

Could we just toss
Out y dimension here
without running into problems?

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

Project data on those 30



Run analysis on those



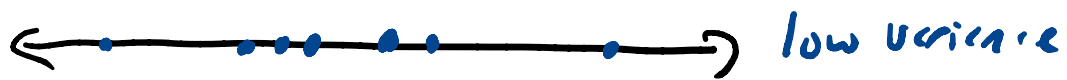
Run analysis on 30-D
data, not 10,000-D.

Principal Component analysis (PCA) is a popular technique for doing this dimensionality reduction.

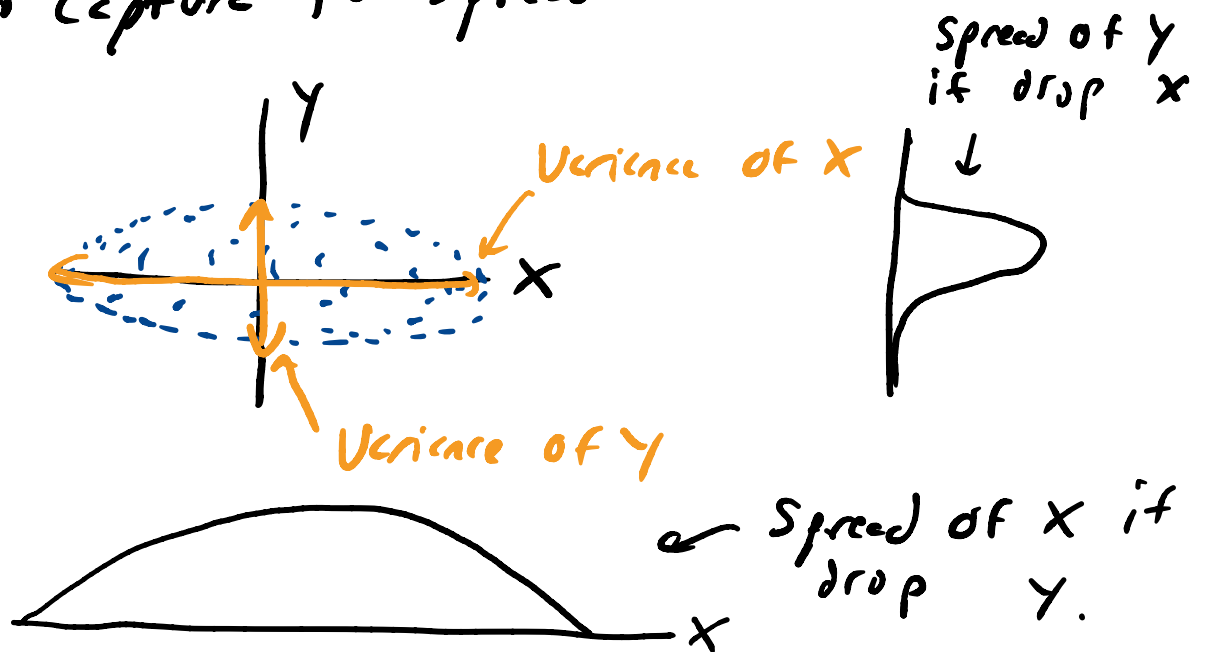
PCA depends on something called the **Covariance matrix** — Let's introduce this concept 1st before diving into PCA.

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



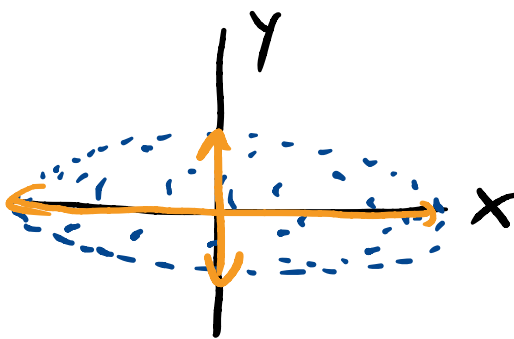
Clearly $\text{Variance}(X) > \text{Variance}(Y)$.

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:

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

$x \quad y$



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

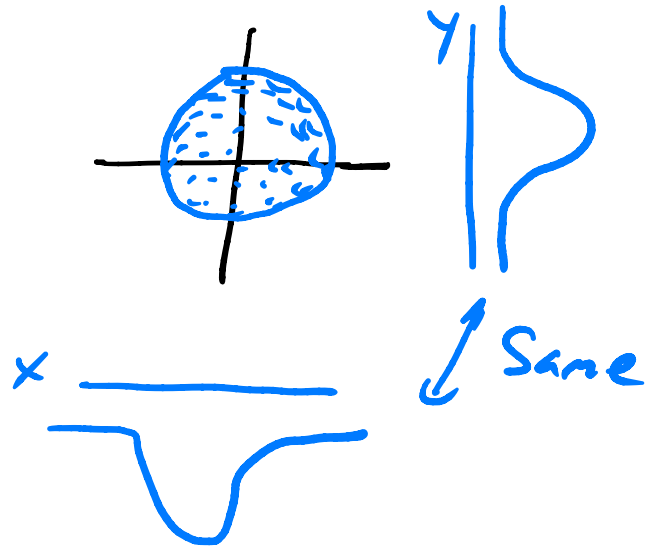
$$\text{COV}(y, x) = 0$$

[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 ?

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



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

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

positive correlation in x & y

