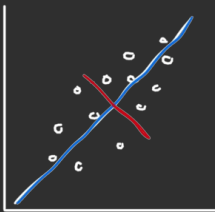# PCA

principal component analysis - extracts essential info
by identifying useful dimensions in multi-dimesional
data


with 2D data, need $2n$ numbers to represent
$n$ data points.

σ² in y
is 0.



✱ even though this is 2D we
could represent the data just
as effectively using $n$ points
due to no variance in y
thus taking out the need for
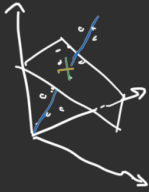an "excess" dimension.


another scenario:



first principal comp.
second principal comp.

direction of the arrows are the
principal components, since they are
the directions of max variance and
thus max information

how in 3D

first
second  } max variance direc.
third

the point :
using those principal component directions we
can create a new set of axes $x', y', z'$.
thus every point can be represented by
$(x', y', z')$. we know that $x$ has the most variance
then $y$ then $z$, thus we can get rid of $z$ and store
much of the same information in $(x', y')$.
huge for sets with lots of dimensions (i.e. $1000 \rightarrow 20$).


## how to calculate PCA

Data matrix:
$$\begin{bmatrix} x_1 & x_2 & \cdots & x_n \\ y_1 & y_2 & \cdots & y_n \\ z_1 & z_2 & \cdots & z_n \end{bmatrix}$$

← this is a 3D matrix
could be more columns

Calculate mean :

$$\mu_x = \frac{1}{n} \sum_i^n x_i$$

$$\mu_y = \frac{1}{n} \sum_i^n y_i$$

$$\mu_z = \frac{1}{n} \sum_i^n z_i$$

$(\mu_x, \mu_y, \mu_z)$ 3D

$(\mu_A, \mu_B, \cdots \mu_m)$ M Dimensional

To find the direction of maximum variance we can calculate the covariance matrix :

$$C = MM^T$$

← $n \times n$ matrix.

such that $M = \begin{bmatrix} (x_1 - \mu_x) & (x_2 - \mu_x) & \cdots & (x_n - \mu_x) \\ (y_1 - \mu_y) & (y_2 - \mu_y) & \cdots & (y_n - \mu_y) \\ \vdots & \vdots & & \ddots \end{bmatrix}$

The diagonal of the covariance matrix are the variances along the X, Y, and Z axes. The off diagonal represent covariance between two dimensions (X/Y, Y/Z, X/Z).



$$C = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$C = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$

$$C = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$$

The eigenvectors of the covariance matrix are the principal components, in order of size.

eigenvalues & eigenvectors:

$$Av = \lambda v$$

→ matrix

↘ scalar

$\lambda$ is the eigenvalue fo the eigenvector, $v$

ex:    consider   a  3×3 matrix

$$A = \begin{bmatrix} 1.04 & 1.101 & 0.83 \\ 1.10 & 1.47 & 1.10 \\ 0.83 & 1.40 & 0.8931 \end{bmatrix}$$

now   consider, v

$$v = \begin{bmatrix} -0.1233 \\ 0.6644 \\ -0.7371 \end{bmatrix}$$

notice that   $Av = \begin{bmatrix} -0.0051 \\ 0.0274 \\ -0.0304 \end{bmatrix} = 0.0412 \begin{bmatrix} -0.01233 \\ 0.6644 \\ -0.7371 \end{bmatrix}$

↑                ↑
A                λ

λ  is the eigen value, and  v  is the eigen vector

throw backs to multi.