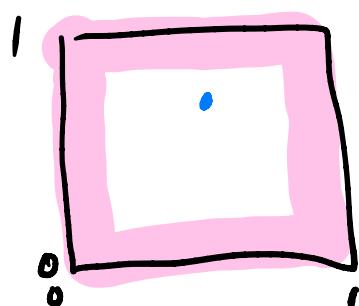


Project 4: Principal Component Analysis

Curse of dimensionality

Data behaves very differently in a high dimensional space (e.g. 10,000-D) compared to 2D/3D.

Prevalence of extreme values



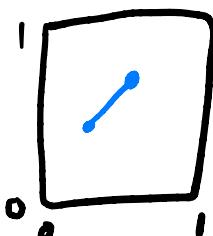
2D unit Square :

$$\text{probability pt within } 0.001 \text{ of border} = 0.000001$$

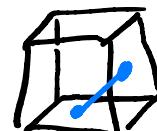
10,000-D hypercube : Same probability >0.9999!

\Rightarrow In high dimensional spaces, almost all points are "extreme".

Distance inflation :



2D unit Square :
avg dist between
2 pts ≈ 0.5



3D unit Cube :
avg dist
 ≈ 0.66

[You can simulate this to verify!]

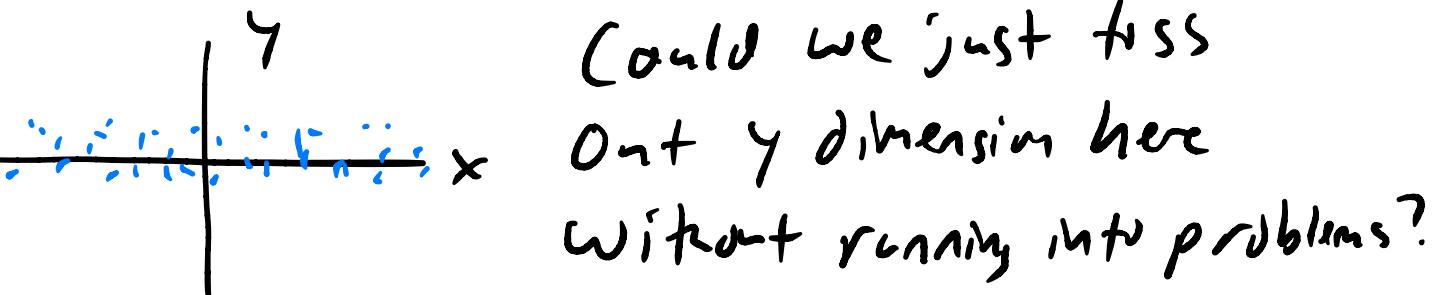
10,000-D hypercube: avg dist between 2 pts
 $\approx \underline{408.25}$!

\Rightarrow Things are much farther apart in high dimensional spaces!
Space is sparsely filled with data

• Curse of dimensionality: Refers to these problems.

Difficult to analyze data in high dimensional spaces — techniques that find patterns in 2D/3D may not generalize to 10,000-D.

Work around for Curse of dimensionality:



- Real data not evenly distributed across all variables — Not all 10,000 vars carry equal weight — maybe only 30 do.

Project data on those 3D



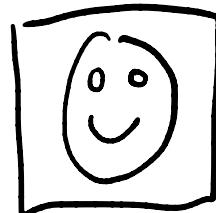
Run analysis on those



Run analysis on 3D-D
data, not 10,000-D.

principal component analysis (PCA) is a popular
technique for doing this dimensionality reduction.

In project, you will apply PCA
to images of human faces:



PCA allows us, without reducing # pixels, to discover intrinsic
variables in image data that are important (e.g., eyes, mouth, head).

⇒ Can express an image as linear combination of these variables:

$$\begin{matrix} \text{Smiley face} \\ \text{image} \end{matrix} = \begin{matrix} \text{Eye} \\ \text{image} \end{matrix} + \begin{matrix} \text{Mouth} \\ \text{image} \end{matrix} + \begin{matrix} \text{Head} \\ \text{image} \end{matrix}$$

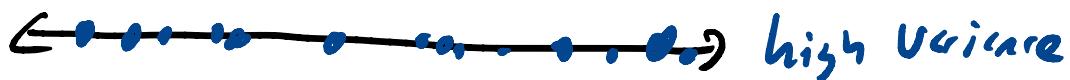
[Show Eigenface examples]

Application: Facial recognition

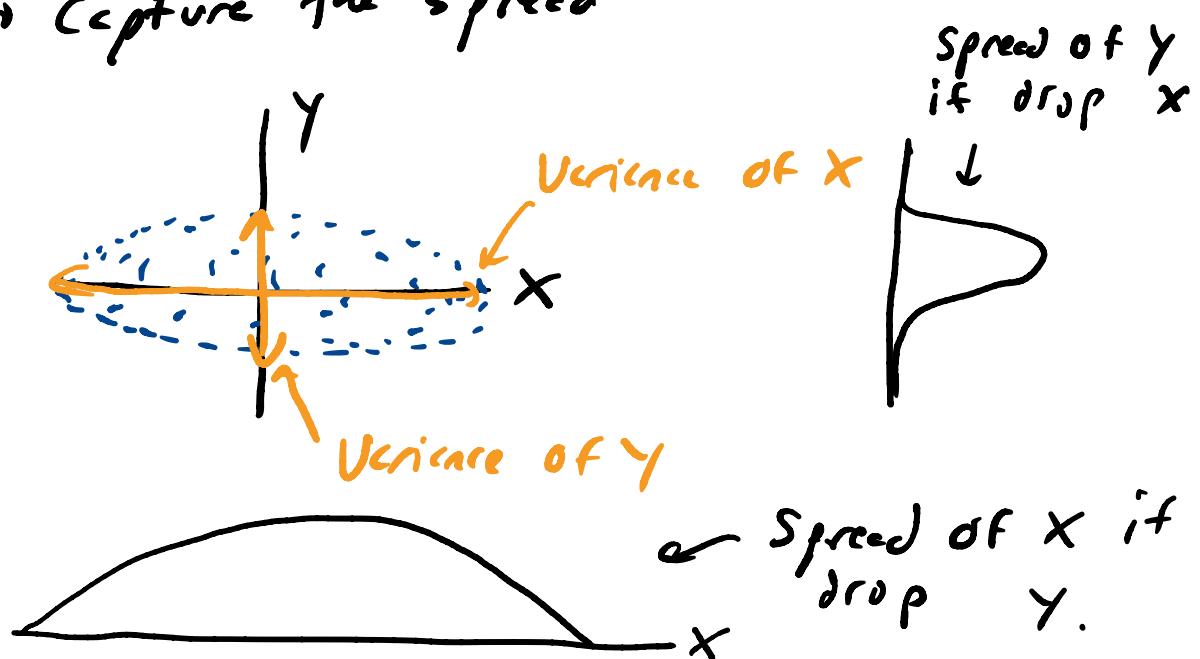
PCA depends on something called the **Covariance Matrix** — Let's introduce this concept 1st before diving into PCA.

Covariance Matrix

For 1D data, the ordinary sample variance describes the spread of data:



For 2D data, we can't just have a single variance to capture the spread

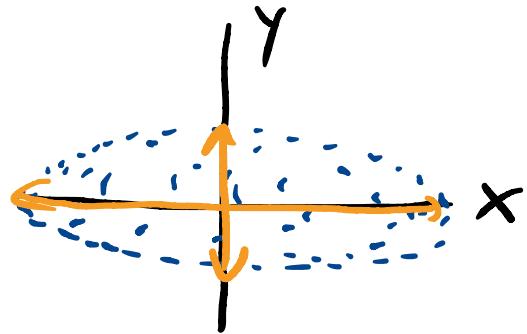


Clearly, $\text{Variance}(X) > \text{Variance}(Y)$.

Instead of a scalar to describe the variance of 2D data or higher dimensional data we use a matrix — called Covariance matrix Σ .

Plausible Covariance matrix for the above 2D data:

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \begin{matrix} x \\ y \end{matrix}$$



Main diagonal: Variance within each variable
 — e.g. 2 is $\text{Var}(x)$
 1 is $\text{Var}(y)$
 $\Rightarrow \text{Var}(x) = 2 \text{Var}(y)$

Off-Main diagonal: Variance between pairs of vars
 : Called Covariance.

e.g. row 0, col 1

$$\text{Cov}(x, y) = 0$$

[x does not vary with y]

e.g. row 1, col 0

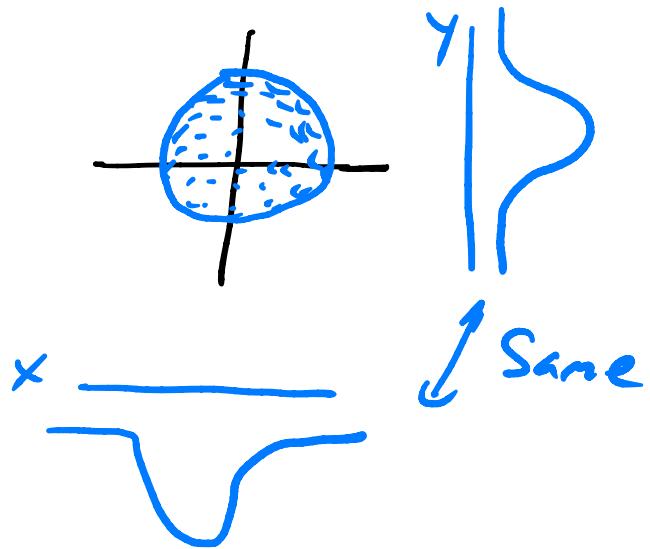
$$\text{Cov}(y, x) = 0$$

[y does not vary with x]

Σ , the Covariance Matrix, is always Symmetric about the main diagonal b/c $\text{Cov}(x,y) = \text{Cov}(y,x)$

Picture for these Σ ?

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



$$\Sigma = \begin{bmatrix} 3 & 1 \\ 0 & 3 \end{bmatrix} \Rightarrow \text{impossible, not symmetric}$$

