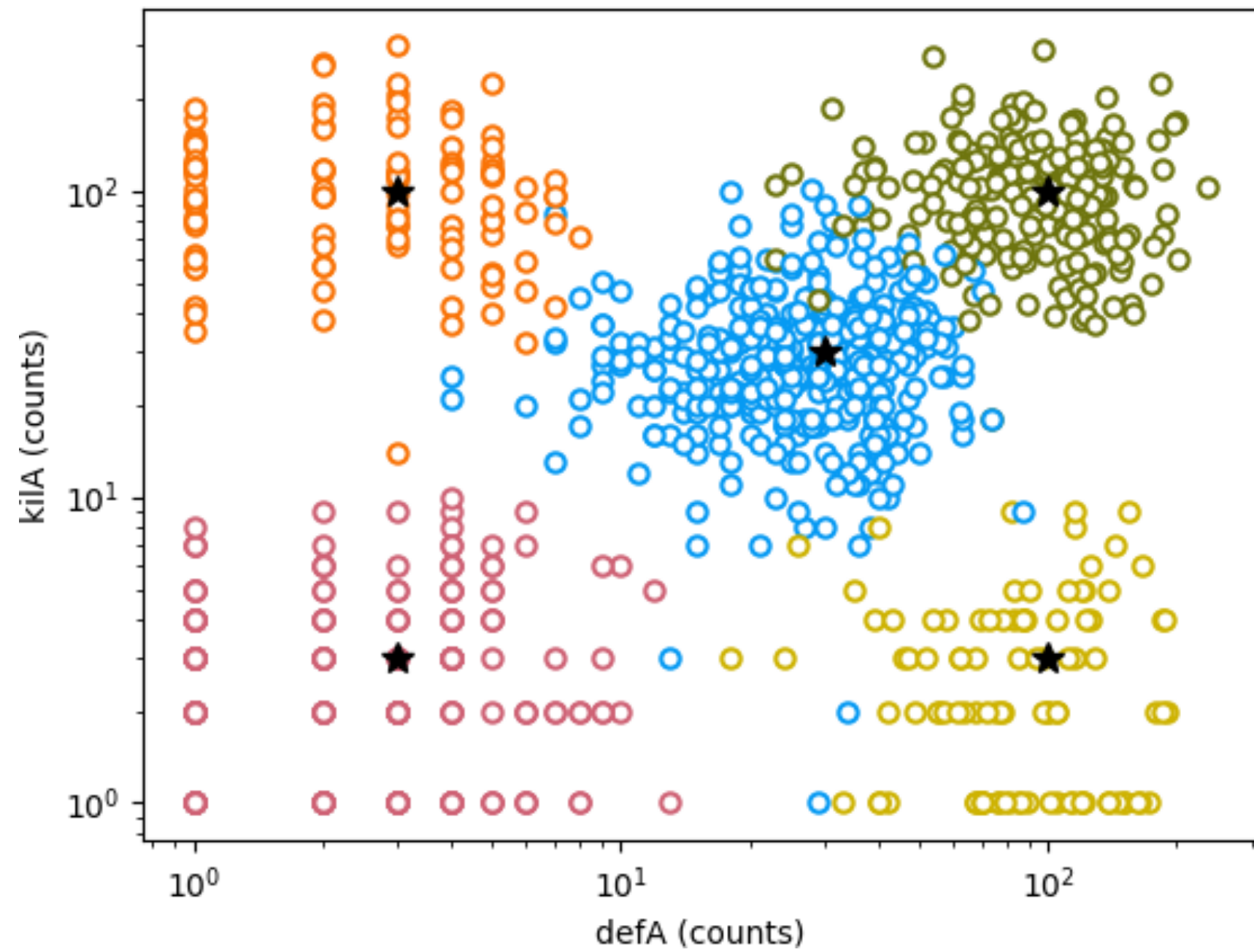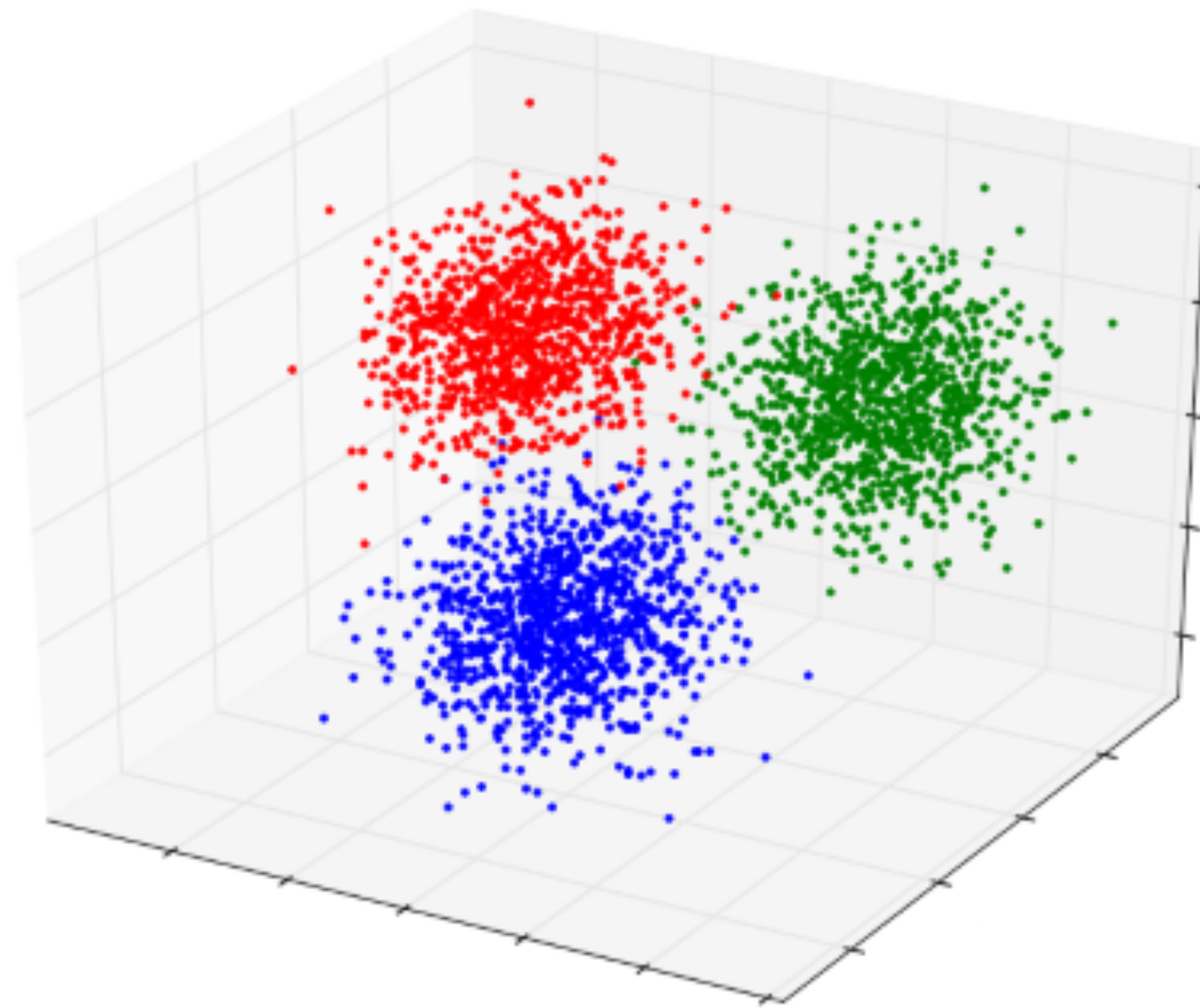# Week 11: PCA

## What is Principal Component Analysis and SVD
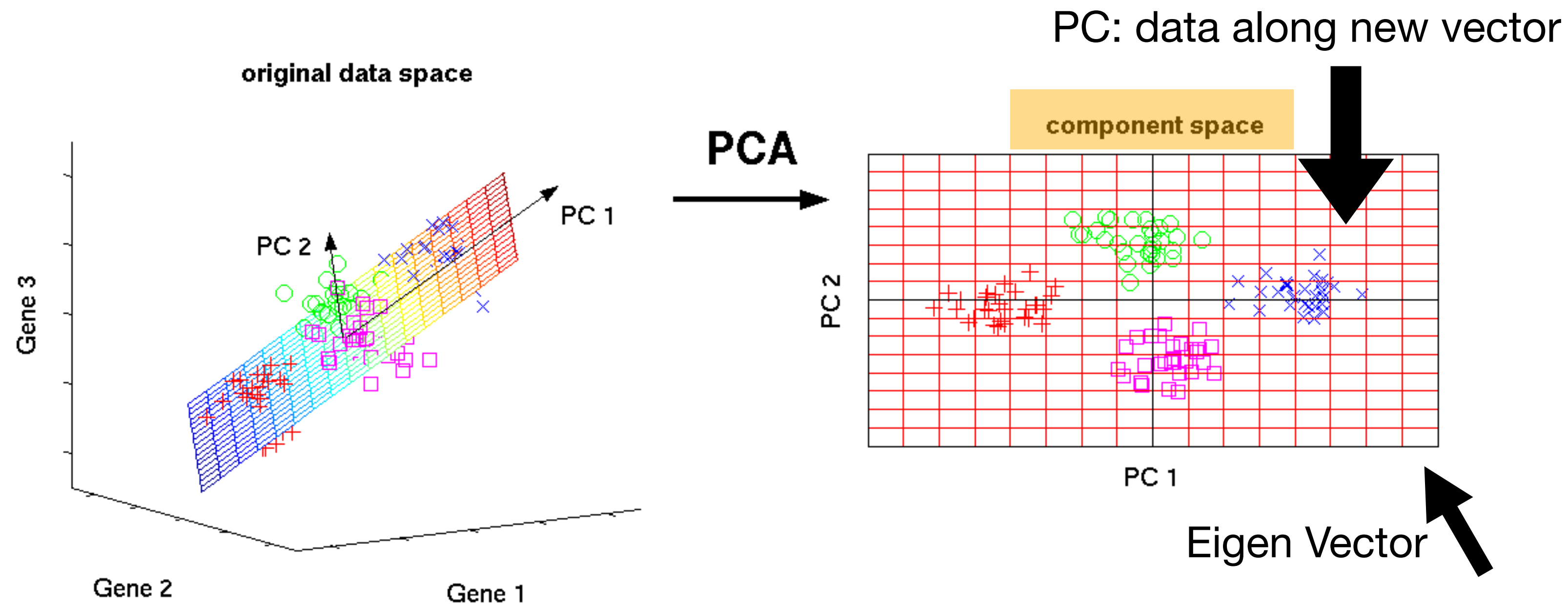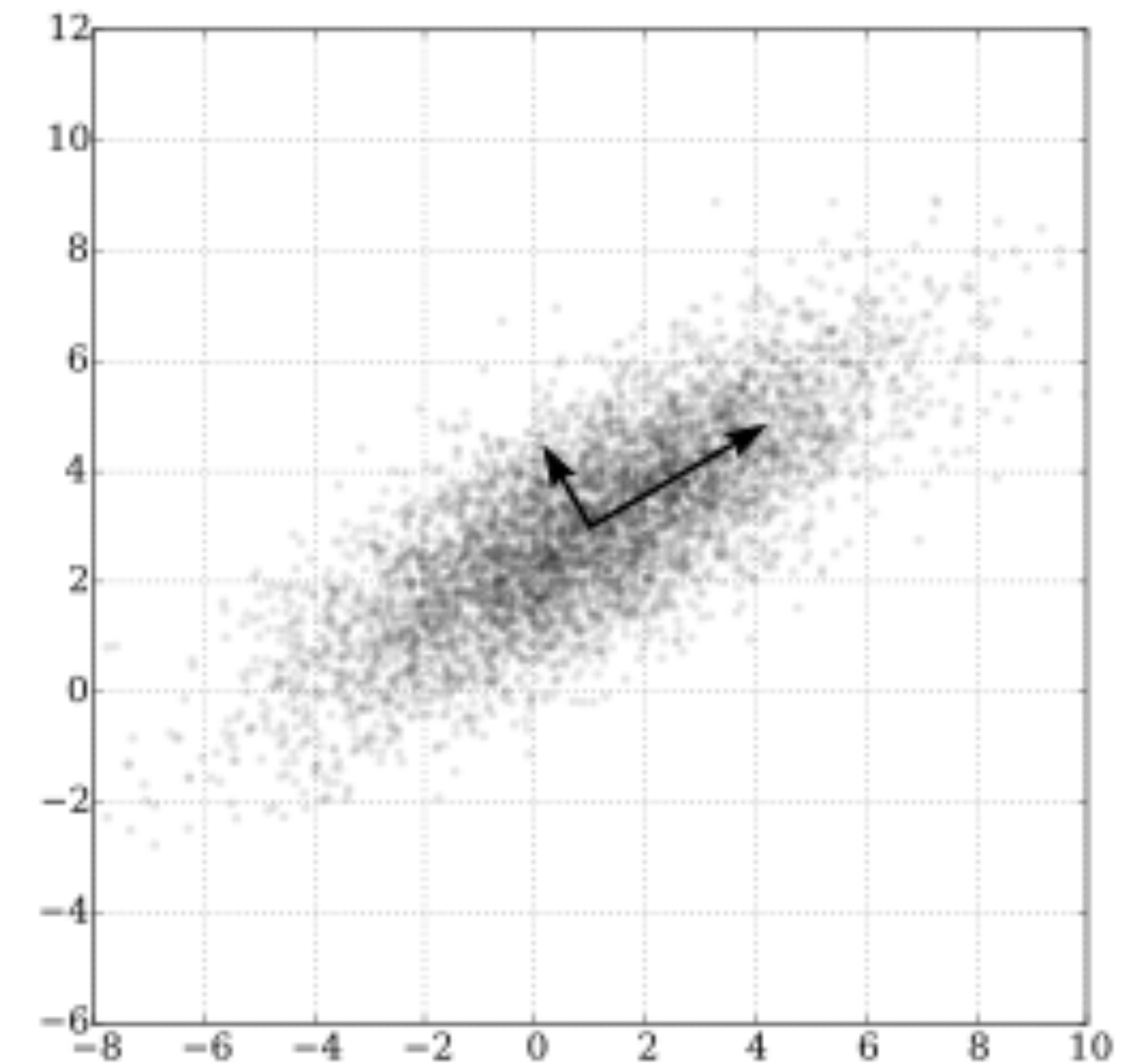
# Higher Dimensional Data



N=2

N=3

N=4

(Each point is a single cell (observation), we are working in 'gene space'

# PCA: A rotation

Multidimensional Gaussian with its 2 PC's marked



original data space

PCA

PC: data along new vector

component space
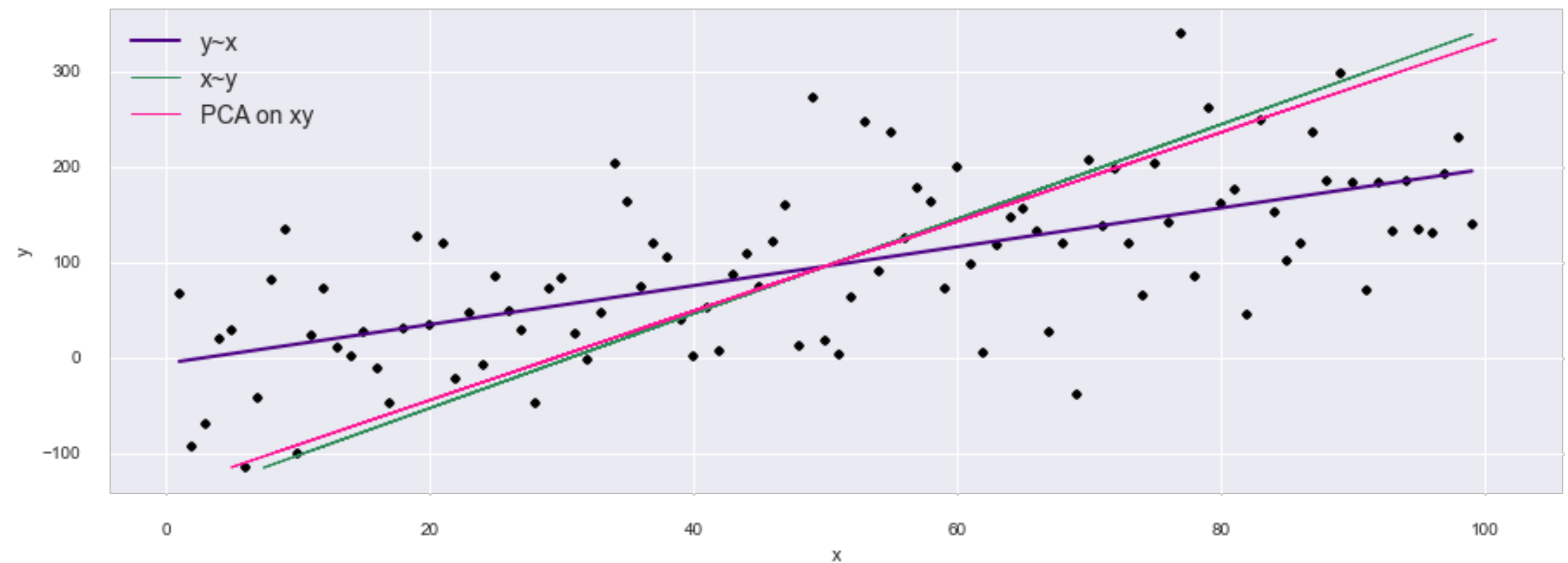
Eigen Vector

Transforms data to a new coordinate system where each axes captures maximum possible variance. We are now in 'component space'.

(Note: Variables covary which gives us information)

# Linear Regression (OLS) vs PC
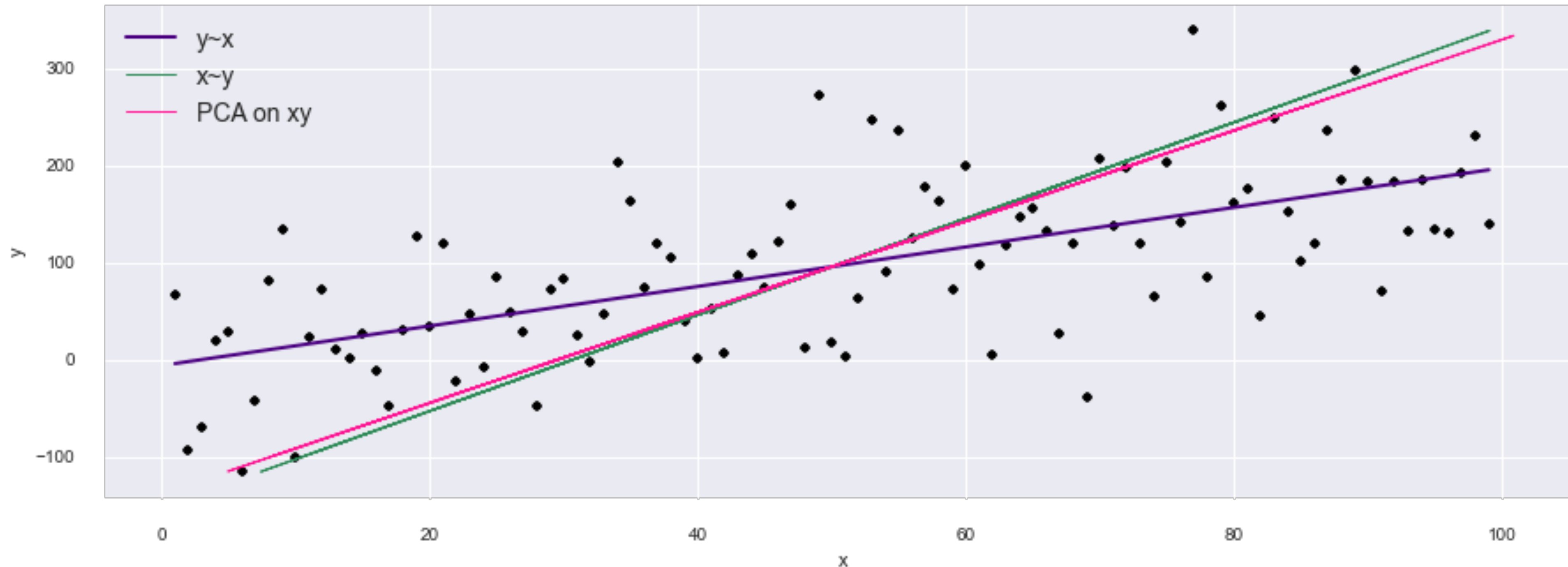


y~x vs. x~y



PCA vs Linear Regression

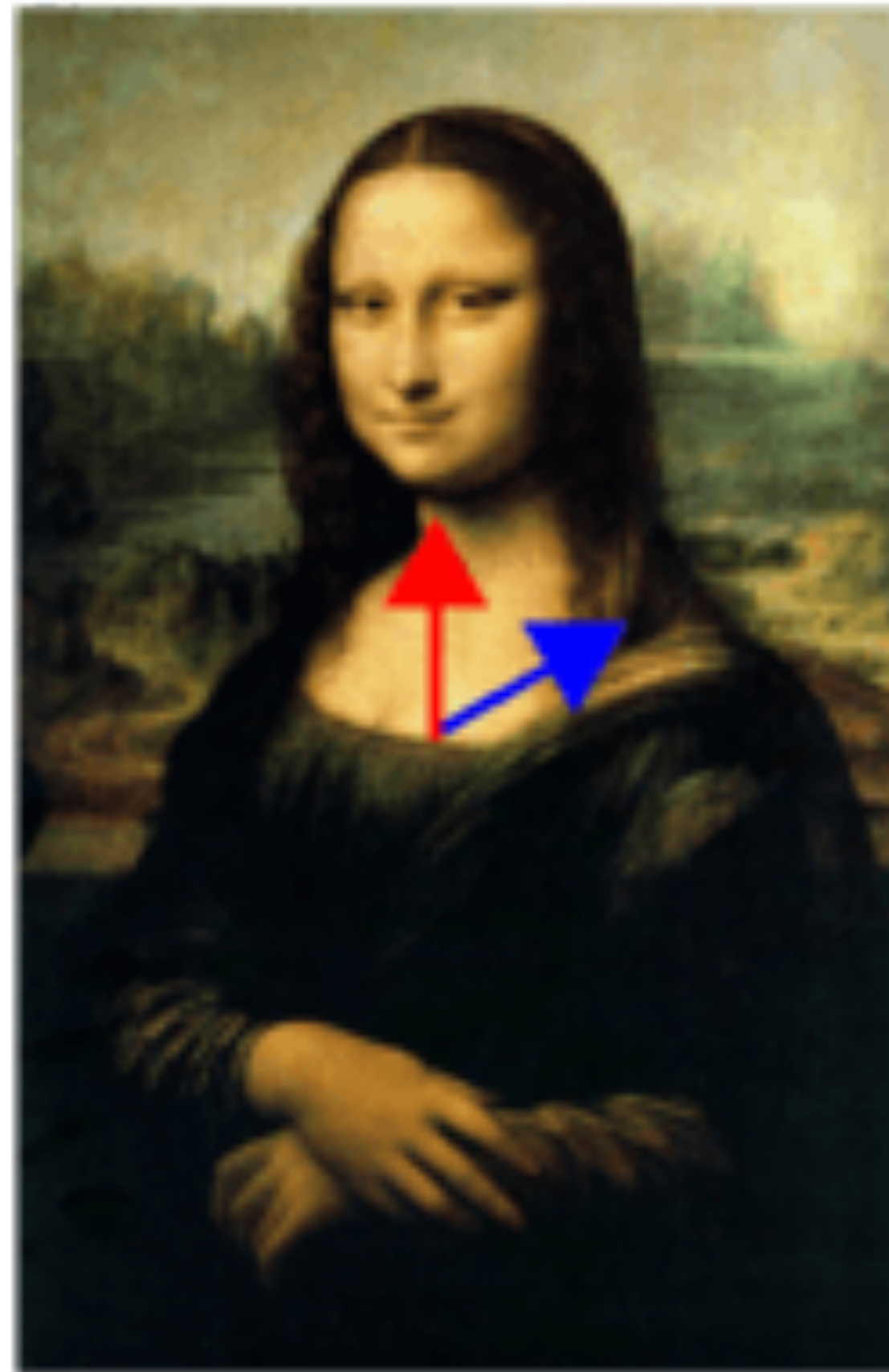# What about covariance across this line?



PCA vs Linear Regression

Maximal Variance along this line (minimal orthogonal to this)

No covariance left along this line (if there is non zero covariance, we can still rotate more to get a better line of best fit, we have not gotten maximal variance along line yet)

# What is an Eigen vector?



Eigen vectors of a covariance matrix provide the direction along which variance varies the most. This is exactly what we want!

Eigen-vectors represent the axes of our new coordinate space. They will no longer change direction in a linear transformation
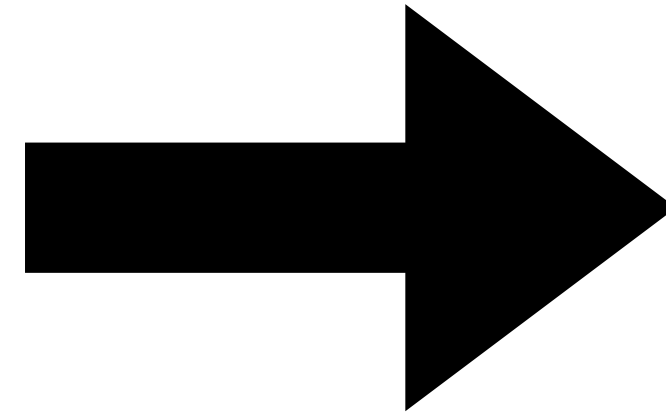
The associated Eigen Values tell us how much variance that particular eigen vector captures.

# Eigen Vectors and Covariance Matrices

Data Matrix: **X** (n points in p dimensions)

(n cells and p genes)

$$\Sigma = \mathbf{W\Lambda W^T}$$

(Eigendecomposition)

Finds eigen values **(W)** with zero covariances

$$\Sigma = \begin{bmatrix} \mathrm{cov}(1,1) & \mathrm{cov}(1,2) & \ldots & \mathrm{cov}(1,p) \\ \mathrm{cov}(2,1) & \mathrm{cov}(2,2) & \ldots & \mathrm{cov}(2,p) \\ \vdots & \vdots & \ddots & \vdots \\ \mathrm{cov}(p,1) & \mathrm{cov}(p,2) & \ldots & \mathrm{cov}(p,p) \end{bmatrix}$$

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 & \ldots & 0 \\ 0 & \lambda_2 & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \ldots & \lambda_p \end{bmatrix}$$

P dimensional data (for ex: p genes)

$\lambda_i$ is eigen value

# Single Value Decomposition and rank

Assume only n observations, but p variables where n<p. Can we fit it?

$$\text{rank of X } r \leq \min(n, p)$$

Can apply SVD only on cantered data

$$x_{ij}^* = x_{ij} - \bar{\mathbf{x}}_j$$

Can decompose it now as:

$$\mathbf{X}^* = \mathbf{U}\mathbf{S}\mathbf{W}^\top.$$

$$\mathbf{W}^\top_{r \text{ x } p}$$

Captures the r principal components, eigen vectors of our covariance matrix

# Plotting data in our lower dimensional space

Assume q PC's are kept

$$\mathbf{Y}_q = \mathbf{X}^*\mathbf{W}_q$$

Data in the new component space

We are rotating the data by our new vectors in **W**