \geq 0$

Though  $f(u_i) = u_i^\top S u_i$  is convex, we try to maximize it, instead of minimizing!  
So this not a convex optimization problem.

Result: For a real symmetric matrix  $A$ ,  $A$  has d-orthonormal eigenvectors.

$S$  is real symmetric. Let  $v_1, v_2, \dots, v_d$  be its orthonormal eigenvectors with eigenvalues  $\lambda_1, \lambda_2, \dots, \lambda_d$ . We know  $v_i^T v_j = \mathbb{1}_{\{i=j\}}$

We've:  $Sv_i = \lambda_i v_i$        $V = \begin{bmatrix} v_1 & \cdots & v_d \end{bmatrix}_{d \times d}$

$$\begin{aligned} SV &= [Sv_1 \ S v_2 \ \dots \ S v_d] \\ &= [\lambda_1 v_1 \ \lambda_2 v_2 \ \dots \ \lambda_d v_d] \\ &= V \begin{bmatrix} \lambda_1 & 0 & 0 & \cdots \\ 0 & \lambda_2 & 0 & \cdots \\ \vdots & \vdots & \ddots & \lambda_d \end{bmatrix}_{d \times d} \end{aligned}$$

Let  $\Sigma = \text{diag}(\lambda_1, \dots, \lambda_d) \therefore S = V\Sigma \quad \textcircled{1}$

Observe that  $V^T V = I$  as  $[V^T V]_{ij} = v_i^T v_j$   
(unitary matrix)

Now,  $SVV^T = V\Sigma V^T \xrightarrow{\text{diagonalization!}}$   
But  $VV^T = I \Rightarrow S = V\Sigma V^T$   
(eigen value decomposition)

Since  $\{v_1, \dots, v_d\}$  form an orthonormal basis for  $\mathbb{R}^d$ , we can write

$$u_1 = \sum_{j=1}^d x_j v_j = Vx \quad \text{where } x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

Returning to our objective fn:

$$\begin{aligned} f(u_1) &= u_1^T S u_1 = x^T \underbrace{V^T V}_{I} \Sigma \underbrace{V^T V}_{I} x \\ &= x^T \Sigma x = \sum_{j=1}^d \lambda_j x_j^2 \end{aligned}$$

Thus our new optimization problem:

$$\max \sum_{j=1}^d \lambda_j x_j^2 \text{ st. } \sum_{j=1}^d x_j^2 = 1 \quad (\text{convert to linear programm})$$

This can be easily solved: Put  $x_k = 1$  where  $k = \arg \max \lambda_i$  and  $x_j = 0 \forall j \neq k$

This can be easily solved: Put  $x_k = 1$  where  $k = \operatorname{argmax}_i \lambda_i$  and  $x_j = 0 \forall j \neq k$

$\therefore$  First principal component,  $u_1 = v_k$  {Usually we arrange in decreasing  $\lambda$ 's  
(largest eigen vector) then  $v_k = u_1 = u_1$ }

Generalization for  $d \rightarrow m$  reduction:

After we've found first Principal component,  $u_1$ , we can find  $u_2$

$$\begin{aligned} & \max u_2^T S u_2 \\ & \text{s.t. } u_2^T u_2 = 1, \quad u_1^T u_2 = 0 \end{aligned}$$

$\Rightarrow u_2 = \text{second largest eigenvectors.}$

For  $m$  Principal components, choose the largest  $m$  eigenvectors.

PCA algorithm:

1. Compute mean of data points,  $\bar{x}$

2. Mean centre the data,  $\tilde{x}$ , compute  $x_i - \bar{x}, \forall i$

3. Compute the covariance matrix,  $S = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$

4. Do eigenvalue decomposition of  $S = V \Sigma V^T$

5. Pick  $m$  top (largest) eigenvectors:  $u_1, \dots, u_m \rightarrow \text{Principal Components}$

6. Create a projection matrix,  $U = [u_1, \dots, u_m]$

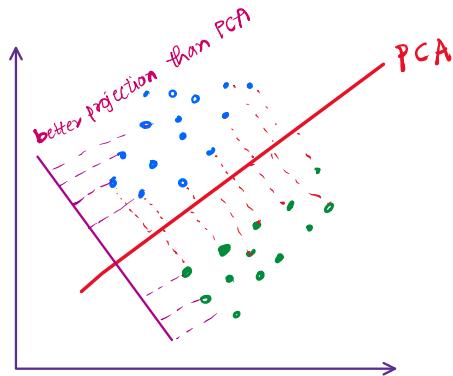
7. For each  $x_i \in D$ ,  $U^T x_i$  will give the projection on  $m$ -dimensional space

$$U^T x = \begin{bmatrix} u_1^T x \\ u_2^T x \\ \vdots \\ u_m^T x \end{bmatrix}$$

## Supervised Learning:

If we use PCA, it fails to retain class information.

We need a projection that separates the classes well



We need a good measure to evaluate these projections (variance not enough)

$$J(u) = |\mu_1^T u - \mu_2^T u| \quad \mu_i: \text{Mean of class } i.$$

One possible measure,  
but this alone is not a good measure.

- Objective:
- The different class means are well separated
  - The data in the same class are not well separated.