Centers that has the best loss.

unsupervised learning $\begin{cases} \rightarrow \text{clustering} \\ \rightarrow \text{dimensionality reduction} \\ \qquad PCA \end{cases}$

$\{ x_1, x_2, --, x_n \}$

image

P: # pixels

$IR^P \rightarrow$

$R^2$
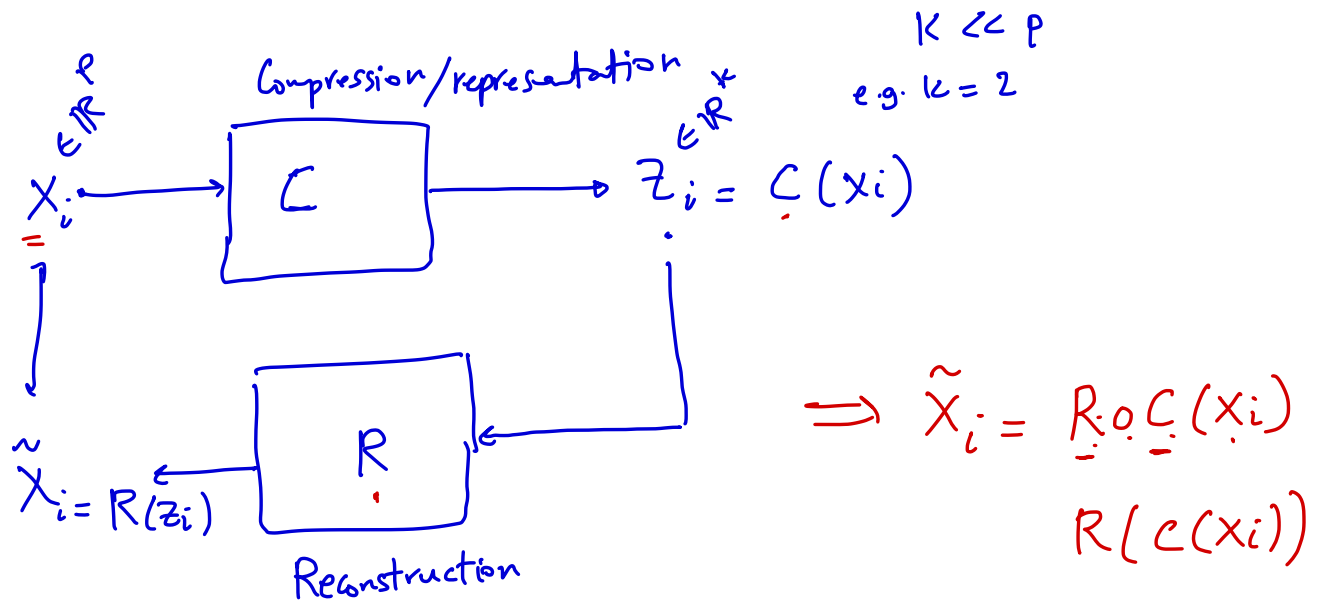
C

- $X_1, X_2, \ldots, X_n \in \mathbb{R}^p$   (P is typically large)

- we'd like to represent the data in a low-dimensional space (e.g. $\mathbb{R}^2$)

$$K \ll p \quad \text{e.g. } k = 2$$

$$X_i \xrightarrow{\quad} \boxed{C} \xrightarrow{\quad} Z_i = C(X_i)$$

Compression/representation

$$\tilde{X}_i = R(Z_i) \xleftarrow{\quad} \boxed{R} \xleftarrow{\quad}$$

Reconstruction

$$\implies \tilde{X}_i = R \circ C(X_i)$$
$$R(C(X_i))$$

Our goal is to design $C, R$ such that $X_i$ and $\tilde{X}_i$ are as close as possible.

objective:

$$\underset{C, R}{\text{minimize}} \quad \sum_{i=1}^{n} \| X_i - \tilde{X}_i \|_2^2$$

$$= \underset{C, R}{\text{minimize}} \quad \sum_{i=1}^{n} \| X_i - R \circ C(X_i) \|_2^2$$

$$x_i = \begin{pmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{pmatrix}^{\in \mathbb{R}^P}$$



$x_i \longrightarrow$ | P { | C | k } | $z_i = C(x_i)$



$x_i$ $\xrightarrow{\quad C \quad}$ $\begin{pmatrix} 0 \\ 0 \end{pmatrix} = z_i$

As usual, let's start with the simplest class of mappings which are linear mappings.

$$z_i = C(x_i) = \underset{\substack{\uparrow \\ \text{matrix } C = [\quad]_{K \times}}}{C} x_i$$

e.g. $k = 2 \longrightarrow C = \begin{bmatrix} \underline{\quad\quad\quad} \\ \underline{\quad\quad\quad} \end{bmatrix}_{2 \times P}$

$$= \begin{bmatrix} \dfrac{U_1^T}{U_2^T} \end{bmatrix}$$

$$\implies z_i = C x_i = \begin{bmatrix} U_1^T x_i \\ U_2^T x_i \end{bmatrix}$$

Also, $R = $ linear mapping

$$\tilde{x}_i = R z_i = \underset{\in \mathbb{R}_{p \times k}}{R} \cdot \underset{\in \mathbb{R}_{k \times p}}{C} \cdot x_i$$

$$[ \quad ]_{p \times k}$$

objective: $\underset{\substack{R, C \\ \in \mathbb{R}^{p \times k} \quad \in \mathbb{R}^{k \times p}}}{\min} \quad \sum_{i=1}^{n} \| x_i - R \cdot C \cdot x_i \|_2^2 \qquad (PCA)$

$C : \mathbb{R}^p \longrightarrow \mathbb{R}^k$

$C = (\text{matrix})_{k \times p}$

<span style="color:red">This problem is called the Principal Component Analysis (PCA)</span>

In order to solve this problem, we need to review an important tool in linear algebra:
( this tool is important for many other branches of Data science )

# The Singular Value Decomposition:

**Theorem:** Every symmetric matrix, $A_{p \times p}$, can be written as

$$A_{p \times p} = U^T_{p \times p} \; \Lambda_{p \times p} \; U_{p \times p}$$

- Where $UU^T = I_{p \times p}$ &larr; identity matrix

$$\Lambda = \text{diagonal}: \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & 0 \\ & & \ddots & \\ 0 & & & \lambda_n \end{pmatrix}$$

- $A = U^T \cdot \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & 0 \\ & & \ddots & \\ 0 & & & \lambda_n \end{pmatrix} U$

- $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \cdots \geq \lambda_n$

- $\lambda_i = $ eigenvalues of $A$.

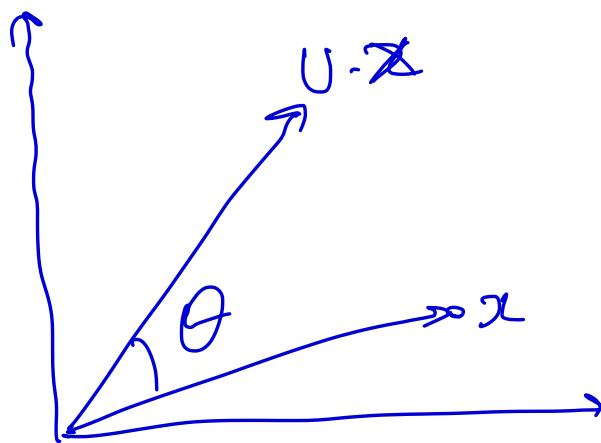- Each row of $U$ is an "eigenvector" of the matrix $A$.
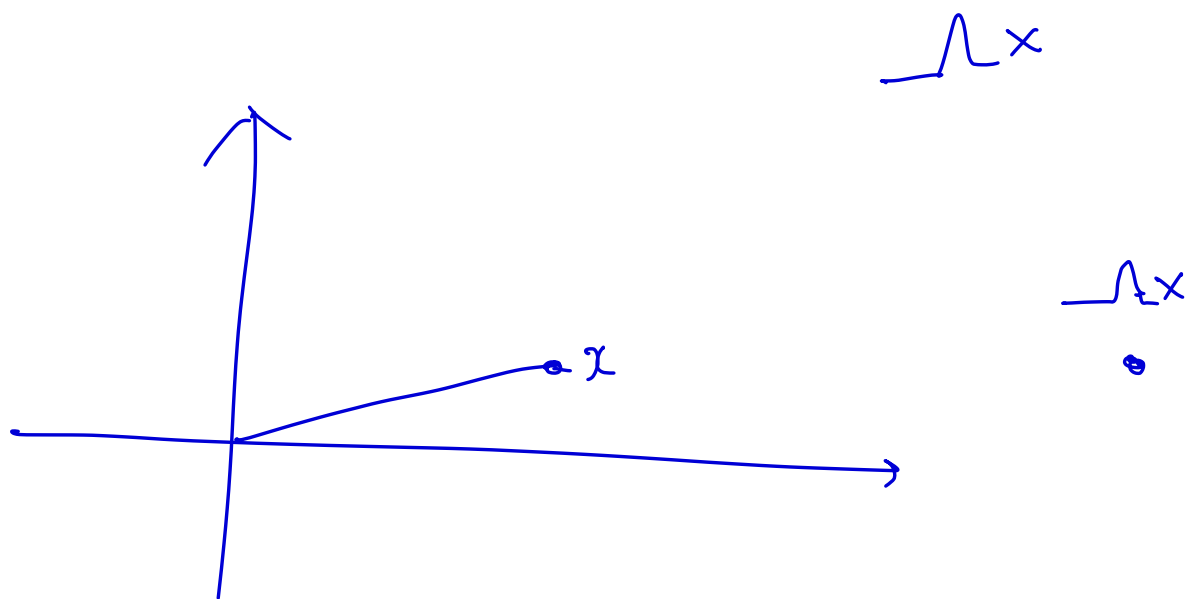
Example:

Assume A is 2x2 matrix.

$$A = \begin{pmatrix} a & c \\ c & b \end{pmatrix}$$

rescaling

$$A = U \cdot \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \cdot V^T$$

Rotation($\theta$)        Rotation($-\theta$)

$$U \cdot U^T = I \longrightarrow U = \begin{pmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{pmatrix}$$
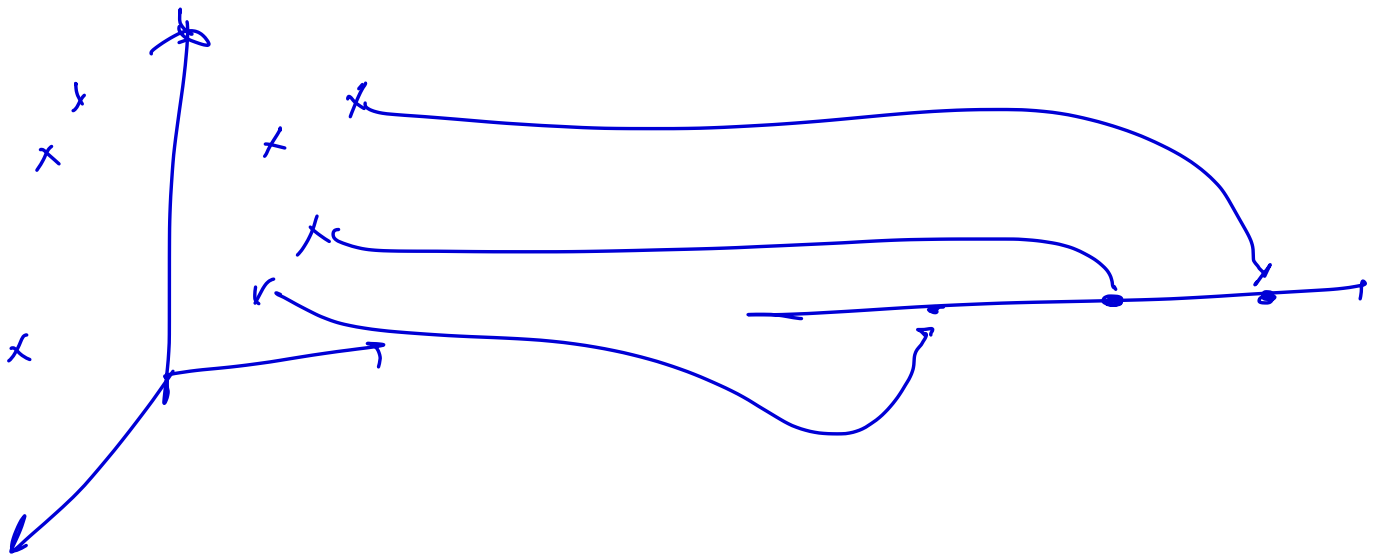
$$\Lambda = \begin{pmatrix} 4 = \lambda_1 \geq 0 & 0 \\ 0 & 1 = \lambda_2 \geq 0 \end{pmatrix} \text{re Scaling}$$

$\Lambda x$

$\Lambda x$



$$A = \text{Rotation}(-\theta) \cdot \text{Rescaling} \cdot \text{Rotation}(\theta)$$

let's go back to our problem.

let's assume for simplicity that $k=1$ ; we're looking for the best 1-dimensional representation of the data.
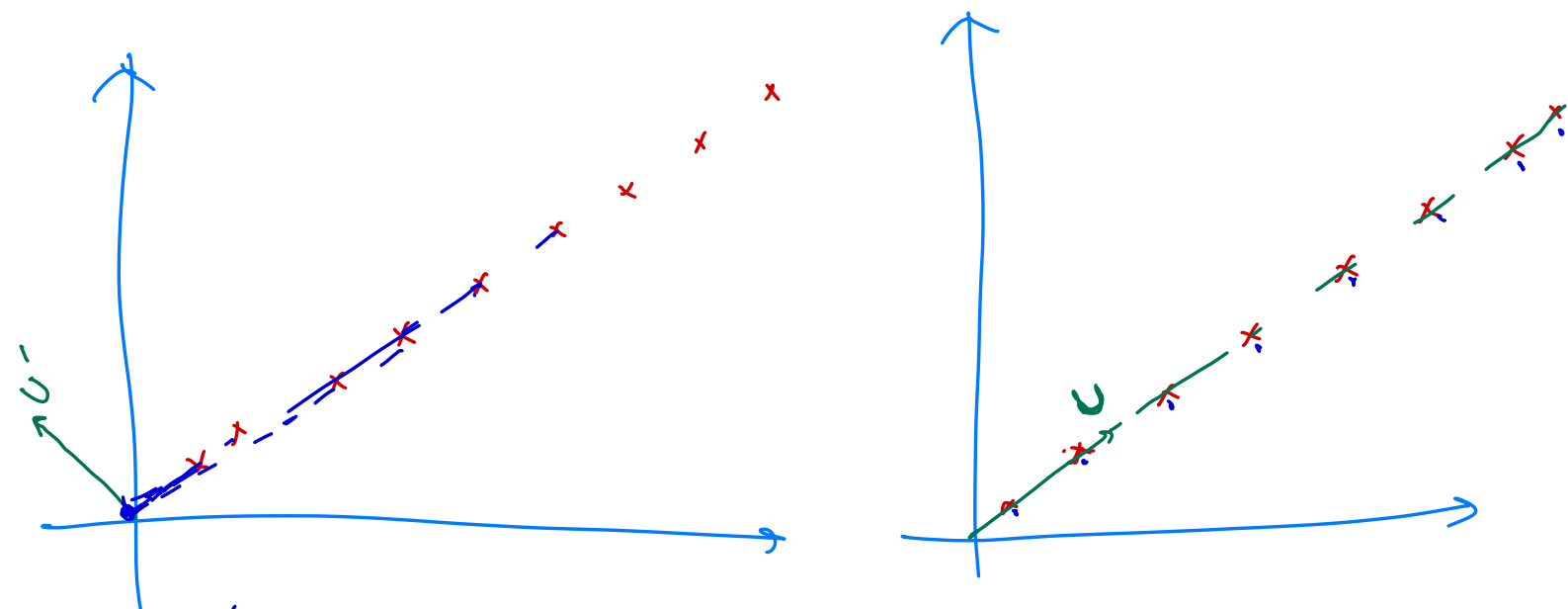
$$C = \begin{bmatrix} & \end{bmatrix}_{k \times p} \longrightarrow C = U^T$$

$$C \cdot X_i = \underset{\underset{\text{Vector}}{=\downarrow}}{U^T X_i}$$

When $k=1$, our goal is to find a single vector $U^T$ such that the representation $z_i = U^T X_i$ is "the most informative" representation.
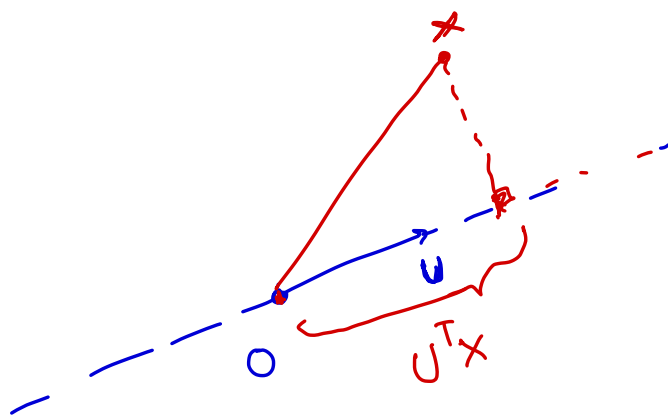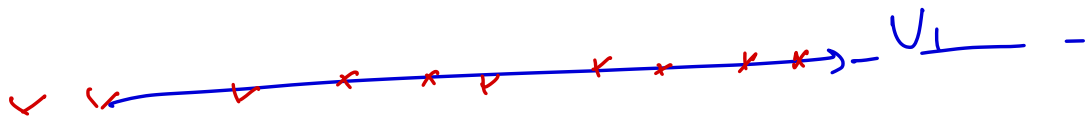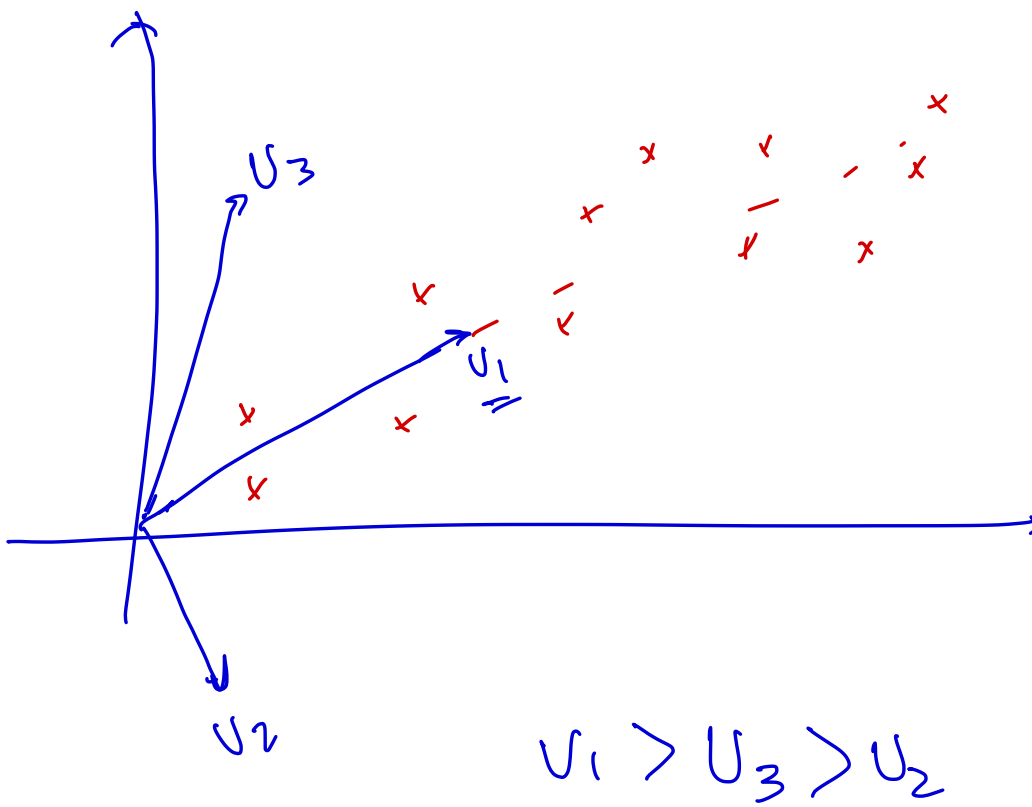
Example:



$U'^T X_i \tilde{\approx} 0$

$U^T X_i$

U is a better direction than U' in terms of Capture more information about the data

Assume $\|U\|_2 = 1$

$U^T X$

O

$U_3$

$U_1$

$U_2$

$$U_1 > U_3 > U_2$$

$U_1$

$U_2$

$U_3$

We're looking for a direction U along which data has the most "variation".

- Mathematically, we're looking for a direction $u$ such the "variance" of $u^T X$ is the most.

- Assume that the data comes from a distribution $X \sim P(x)$. Also assume for simplicity that $E[X] = 0$.

$$\max_{\substack{u^T \in \mathbb{R}^{1 \times p} \\ \|u\|_2 = 1}} Var(\underline{u}^T X)$$

$$= \max_{u : \|u\|_2 = 1} E[(u^T X)^2]$$

$$= \max_{: \|u\|_2 = 1} E[u^T X^T X u]$$

$$= \max_{: \|u\|_2 = 1} u^T \underbrace{E[XX^T]}_{\substack{\text{covariance} \\ \text{matrix of} \\ \text{the data}}} u$$

$$\left( \begin{array}{l} X \in \mathbb{R}^P \\ E[XX^T] \\ = \Sigma_{P \times P} \end{array} \right)$$

Let $\Sigma \triangleq E[xx^T]$

$$= \max_{\substack{u: \\ \|u\|_2 = 1}} u^T \Sigma u$$

The solution of the problem above can be found from the singular value decomposition of $\Sigma$.

$$\Lambda = \begin{pmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & \cdots & \\ & & & \lambda_p \\ 0 & & & \end{pmatrix}$$

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p$$

$$\text{SVD} \longrightarrow \underline{\Sigma} = \underline{U}^T \underline{\Lambda} U$$

It can be shown the the first row of the matrix $U$, is the maximizer of the following problem:

$$\max_{\substack{u: \\ \|u\|=}} u^T \Sigma u$$

$$U = \begin{pmatrix} \dfrac{U_1^T}{} \\ \dfrac{U_2^T}{} \\ \vdots \\ \overline{U_p^T} \end{pmatrix}_{p \times p}$$

- $U_r^T$ is the solution to the problem above.

- $U_1^T$ is the principal component of the data

- $v_i^T$ is the vector along which the data has the most "variation".

---

$$\Sigma = E\left[ X X^T \right]$$

- We know that the data

$$X_1, X_2, \ldots, X_n \overset{iid}{\sim} \text{distribution} \\ X \sim P(x)$$

in practice:
- We cannot compute the matrix $\Sigma$ because we do not know the distribution of the data

$$\Sigma \approx \frac{1}{n} \sum_{i=1}^{n} X_i X_i^T$$

$\underbrace{\qquad}_{\text{unbiased estimator}}$

# Algorithm ( PCA ) : $(k=1)$

- Input: $X_1, X_2, \ldots, X_n \in \mathbb{R}^p$

- Compute the empirical covariance matrix of data:
$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^{n} X_i X_i^T$$

- Compute the Singular value decomposition of $\hat{\Sigma}$ :
$$\hat{\Sigma} = U^T \Lambda U$$

- Use the first row of $U$, $u_1$, as the principal component of the data:
$$Z_i = U_q^T X_i .$$

---

$*$ let's now find the best 2-dimensional representation of the data.

$$\longrightarrow \quad \max_{u_1, u_2} \quad (= u_1^T \Sigma u_1 + u_2^T \Sigma u_2$$

$$\|u_1\|, \|u_2\| = 1$$

$$(\text{orthogonal}) \quad u_1^T u_2 = 0$$

It can be shown that $u_1, u_2$ are the first two rows of the matrix $U$.

* In general if we want to find a matrix $C = \begin{bmatrix} \frac{u_1^T}{u_2^T} \\ \vdots \\ u_k^T \end{bmatrix}_{k \times n}$ then

  · we need to solve the following optimization problem:

$$\sum_{i=1}^{k} u_i^T \Sigma u_i$$

<span style="color:red">Cumulative variance of the data along the directions $u_1, \dots, u_k$</span> ←

$$\text{s.t. } \|u_i\| = 1$$
$$\text{and } u_i^T u_j = 0 \text{ for } i \neq j$$

* it can be shown that the solution to the above optimization problem is the first $k$ rows of the matrix $U$.

# Algorithm PCA for general $k$:

- Input: $X_1, \dots, X_n \in \mathbb{R}^p$

- Compute: $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^{n} X_i X_i^T$

- Compute: (SVD) $\hat{\Sigma} = U^T \Lambda U$

- $C = \begin{bmatrix} U_1^T \\ U_2^T \\ \vdots \\ U_k^T \end{bmatrix}$ where $U_1, \dots, U_k$ are the first $k$ rows of the matrix $U$.

It can be proven that PCA gives the optimal linear representation of the data (it solves equation (PCA) exactly).