## Principal Component Analysis

or

Dimensionality Reduction

n_components = 6

or

Feature Reduction/Column Reduction

## Step by Step Breakdown of PCA

1. Scale the data using standardization technique. ✓
2. Calculates the Variance. ✓
3. Calculates the Co Variance ✓
4. Reduces the data in a Covariance Matrix. ✓
5. Calculates the Eigen Values of the reduced Cov. Matrix ✓

Step 1. Calculation of Z Scores

$$Z = X - \bar{X}$$

Discrete values of X

Mean of X

Standard Deviation $\sqrt{\dfrac{\Sigma(x-\bar{x})^2}{N}}$

Step 2.

$$\text{Variance}(x) = \left(\sigma\right)^2 = \frac{\sum \left(x - \bar{x}\right)^2}{}$$

No. of Data Points

Step 3

Formula for Co variance    Joint variance measure of two different features

Population Covariance Formula

$$Cov(x,y) = \sum \left( x - \bar{x} \right) . \left( y - \bar{y} \right)$$

Cov(y,x) & Cov(x,y) remain the same

Population Covariance means calculating co variance of the entire dataset.

# Sample Covariance Formula

$$\frac{\sum (X - \bar{X}) \cdot (X - Y)}{N - 1}$$

N → sample size

# Covariance Matrix

## Diff b/w Corelation and Covariance

Correlation ranges from -1 to +1

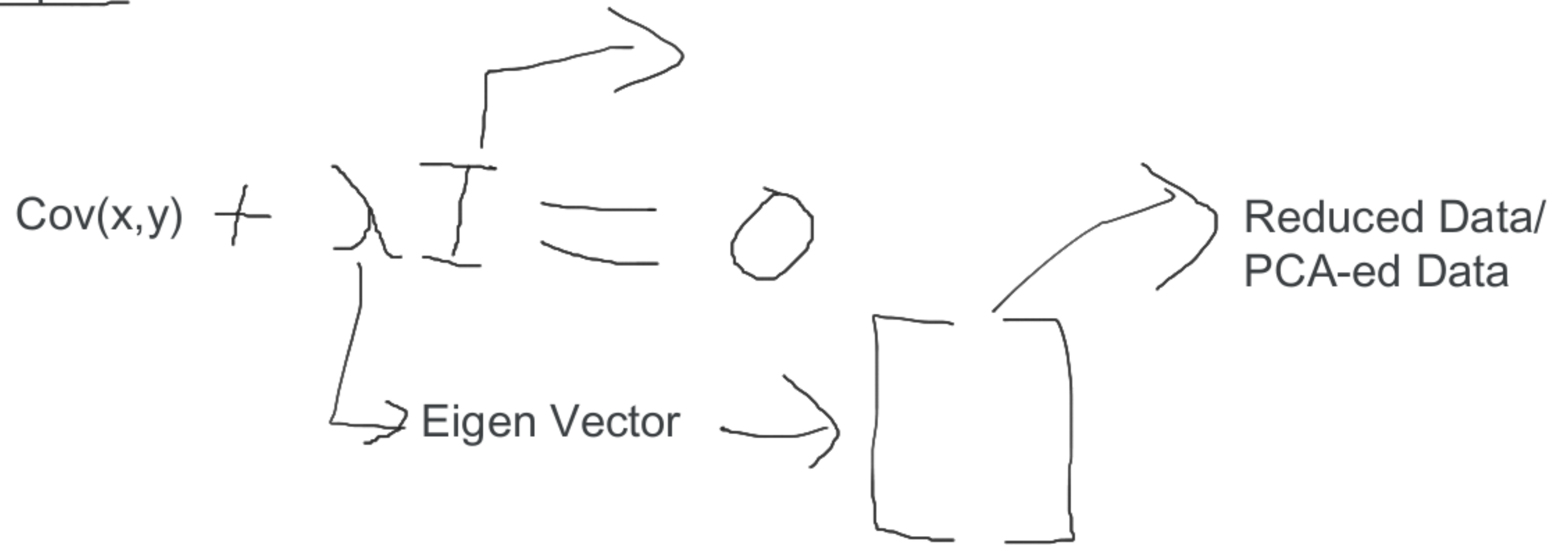Covariance ranges from $-\infty$ to $+\infty$

Covariance Matrix $=$
$$\begin{bmatrix} Var(x) & Cov(x,y) \\ Cov(y,x) & Var(y) \end{bmatrix}$$

Cov(x,y) $\ne$ Eigen Vector = Reduced Data/ PCA-ed Data

a = [1,1,1]
b = np.diag(a)

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

3 × 3

Pipeline

Input(Data)

Scaling of Data by
standardization
technique

PCA

Random
Forest
Classifier

Pipeline avoids data leakage to give us
right/accurate predictions.