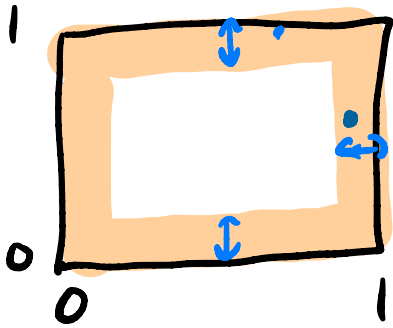<u>Lecture 17</u>

<u>Example:</u> 10,000-D data.
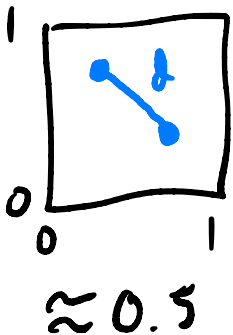


prob a point is within 0.001
of border. = 0.004

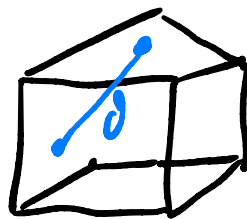<u>But</u> in 10,000 D hypercube $\Rightarrow$ $\underline{> 0.9999}$!

$\Rightarrow$ Virtually all data becomes very extreme!

Avg distance between 2 random pts

2D:



$\approx 0.5$

3D:



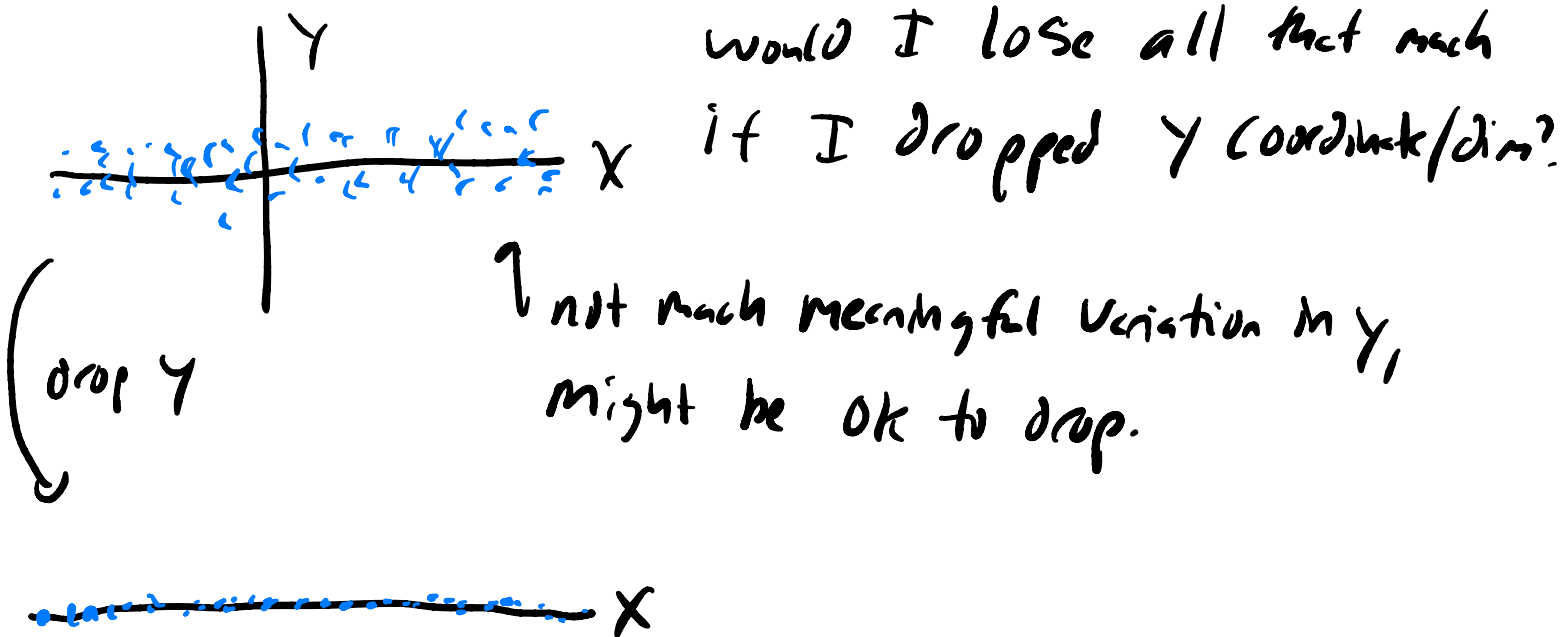$\approx 0.66$

10,000 D hypercube
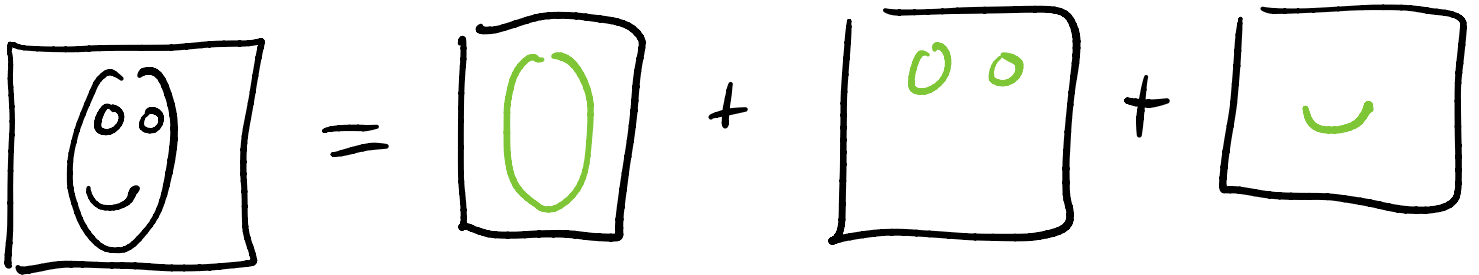
$\approx 408.25$!!

higher dimension, on avg, things are more spaced out

**Sparse**

Curse of dimensionality: Techniques that work in lower dimensions fail in high dimensional spaces.

Work around:

would I lose all that much if I dropped Y coordinate/dim?

↑ not much meaningful variation in Y, might be ok to drop.

(drop Y)
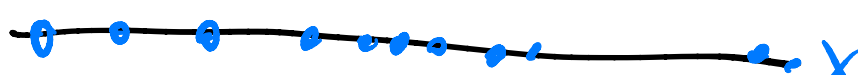
When we have 10,000 Dimension, maybe 30 variables are actually meaningful — find out how to reduce 10,000 D $\xrightarrow{\text{project}}$ 30 D ⌐ run analysis on 30 meaningful vars

principal component analysis (PCA): popular technique to reduce dimensionality — algorithm to tell us which variables are important.
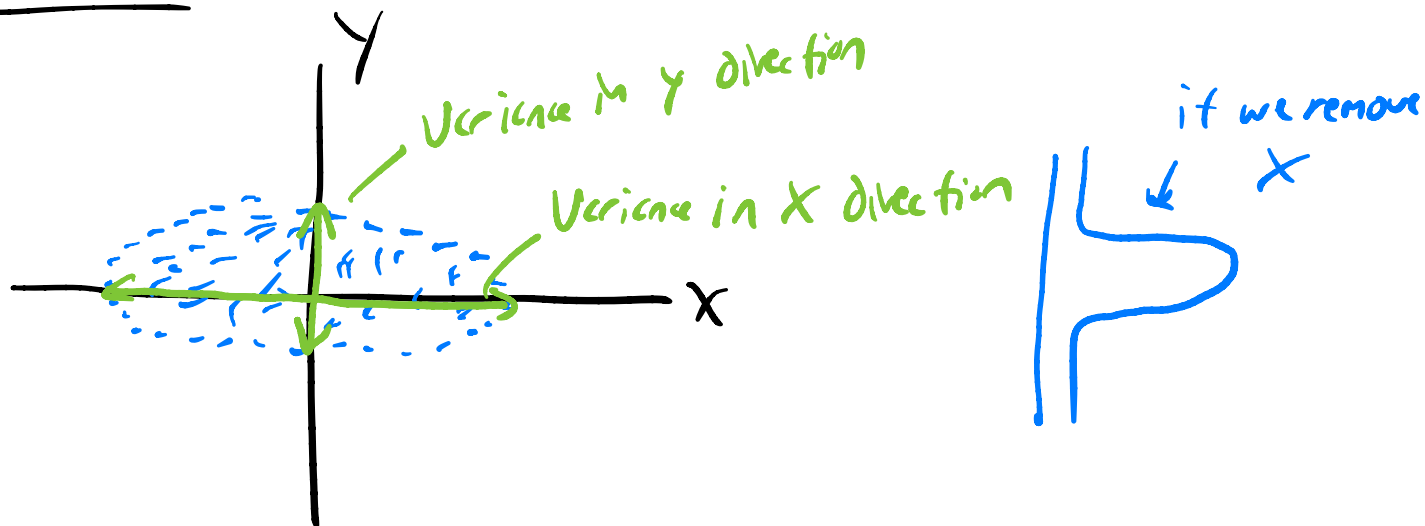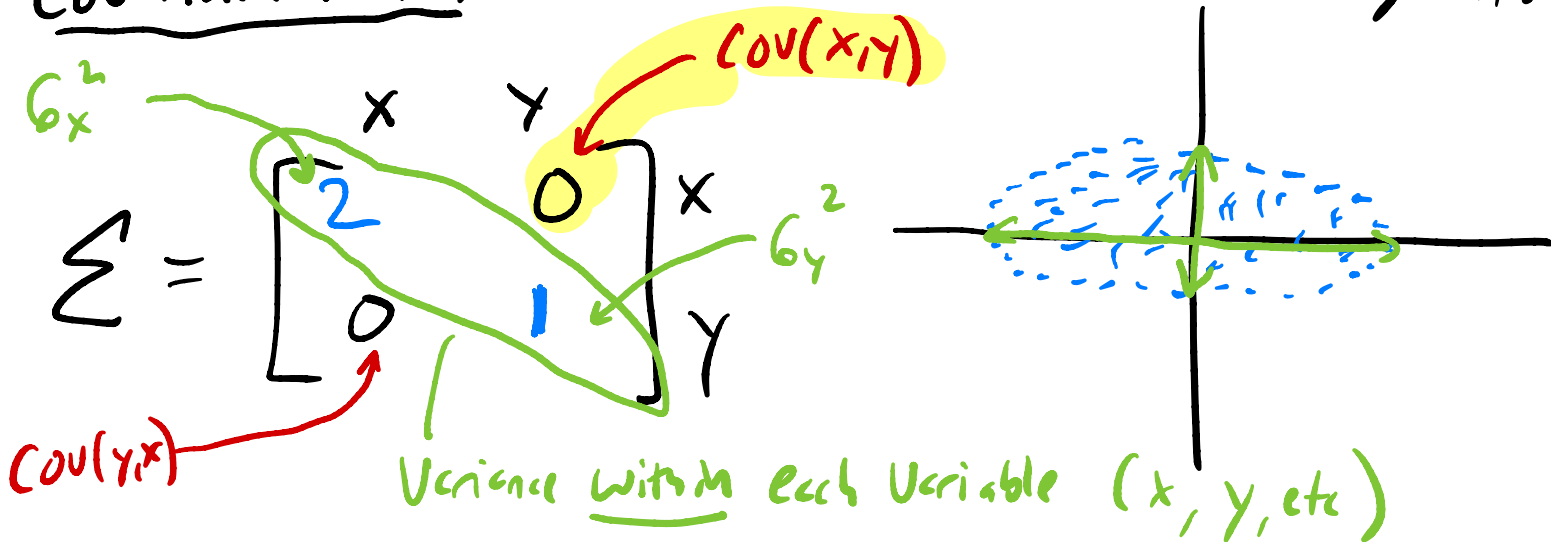


---

☑ Covariance matrix

1)  x    low var

2)  x    high var

Variance: measure of spread/dispersion in data
Scaler value

# 2D data



Variance in Y direction

Variance in X direction

if we remove X

if we remove Y

$$Var(x) > Var(y)$$

Covariance matrix: Store info about Variances among Variables

$Cov(x,y)$

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{matrix} X \\ Y \end{matrix}$$

$\sigma_x^2$

$\sigma_y^2$

$Cov(y,x)$

Variance within each Variable $(x, y, etc)$
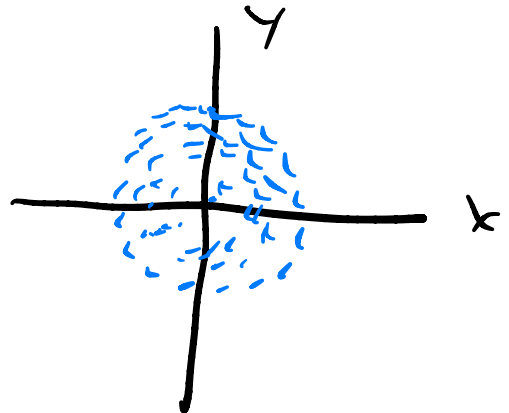
Off diagonals: How vars vary together jointly

Covariance: $Cov(x,y)$

above it is 0

$$\not{\star} \quad \boxed{Cov(x,y) = Cov(y,x)}$$

↳ Covariance matrix always must be Symmetry

$$\Sigma = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix}$$



$$\Sigma = \begin{bmatrix} 3 & 1 \\ 0 & 3 \end{bmatrix}$$

✓ impossible! not Symmetry

$$\Sigma = \begin{bmatrix} 20 & +5 \\ +5 & 1 \end{bmatrix}$$