# UNSUPERVISED LEARNING

└ REPRESENTATION LEARNING.

**Goal:** Given a set of "data points", "understand" something "useful" about them.

Data points → Vectors in $\mathbb{R}^d$ $\begin{bmatrix} \text{height} \\ \text{weight} \\ \text{age} \end{bmatrix} \in \mathbb{R}^3$

Running theme: "COMPREHENSION IS COMPRESSION" (George Chaitin)

└→ understanding
└→ learning

**Problem:**

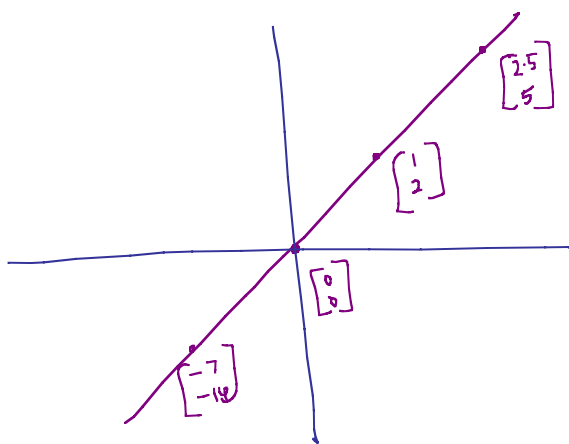Input: $\{x_1, x_2, \dots, x_n\}$     $x_i \in \mathbb{R}^d \leftarrow$ # of features

Output:   Some "Compressed" representation of the dataset

**Example:**
$$\left\{ \overset{x_1}{\begin{bmatrix} -7 \\ -14 \end{bmatrix}}, \overset{x_2}{\begin{bmatrix} 2.5 \\ 5 \end{bmatrix}}, \overset{x_3}{\begin{bmatrix} 0.5 \\ 1 \end{bmatrix}}, \overset{x_4}{\begin{bmatrix} 0 \\ 0 \end{bmatrix}} \right\}$$

**Question:** How many real numbers are needed to store this dataset  $\boxed{8}$

Representative

$\begin{bmatrix} 1 \\ 2 \end{bmatrix} \in \mathbb{R}^2$

Co-efficients

$\{-7, 2.5, 0.5, 0\}$   $\boxed{6}$

NOTE:
using representative
& co-effecients
can "RECONSTRUCT"
the dataset exactly.

**Rep:**

$\begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$\overline{\text{Co-eff}}$

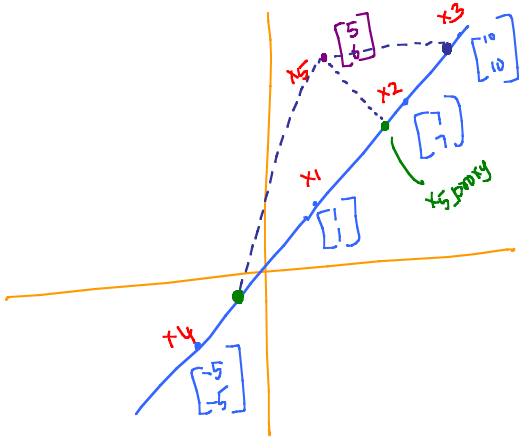$\{-7,\ 2.5,\ 0.5, 0\}$

$\begin{bmatrix} 1/\sqrt{5} \\ 2/\sqrt{5} \end{bmatrix}$

$\{-7\sqrt{5},\ 2.5\sqrt{5},\ 0.5\sqrt{5}, 0\}$

**NOTE:**

Any vector along the purple line can be chosen as a representative $\left[\text{except } \begin{bmatrix} 0 \\ 0 \end{bmatrix}\right]$

**Original:** # real numbers $= dn$

# real numbers in compressed representation $= d + n$

Rep

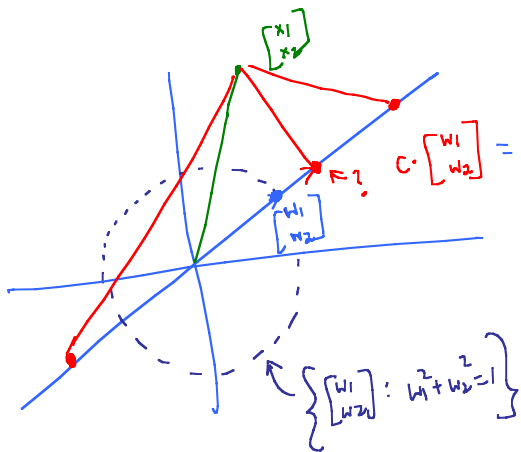$$\left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix} \right\}$$

Co-eff

$$\left\{ (-5,-5), (1,1), (7,7), (5,5), \atop (10,10) \right\}$$

$$x_3 = 10 \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} + 10 \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 10 \\ 10 \end{bmatrix}$$

Q: Who can "pretend" to be a proxy for $x_5$ along the blue line?

Ans: Projection of $x_5$ onto the blue line.

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

R?

$$c \cdot \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \left( \frac{x_1 w_1 + x_2 w_2}{w_1^2 + w_2^2} \right) \cdot \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$= \left( \frac{x^T w}{\|w\|^2} \right) \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

$$\left\{ \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} : w_1^2 + w_2^2 = 1 \right\}$$

$$\min_c \quad \text{length}^2 \left( \underline{\text{error vector}} \right)$$

$$\hookrightarrow \quad \begin{bmatrix} x_1 - c\, w_1 \\ x_2 - c\, w_2 \end{bmatrix}$$

$$\min_c \quad (x_1 - c w_1)^2 + (x_2 - c w_2)^2$$

$$c^* = \left( \frac{x_1 w_1 + x_2 w_2}{w_1^2 + w_2^2} \right) \qquad \text{(scalar)}$$

Inner product / dot product of $x$ and $w$.

$$\text{length}^2 \left( \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \right) \leftarrow$$

NOTE: Can always pick $\begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$ s.t

$$\text{length}\left( \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \right) = 1$$

$$\Rightarrow \quad c^* = \left( x^T w \right) \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$$

**Goal:** Develop a way to find a "compressed" representation of data when data points not-necessarily fall on line

**Goal:** Find the line that has the least "reconstruction" error.

Dataset : $\{x_1, x_2, \dots, x_n\}$    $x_i \in \mathbb{R}^d$

$$\|z\|^2 = z^T z$$

$$\text{ERROR}\left(\text{line, dataset}\right) = \sum_{i=1}^{n} \text{error}\left(\text{line, } x_i\right)$$

Represented using $w$
s.t $\|w\|^2 = 1$

$$= \sum_{i=1}^{n} \text{length}^2\left(x_i - (x_i^T w)w\right)$$

$$= \sum_{i=1}^{n} \|x_i - (x_i^T w)\cdot w\|^2$$

$$f(w) = \frac{1}{n}\sum_{i=1}^{n} \|\underbrace{x_i}_{\in \mathbb{R}^d} - \underbrace{(x_i^T w)}_{\in \mathbb{R}}\cdot \underbrace{w}_{\in \mathbb{R}^d}\|^2$$

$$= \frac{1}{n}\sum_{i=1}^{n} \left(x_i - (x_i^T w)\cdot w\right)^T \left(x_i - (x_i^T w)\cdot w\right)$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left[ x_i^T x_i - (x_i^T w)^2 - (x_i^T w)^2 + (x_i^T w)^2 \cdot 1 \right]$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left( \underbrace{x_i^T x_i} - (x_i^T w)^2 \right)$$

$$\min_{\substack{w: \\ \|w\|^2=1}} \quad g(w) = \frac{1}{n} \sum_{i=1}^{n} -(x_i^T w)^2$$

$$\max_{\substack{w: \\ \|w\|^2=1}} \quad \frac{1}{n} \sum_{i=1}^{n} (x_i^T w)^2 \quad = \quad \frac{1}{n} \sum_{i=1}^{n} \underset{1 \times d}{(w^T x_i)} \underset{d \times 1}{} \underset{1 \times d}{(x_i^T w)} \underset{d \times 1}{}$$

$$= \frac{1}{n} \sum_{i=1}^{n} \underline{w^T} \left( x_i x_i^T \right) \underline{w}$$

$$= w^T \left( \underbrace{\frac{1}{n} \sum_{i=1}^{n} x_i x_i^T}_{C \quad d \times d \quad \text{matrix}} \right) w$$
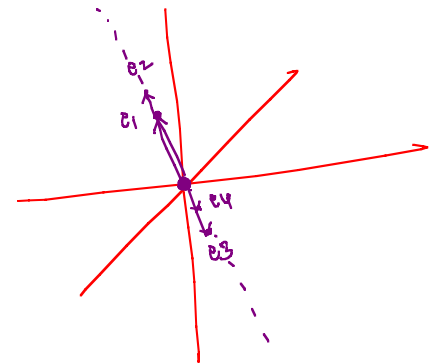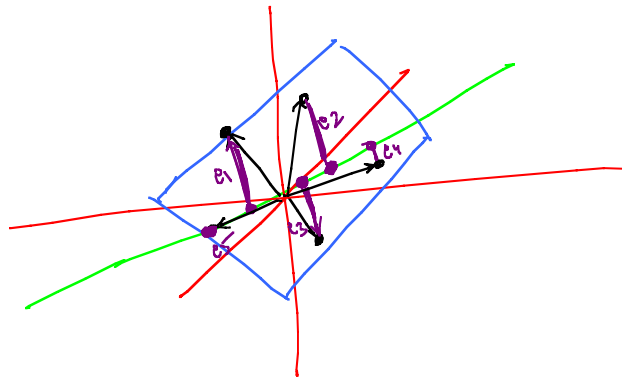
equivalently

$$\max_{\substack{w: \\ \|w\|^2 = 1}} w^T C w$$

$$C = \frac{1}{n} \sum_{i=1}^{n} x_i x_i^T$$

↳ Co-variance matrix

Soln: $w$ is the eigenvector corresponding to the maximum eigen value of C

$x \in \mathbb{R}^d$

$\downarrow$ Find $w$

$(x^T w) \cdot w$

$\downarrow$ Residue/error

$x - (x^T w) \cdot w$
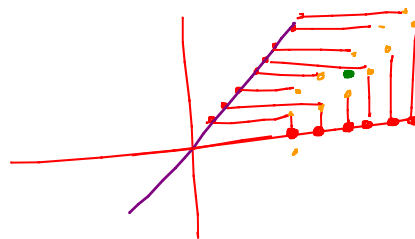
Might not be
error but
has "information"

Input: $\{ x_1, \cdots, x_n \}$   $x_i \in \mathbb{R}^d$

$\rightarrow$ $\mu = \frac{1}{n} \sum_{i=1}^{n} x_i$    $x_i = x_i - \mu$   $\forall i$

$\rightarrow$ Find "best" line  $w_1 \in \mathbb{R}^d$

$\rightarrow$ Replace  $x_i \leftarrow x_i - (x_i^T w) w$

$\rightarrow$ Repeat  to obtain  $w_2$



ISSUE:  Data  may not
         be  centered.

# Questions

→ How to solve $\max\limits_{\substack{w \\ \|w\|^2=1}} w^T C w$ ?

→ How many times to repeat the procedure?

→ Where exactly is "compression" is happening?

→ What "representations" are we learning?

$$D = \{x_1, x_2, \cdots, x_n\} \qquad x_i \in \mathbb{R}^d.$$

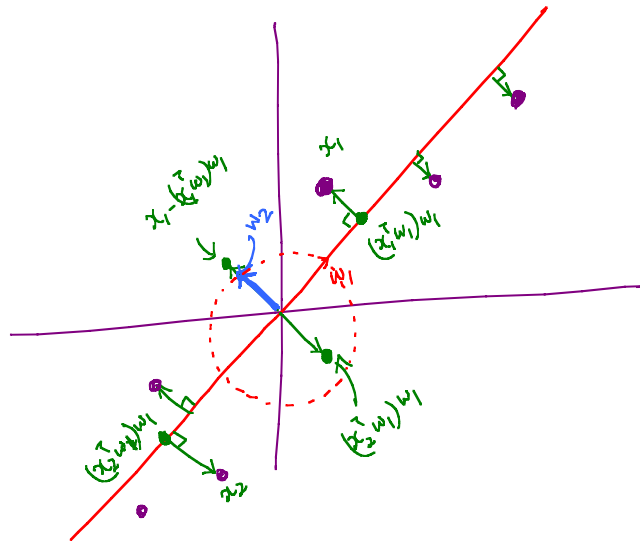$$W_1 = \underset{\substack{w: \\ \|w\|^2 = 1}}{\text{argmax}} \quad w^T C w \qquad\qquad C = \frac{1}{n} \sum_{i=1}^{n} x_i x_i^T$$

$$\boxed{\begin{array}{c} x_1' \\ x_2' \\ \vdots \\ x_n' \end{array}} \begin{array}{l} \rightarrow \quad x_1 - (x_1^T w_1) w_1 \\ \rightarrow \quad x_2 - (x_2^T v_1) w_1 \\ \\ \rightarrow \quad x_n - (x_n^T w_1) w_1 \end{array}$$

$$W_2 = \underset{\substack{w: \\ \|w\|_2^2 = 1}}{\text{argmax}} \quad w^T C' w$$

$$C' = \frac{1}{n} \sum_{i=1}^{n} x_i' x_i'^T$$

**Question:**  What can we say about $w_1$ & $w_2$?



$$\Rightarrow \boxed{w_2^T w_1 = 0}$$

**Observation:**

→ All residues are orthogonal to $w_1$

→ Any line which minimizes sum of errors w.r.t residues must also be orthogonal to $w_1$   [ARGUE WHY]

By continuing this procedure, we get $\{w_1, w_2, \ldots, w_d\}$ s.t $\|w_k\|_2^2 = 1 \quad \forall k$ and $w_i^T w_j = 0 \quad \forall i \neq j$

$\downarrow$

ORTHONORMAL

VECTORS

Residue after round 1

$\left\{ x_1 - (x_1^T w_1) w_1, \ldots, x_n - (x_n^T w_1) w_1 \right\}$

all vectors in $\in \mathbb{R}^d$

- $w_2 \rightarrow$ Best line that fits residues.

- $w_1^T w_2 = 0$

Residues    after   round 2

$$\left\{ \quad x_1 - (x_1^T \omega_1)\omega_1 \quad - \left( \left( x_1 - (x_1^T \omega_1)\omega_1 \right)^T \omega_2 \right) \omega_2 \quad , \quad \cdots \qquad \right\}$$

$$= \left\{ \quad x_1 - (x_1^T \omega_1)\omega_1 \quad - \left( x_1^T \omega_2 - \underbrace{(x_1^T \omega_1) \cdot \omega_1^T \omega_2}_{=0} \right)\omega_2 \quad , \quad \cdots \right\}$$

$$= \left\{ \quad x_1 - (x_1^T \omega_1)\omega_1 \quad - (x_1^T \omega_2)\omega_2 \quad , \quad \cdots \right\}$$

Residues    after   $d$-rounds.

$$\forall i \qquad x_i - \left( (x_i^T \omega_1)\omega_1 + (x_i^T \omega_2)\omega_2 + \cdots + (x_i^T \omega_d)\omega_d \right) \quad = \vec{0} \in \mathbb{R}^{d}.$$

$$\forall i \quad x_i = (z_i^T w_1) \underline{w_1} + (z_i^T w_2) \underline{w_2} + \cdots + (z_i^T w_d) \underline{w_d}.$$

## What have we gained?

— If data lives in a "low" dimensional linear sub-space, then residues become 0 much earlier than d rounds.

Example: Say Dataset is such that after 3-rounds, residues become 0.

$$\text{Dataset} = \{x_1, \cdots, x_n\} \qquad x_i \in \mathbb{R}^{100}$$

$$\forall i \quad x_i = (x_i^T w_1) \underline{w_1} + (x_i^T w_2) \underline{w_2} + (x_i^T w_3) \underline{w_3}$$

$\{w_1, w_2, w_3\}$
$\in \mathbb{R}^{100}$

Rep

$\{w_1, w_2, w_3\}$

↳ Common for
    dataset

Co-efficients

$x_i \rightarrow [x_i^T w_1 \quad x_i^T w_2 \quad x_i^T w_3] \in \mathbb{R}^3$

↳ Data point specific

Original :     $100 \times n$
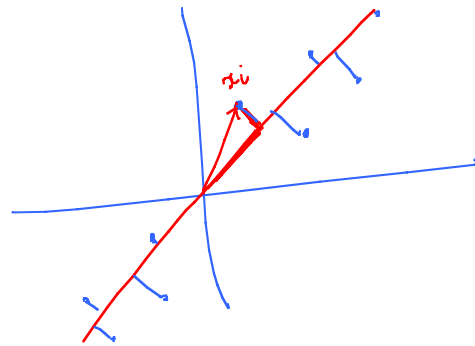
$100 \times 100 = 10000$

$\boxed{d \times n}$

now :    $3 \times 100 + 3n$

$3 \times 100 + 3 \times 100 = \underline{600}$

$\boxed{d \times k + k \times n}$

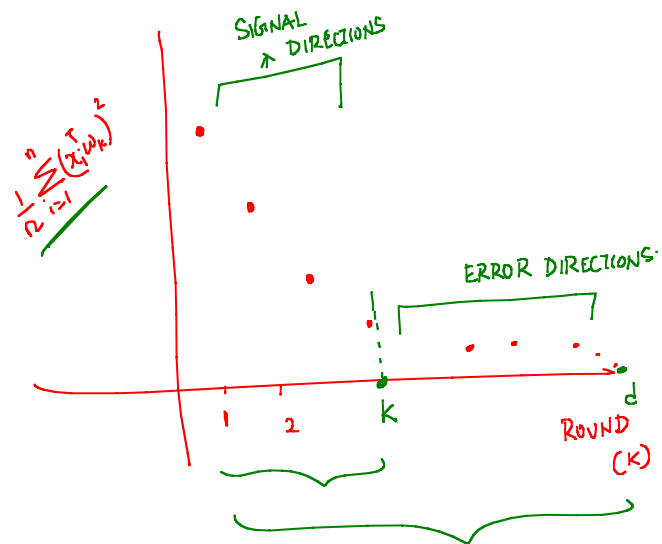Question:    What if data "approximately" lies in a low-dimensional space?

PYTHAGOROUS

For any $w \in \mathbb{R}^d$, s.t $\|w\|_2^2 = 1$

$$\forall i \quad \|x_i\|^2 = \|x_i - (x_i^T w) w\|^2 + \|(x_i^T w) w\|^2$$

$$\frac{1}{n} \sum_{i=1}^{n} \|x_i\|^2 = \frac{1}{n} \sum_{i=1}^{n} \|x_i - (x_i^T w) w\|^2 + \frac{1}{n} \sum_{i=1}^{n} \|(x_i^T w) w\|^2$$

As large as possible.

"Larger the value of $\frac{1}{n} \sum_{i=1}^{n} (x_i^T w)^2$, the better the fit."

$\frac{1}{2} \sum_{i=1}^{n} (z_i^T w_k)^2$

SIGNAL DIRECTIONS

ERROR DIRECTIONS.

1    2    k    d

ROUND (K)

ENTER   LINEAR ALGEBRA

$$\max_{w:\ \|w\|_2^2 = 1} w^T C w \qquad C = \frac{1}{n} \sum_{i=1}^{n} z_i z_i^T$$

$C \rightarrow$ covariance matrix.

<u>Soln</u>:  $w_1$  is  eigenvector  corresponding  to  the  "largest"  eigenvalue  of  C    [ HILBERT'S min-max Theorem ]

In fact  $\{ w_1, \ldots, w_d \}$ , the eigenvectors of C  form  an  orthonormal basis.

$w_R \rightarrow$  Best line  one can  obtain  in round k

What do eigenvalues of $c$ mean?

we know

$$C w_1 = \lambda_1 w_1$$

$$w_1^T C w_1 = w_1^T (\lambda_1 w_1) = \lambda_1$$

$$\lambda_1 = w_1^T C w_1 = w_1^T \left( \frac{1}{n} \sum_{i=1}^{n} x_i x_i^T \right) w_1$$
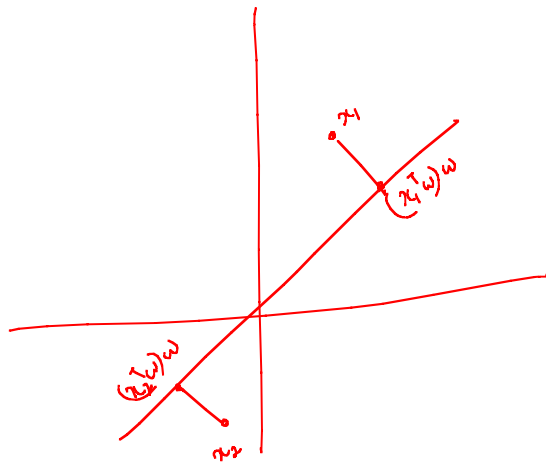
$$\boxed{\lambda_1 = \frac{1}{n} \sum_{i=1}^{n} (x_i^T w_1)^2} \quad \longleftarrow \text{term we used earlier}$$



Rule of thumb for # dimensions

$$\left( \sum_{i=1}^{k} \lambda_i(c) \Big/ \sum_{i=1}^{d} \lambda_i(c) \right) \geq 0.95$$

$\hookrightarrow$ usualy in practice

$$\left\{ (x_1^T \omega), \quad \cdots \quad , (x_n^T \omega) \right\}$$

<u>Average?</u>

$$\frac{1}{n} \sum_{i=1}^{n} (x_i^T \omega) = \left( \frac{1}{n} \sum_{i=1}^{n} x_i \right)^T \omega$$

$$\underbrace{\qquad\qquad}_{0^T \omega} \quad \left[ \begin{array}{c} \text{for a} \\ \text{centered} \\ \text{dataset} \end{array} \right].$$

$$= 0$$

<u>Variance</u>

$$\frac{1}{n} \sum_{i=1}^{n} (x_i^T \omega - \underbrace{mean}_{=0})^2 = \boxed{\frac{1}{n} \sum_{i=1}^{n} (x_i^T \omega)^2}$$

$$\begin{array}{ccc} \text{ERROR} \quad \text{MINIMIZATION} & <=> & \text{VARIANCE} \\ \text{on} & & \text{MAXIMIZATION} \\ \text{CENTERED DATASET} & & \end{array}$$

Want directions where

Projections don't "crowd-up"

i.e., Variance is not small.

ONE MORE EXAMPLE

$\frac{W-h}{\sqrt{2}}$

Weight (centered)

$\frac{h+w}{\sqrt{2}}$

$\begin{bmatrix} -1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$

$\begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$

height (centered)

height gives some information about weight.

$\frac{-height + weight}{\sqrt{2}}$

$\frac{height + weight}{\sqrt{2}}$

$(h+w)/\sqrt{2}$ does not give any information about $\frac{w-h}{\sqrt{2}}$.

"de-correlated"

PRINCIPAL COMPONENT ANALYSIS    —    $\{ W_1, \cdots, W_R \}$

PRINCIPAL COMPONENT

"DIMENSIONALITY REDUCTION"

PCA Finds combination of features that are de-correlated. [loosely speaking independent of each other].

" EIGEN FACES"