

Correlation and Cosine Similarity: An Introduction to Thinking Geometrically

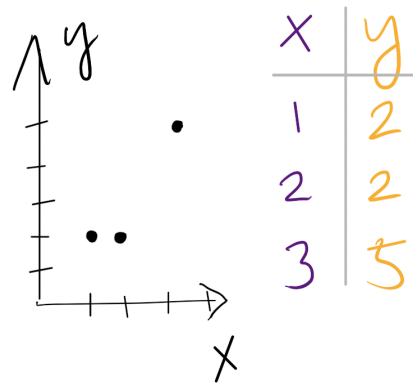
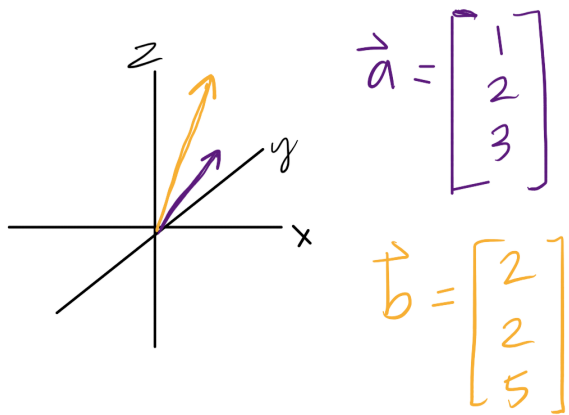
MOTIVATIONS

Why go high dimensional?

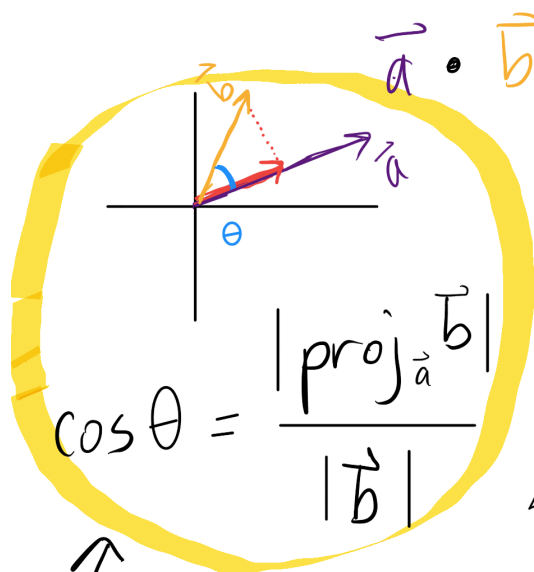
- patterns
- intuition
- novel techniques
- THE WORLD IS MULTI DIMENSIONAL! ☺

DATA + QUESTIONS

How similar or related are two sets of data?



DOT PRODUCT REVIEW



$$\vec{a} \cdot \vec{b} = |\text{proj}_{\vec{a}} \vec{b}| \cdot |\vec{a}|$$

$$= |\vec{b}| \cos \theta \cdot |\vec{a}|$$

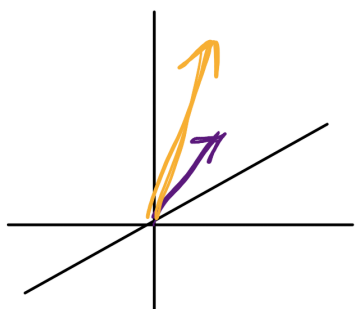


Force and Displacement

$$\text{Work} = \vec{F} \cdot \Delta \vec{x}$$

In other words, proportion of \vec{b} in direction of \vec{a} .

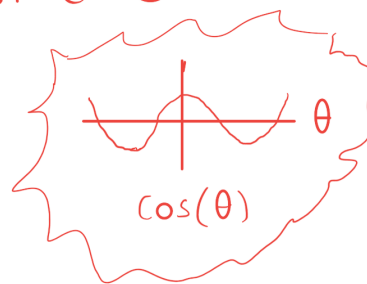
DATA! (HOW SIMILAR ARE my VECTORS?)



$$\cos \theta = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|}$$

Similarity
Metric:

cosine θ

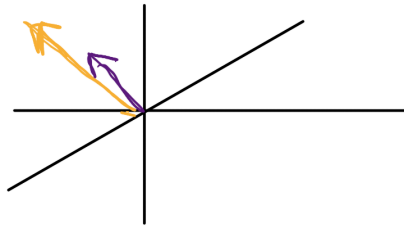


LET'S CALCULATE!

let's ^①mean-center the vectors
and ^②take cosine similarity

$$\textcircled{1} \quad \vec{a} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \quad \vec{a}_c = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix}$$

$$\vec{b} = \begin{bmatrix} 2 \\ 2 \\ 5 \end{bmatrix} \quad \vec{b}_c = \begin{bmatrix} -1 \\ -1 \\ 2 \end{bmatrix}$$



$$\begin{aligned}
 \textcircled{2} \\
 \cos \theta &= \frac{\vec{a}_c \cdot \vec{b}_c}{|\vec{a}| |\vec{b}|} = \frac{\begin{bmatrix} -1 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} -1 \\ -1 \\ 2 \end{bmatrix}}{\sqrt{(-1)^2 + (1)^2} \sqrt{(-1)^2 + (-1)^2 + (2)^2}} \\
 &= \frac{\overbrace{(-1)(-1)}^1 + \overbrace{0}^0 + \overbrace{(1)(2)}^2}{\sqrt{2} \sqrt{6}} \\
 &= \sqrt{3}/2 = .86
 \end{aligned}$$

DATA! (HOW SIMILAR ARE THESE DATASETS?)

x	y
1	2
2	5
3	5

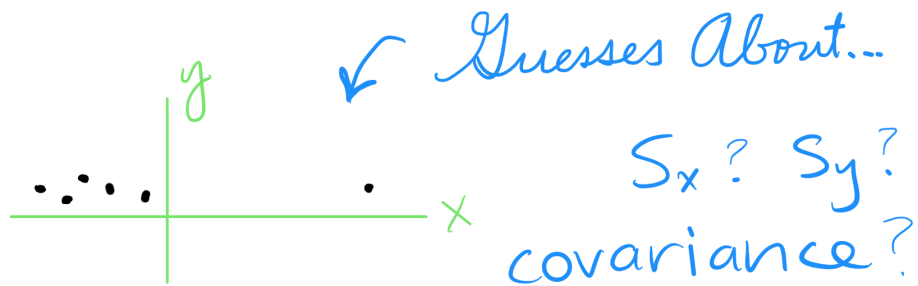
Similarity
measure:
Pearson R

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

CORRELATION REVIEW!

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

s_x \swarrow \nwarrow s_y \nwarrow covariance



LET'S CALCULATE!

x	1	2	3	$\bar{x} = 2$
y	2	2	5	$\bar{y} = 3$

$$r = \frac{\overbrace{(1-2)(2-3)}^1 + \overbrace{(2-2)(2-3)}^0 + \overbrace{(3-2)(5-3)}^2}{\sqrt{(1-2)^2 + (2-2)^2 + (3-2)^2} \sqrt{(2-3)^2 + (2-3)^2 + (5-3)^2}}$$

$$r = \frac{3}{\sqrt{12}} = \frac{\sqrt{3}}{2} = \boxed{.86}$$

CONCLUSION

Pearson's r \equiv Cosine Similarity of mean-centered vectors

One example of the many ways we can think about our data geometrically ☺