

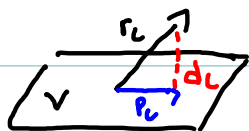
Lesson 24

Read chapter 8

PCA

Th: Given $A \in \mathbb{R}^{m \times n}$ of rank r
 the best fit k dimensional subspace of \mathbb{R}^n
 to the rows of A for $k=1 \dots r$ is $\text{span}(v_1 \dots v_k)$
 where $v_1 \dots v_k$ are the first k right singular
 vectors of A .

Best fit subspace V is the subspace that
 minimizes $\sum_{l=1}^m (\text{distance of row } l \text{ of } A \text{ from } V)^2$



$$\|d_l\|^2 = \|r_l\|^2 - \|p_l\|^2$$

Look for orthonormal set of vectors $b_1 \dots b_k$

s.t $\sum_{l=1}^m \|\text{Projection } r_l \text{ on } \text{span}(b_1 \dots b_k)\|^2$ is max

$$p_l = r_l^T b_1 b_1 + r_l^T b_2 b_2 + \dots + r_l^T b_k b_k$$

$$\|p_l\|^2 = (r_l^T b_1)^2 + \dots + (r_l^T b_k)^2$$

want to maximize

$$\begin{aligned} & (r_1^T b_1)^2 + \dots + (r_1^T b_k)^2 + \\ & (r_2^T b_1)^2 + \dots + (r_2^T b_k)^2 + \\ & \dots \dots \dots \\ & (r_m^T b_1)^2 + \dots + (r_m^T b_k)^2 \end{aligned}$$

$$\|A b_1\|^2 + \|A b_2\|^2 + \dots + \|A b_k\|^2$$

$$\begin{matrix} (r_1^T b_1)^2 + (r_1^T b_2)^2 + \dots + (r_1^T b_k)^2 \\ (r_2^T b_1)^2 + (r_2^T b_2)^2 + \dots + (r_2^T b_k)^2 \\ \vdots \\ (r_m^T b_1)^2 + (r_m^T b_2)^2 + \dots + (r_m^T b_k)^2 \end{matrix}$$

↓

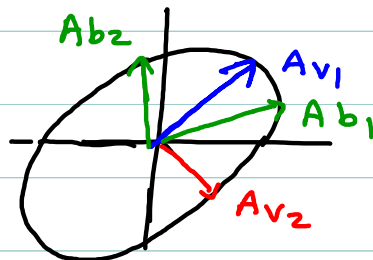
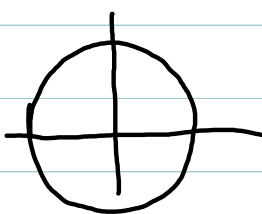
↓ Taking $b_2 = v_2$ maximizes this column $\|A b_2\|^2$

Taking $b_1 = v_1$ maximizes this column $\|A b_1\|^2$

but how do I know I cannot choose

b_1 and b_2 with $\|A b_1\| < \|A v_1\|$

but $\|A b_2\| > \|A v_2\|$ so that $\|A b_1\|^2 + \|A b_2\|^2 > \|A v_1\|^2 + \|A v_2\|^2$

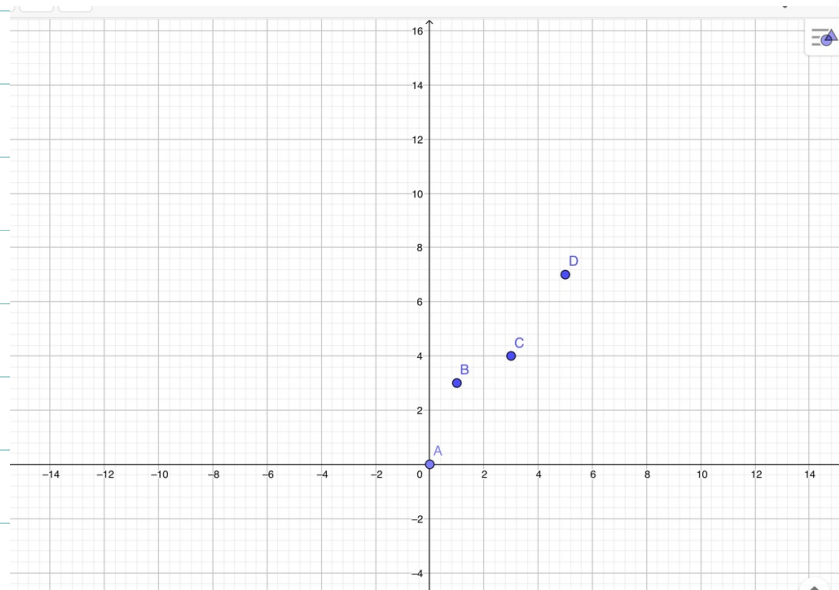


Th: Given $A \in \mathbb{R}^{m \times n}$ of rank r
the best fit k dimensional subspace of \mathbb{R}^m
to the columns of A for $k=1 \dots r$ is $\text{span}(u_1, \dots, u_k)$
where u_1, \dots, u_k are the first k left singular
vectors of A .

Ex: In hw 5 you found
linear regression line to:

$$(0,0) (1,3) (3,4) (5,7) : y = 0.68 + 1.25x$$

What is the best fit line?



Note: maybe no line through origin
is very good fit to data, if data
is not centered around origin, so we
may want to shift our data first.

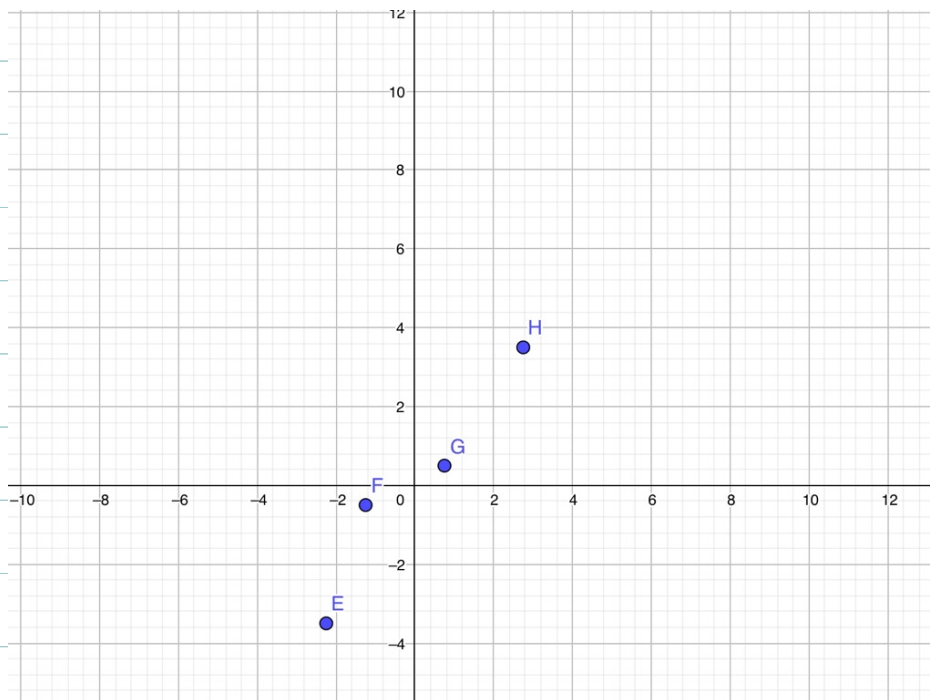
Data is in rows : Subtract mean of each column
to center around origin .

$$A = \begin{bmatrix} 0 & 0 \\ 1 & 3 \\ 3 & 4 \\ 5 & 7 \end{bmatrix}$$

$$x \text{ mean} = \frac{0+1+3+5}{4} = 9/4$$

$$y \text{ mean} = \frac{0+3+4+7}{4} = \frac{14}{4} = \frac{7}{2}$$

$$B = \begin{bmatrix} -9/4 & -7/2 \\ 1-9/4 & 3-7/2 \\ 3-9/4 & 4-7/2 \\ 5-9/4 & 7-7/2 \end{bmatrix} = \begin{bmatrix} -9/4 & -7/2 \\ -5/4 & -1/2 \\ 3/4 & 1/2 \\ 11/4 & 7/2 \end{bmatrix} \begin{matrix} E \\ F \\ G \\ H \end{matrix}$$



```

[1] top-level scope
@ none:1

julia> using LinearAlgebra

julia> a=[-9/4 -7/2; -5/4 -1/2; 3/4 1/2; 11/4 7/2]
4x2 Matrix{Float64}:
-2.25 -3.5
-1.25 -0.5
 0.75  0.5
 2.75  3.5

julia> svd(a)
SVD{Float64, Float64, Matrix{Float64}}
U factor:
4x2 Matrix{Float64}:
-0.663581  0.398393
-0.184732 -0.8405
 0.136306  0.357277
 0.712008  0.0848297
singular values:
2-element Vector{Float64}:
 6.25074081965229
 0.8235527946237668
Vt factor:
2x2 Matrix{Float64}:
 0.605404  0.795918 * v1
 0.795918 -0.605404

julia>

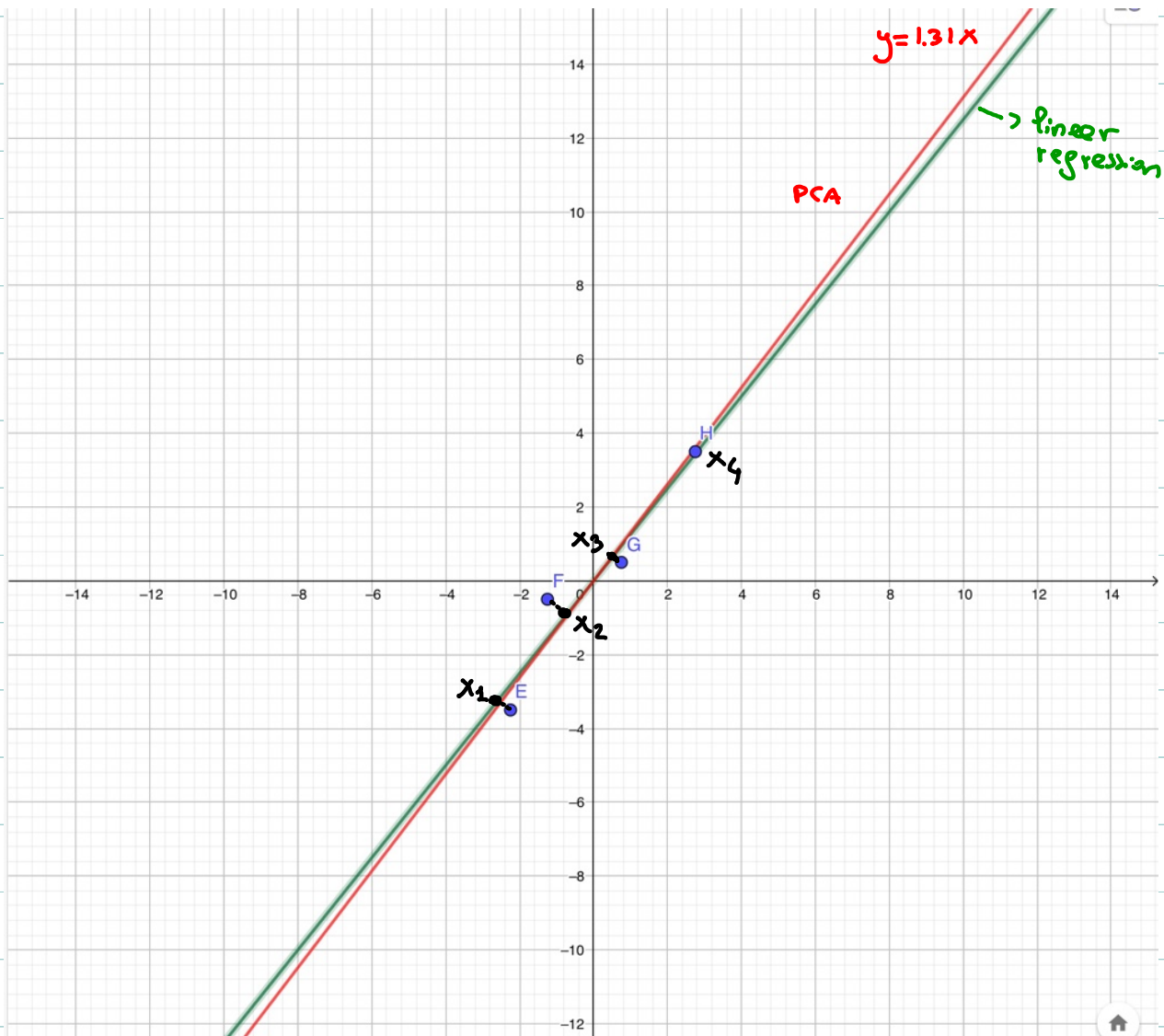
```

$$\begin{bmatrix} x \\ y \end{bmatrix} = \underbrace{\begin{bmatrix} 0.605 \\ 0.796 \end{bmatrix}}_{v_1}$$

$$y = \frac{0.796}{0.605} x$$

$$y = 1.31 x$$

Best fit line is line through origin in the direction of v_1 .



$y + \frac{7}{2} = 0.68 + 1.25(x + 9/4)$ Linear regression line

$$y = 1.25x$$

$$y = 1.31x \quad \text{Best fit line}$$

Th Given $A = \sum_{l=1}^r \sigma_l u_l v_l^T$ and

$$A_k = \sum_{l=1}^k \sigma_l u_l v_l^T, \text{ then}$$

the rows of A_k are the projections of
the rows of A onto $V_k = \text{span}(v_1, v_2, \dots, v_k)$

$$A = \begin{bmatrix} r_1^T \\ \vdots \\ r_m^T \end{bmatrix} \quad A_k = \begin{bmatrix} (\text{Proj}_{V_k} r_1)^T \\ \vdots \\ (\text{Proj}_{V_k} r_m)^T \end{bmatrix}$$

Proof: the l^{th} row of A_k is

$$e_l^T A_k = e_l^T \sum_{j=1}^k \sigma_j u_j v_j^T = \sum_{j=1}^k \sigma_j e_l^T u_j v_j^T$$

$$l^{\text{th}} \text{ row of } A \text{ is: } r_l^T = e_l^T A$$

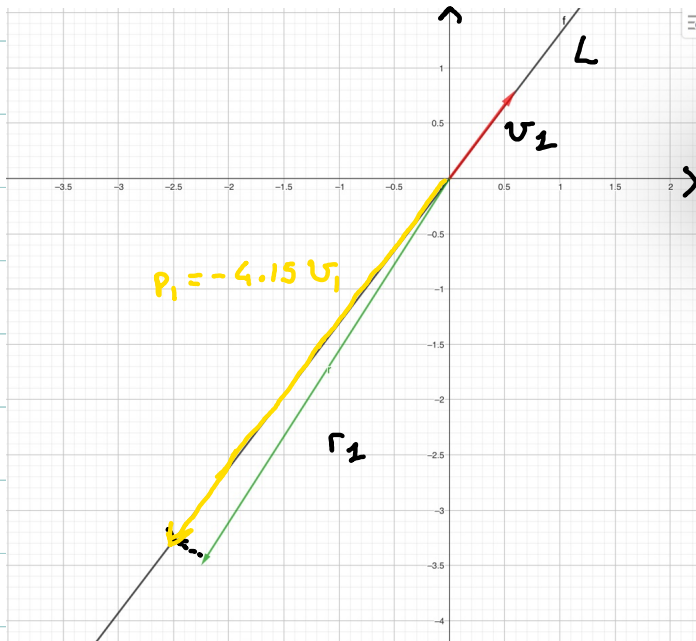
$$P_{V_k} r_l = \sum_{j=1}^k \underbrace{r_l^T v_j}_{\text{scalar}} v_j$$

(projection formula
starting from
orthonormal basis.
see hw)

We are transposing to make this into a 1xn vector

$$\begin{aligned} (P_{V_k} r_l)^T &= \sum_{j=1}^k e_l^T A v_j v_j^T = \sum_{j=1}^k e_l^T U \Sigma V^T v_j v_j^T = \\ &= \sum_{j=1}^k e_l^T U \sigma_j e_j v_j^T = \sum_{j=1}^k \sigma_j e_l^T u_j v_j^T \end{aligned}$$

Back to our example



$$A = \begin{bmatrix} -9/4 & -7/2 \\ -5/4 & -1/2 \\ 3/4 & 1/2 \\ 11/4 & 7/2 \end{bmatrix} \begin{matrix} E \\ F \\ G \\ H \end{matrix}$$

From SVD: $v_1 = \begin{bmatrix} 0.605 \\ 0.796 \end{bmatrix}$, $u_1 = \begin{bmatrix} -0.66 \\ 0.18 \\ 0.14 \\ 0.71 \end{bmatrix}$, $\sigma_1 = 6.25$.

$L = \text{span}(v_1)$

$$\text{Proj}_L \left(\begin{bmatrix} -9/4 \\ -7/2 \end{bmatrix} \right) = x_1 v_1$$

$$x_1 = \begin{bmatrix} -9/4 & -7/2 \end{bmatrix} \underbrace{\begin{bmatrix} 0.605 \\ 0.796 \end{bmatrix}}_{v_1} \approx -4.15$$

$$\text{Proj}_{v_1}(r_1) = -4.15 \begin{bmatrix} 0.605 \\ 0.796 \end{bmatrix} \approx \begin{bmatrix} -2.51 \\ -3.3 \end{bmatrix}$$

$$A_1 = 6.251 \begin{bmatrix} -0.664 \\ -0.18 \\ 0.14 \\ 0.71 \end{bmatrix} \begin{bmatrix} 0.605 & 0.796 \end{bmatrix} = 6.251 \begin{bmatrix} -0.664 \times 0.605 & -0.664 \times 0.796 \\ \dots & \dots \\ \dots & \dots \\ \dots & \dots \end{bmatrix}$$

$4 \times 1 \quad 1 \times 2$

$$= \begin{bmatrix} -2.51 & -3.3 \\ \dots & \dots \\ \dots & \dots \\ \dots & \dots \end{bmatrix}$$

Two aspects of PCA :

- 1) dimension reduction
- 2) maximum variation of data

Example of PCA:

In a study, 200 gene expressions were studied in 96 tissue samples
Data are 96 vectors in \mathbb{R}^{200}

$A = 96 \times 200$ matrix of data

To visualize the data we want to project it onto 2 dimensional subspace: best fit 2 subspace.

1) Preprocessing like centering data $A \rightarrow A'$

2) $A' = U \Sigma V^T$

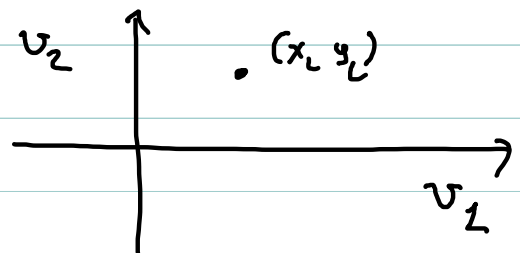
3) Best fit 2 dimensional space = $\text{span}(v_1, v_2) = V_2$

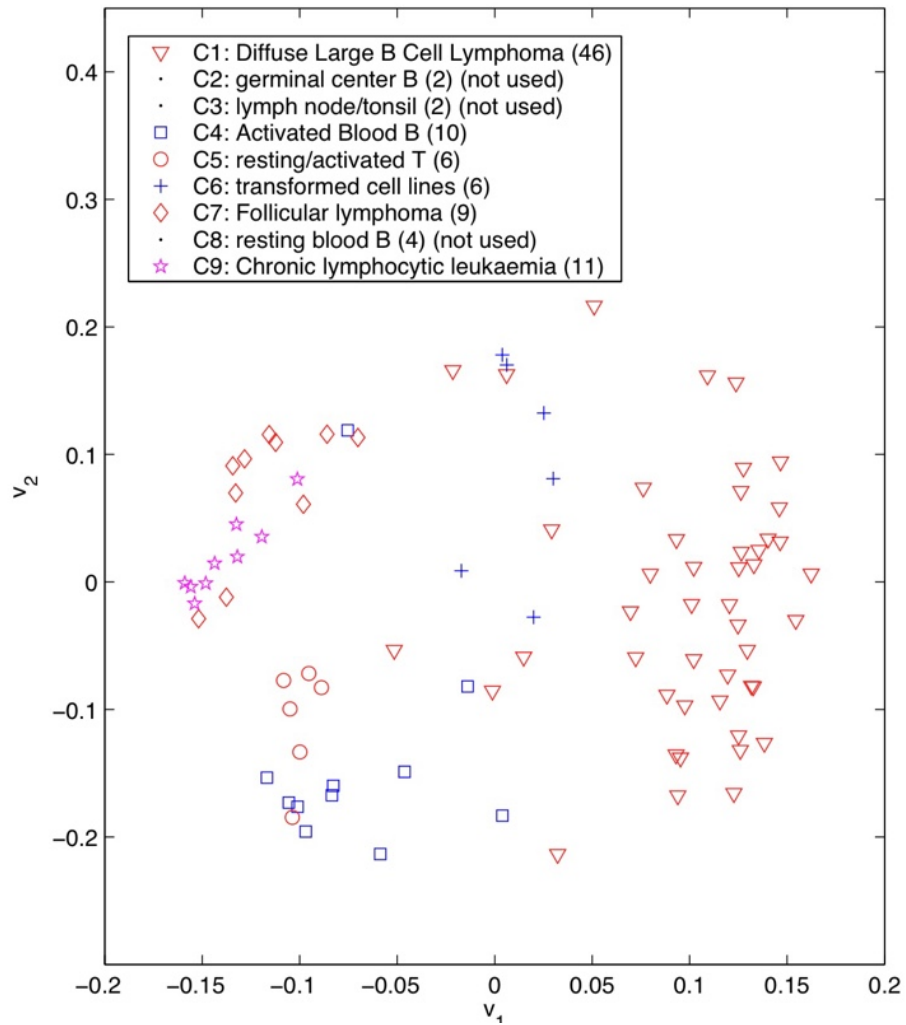
4) Project the rows of A' onto V_2 :

$$\text{Proj}_{V_2} r_i = x_i v_1 + y_i v_2$$

$$x_i = r_i^T v_1, \quad y_i = r_i^T v_2.$$

5) plot all points (x_i, y_i)
in 2D space





Plot of the points (x_i, y_i)

Maximum variation

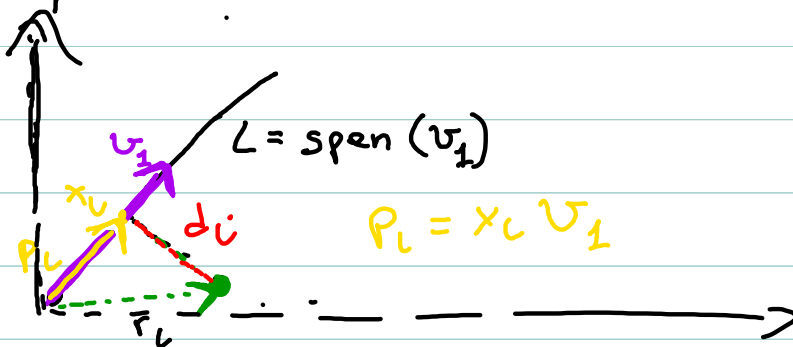
See animation at: setosa.60

data centered around origin

P_L projection of r_L on v_1 (principal

component = first right singular vector)

$$P_L = x_L v_1 \quad x_L = r_L^T v_1 \quad x_L = \|P_L\|$$



$$\textcircled{1} \sum_{L=1}^m x_L = \sum_{L=1}^m r_L^T v_1 = 0 \quad \text{since } \sum_{L=1}^m r_L^T = 0$$

(data is centered) so x_L are also centered around 0. How spread out are they?

$\textcircled{2}$ In statistics variance is a measure

of how data is spread out around mean

$$\frac{1}{m-1} \sum_{L=1}^m (x_L - \bar{x})^2 = \frac{1}{m-1} \sum_{L=1}^m x_L^2 = \frac{1}{m-1} \sum_{L=1}^m \|P_L\|^2$$

since we best fit line maximizes this sum

it maximizes the variance of the data

Example to explain proof of Eckert Young th.

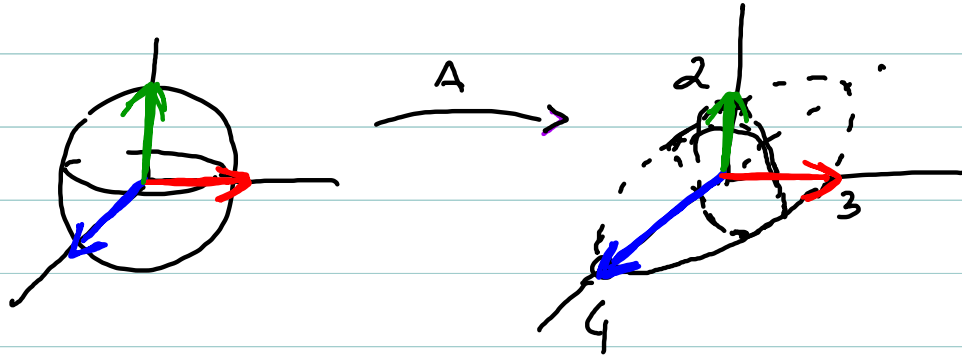
$$A = \begin{bmatrix} 4 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 2 \end{bmatrix} = I \begin{bmatrix} 4 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 2 \end{bmatrix} I^T =$$

$$= 4 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} [100] + 3 \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} [010] + 2 \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} [001]$$

$$A_1 = \begin{bmatrix} 4 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$B = \begin{bmatrix} 3.5 & 3.5 & 0 \\ 3.5 & 3.5 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

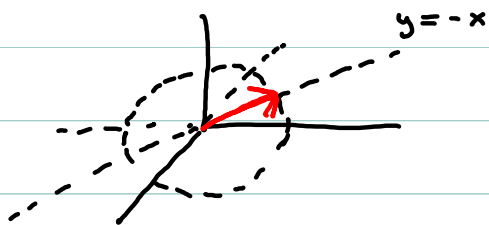
We know $\|A - A_1\|_2 = 3$ so we must
have $\|A - B\|_2 \geq 3$



want w , $\|w\|=1$ s.t. $\|Aw - Bw\| \geq 3$

$\text{Null}(B)$ is the plane $x+y=0$
 $\text{Span}\left(\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}\right)$ is the plane $z=0$

want w on both of these planes



$$w = \begin{bmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \\ 0 \end{bmatrix}$$

$$\| \begin{bmatrix} 4 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 2 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \\ 0 \end{bmatrix} \| = \| \begin{bmatrix} 4/\sqrt{2} \\ -3/\sqrt{2} \\ 0 \end{bmatrix} \| = \sqrt{\frac{16}{2} + \frac{9}{2}} = \frac{5}{\sqrt{2}}$$

$$\approx 3.54$$

