# Principal Component Analysis

# Which features can we ignore?

**Constant:** e.g., number of cabins: 1, 1, 1, 1, 1

**Constant with Noise:** e.g., hair thickness: .008, .003, .005

**Linearly Dependent:** e.g., weight and height

# Variance

# Covariance

$$\text{var}(X) = \frac{\sum_{i=1}^{n} (X_i - \overline{X})(X_i - \overline{X})}{(n-1)}$$

$$\text{cov}(X, Y) = \frac{\sum_{i=1}^{n} (X_i - \overline{X})(Y_i - \overline{Y})}{(n-1)}$$

# Covariance Matrix

$$C = \begin{bmatrix} \text{cov}(X,X) & \text{cov}(X,Y) & \text{cov}(X,Z) \\ \text{cov}(Y,X) & \text{cov}(Y,Y) & \text{cov}(Y,Z) \\ \text{cov}(Z,X) & \text{cov}(Z,Y) & \text{cov}(Z,Z) \end{bmatrix}$$

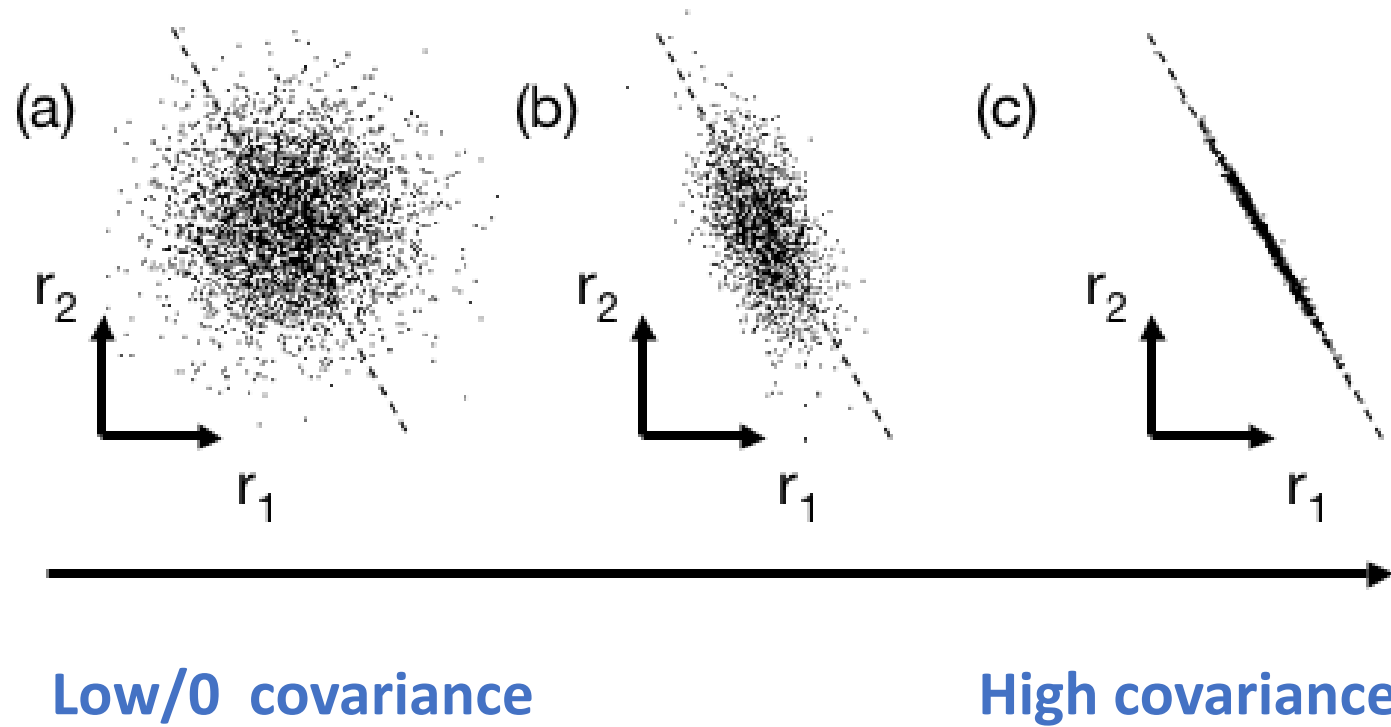# Covariance

**Positive:** both features move together

**Negative:** when one increases the other decreases

**Independent:** no relationship (covariance is 0)
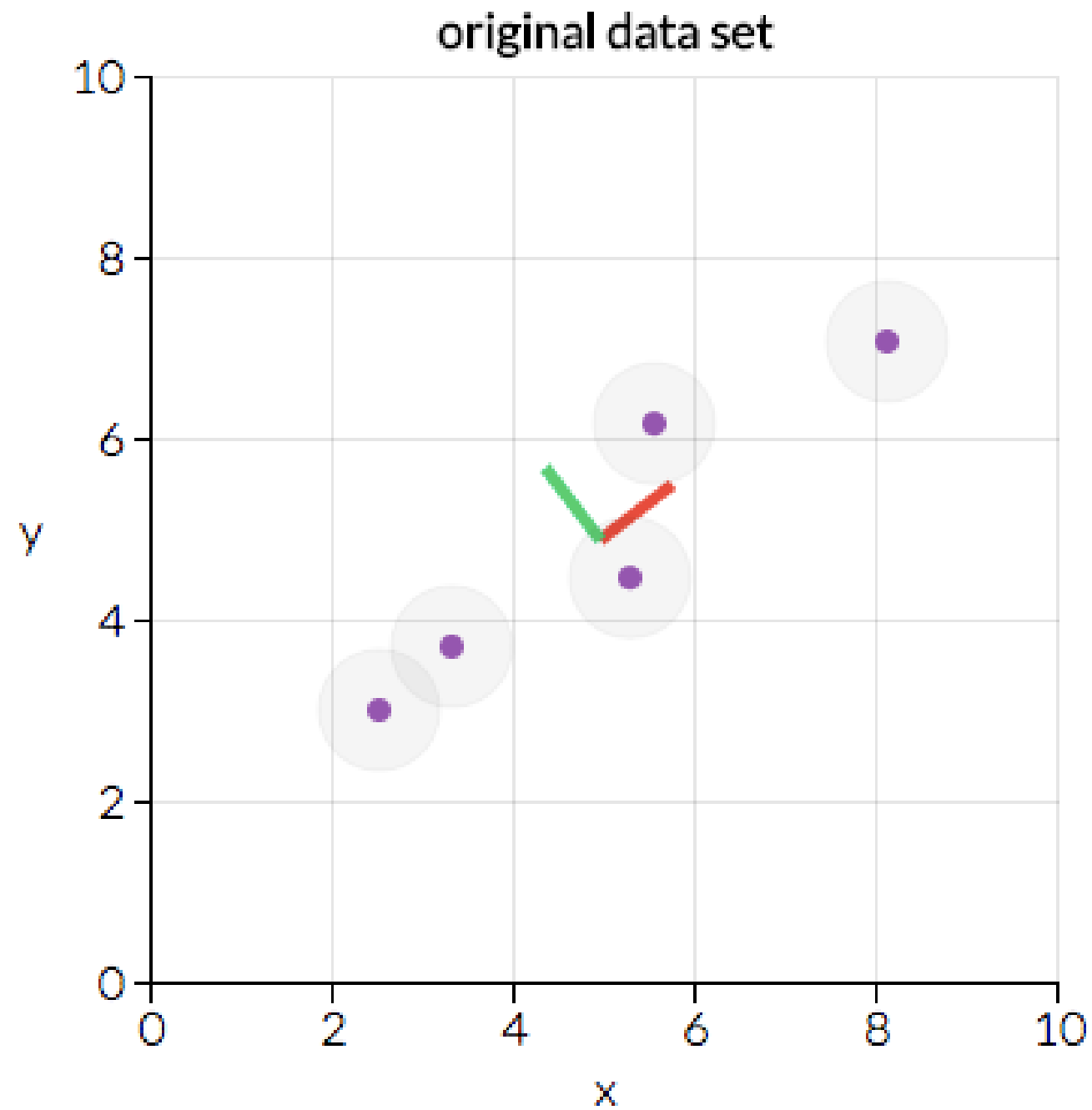
# Covariance

# Which features do we want to keep?

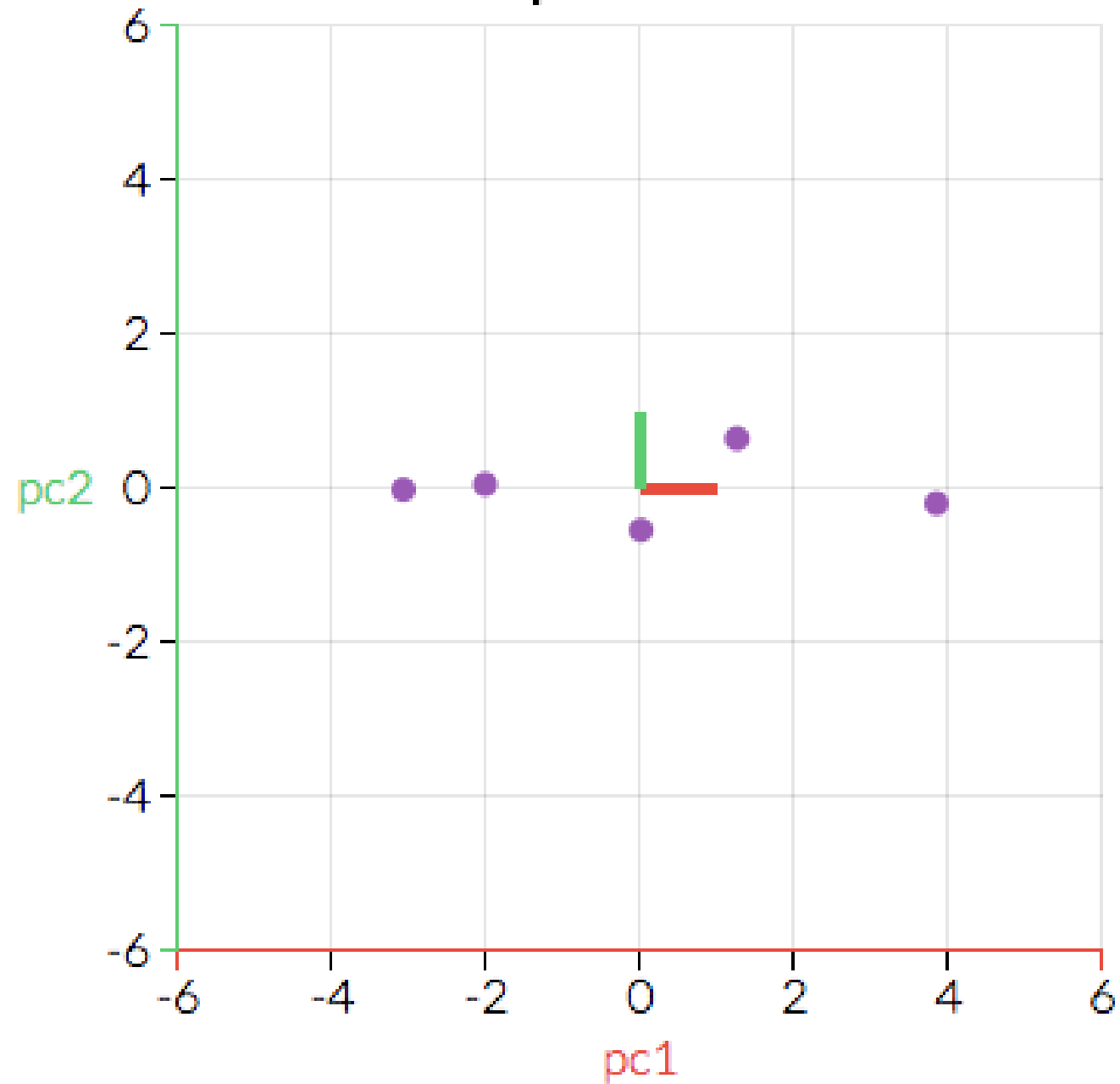**Low covariance:** e.g., hours of sleep the previous night
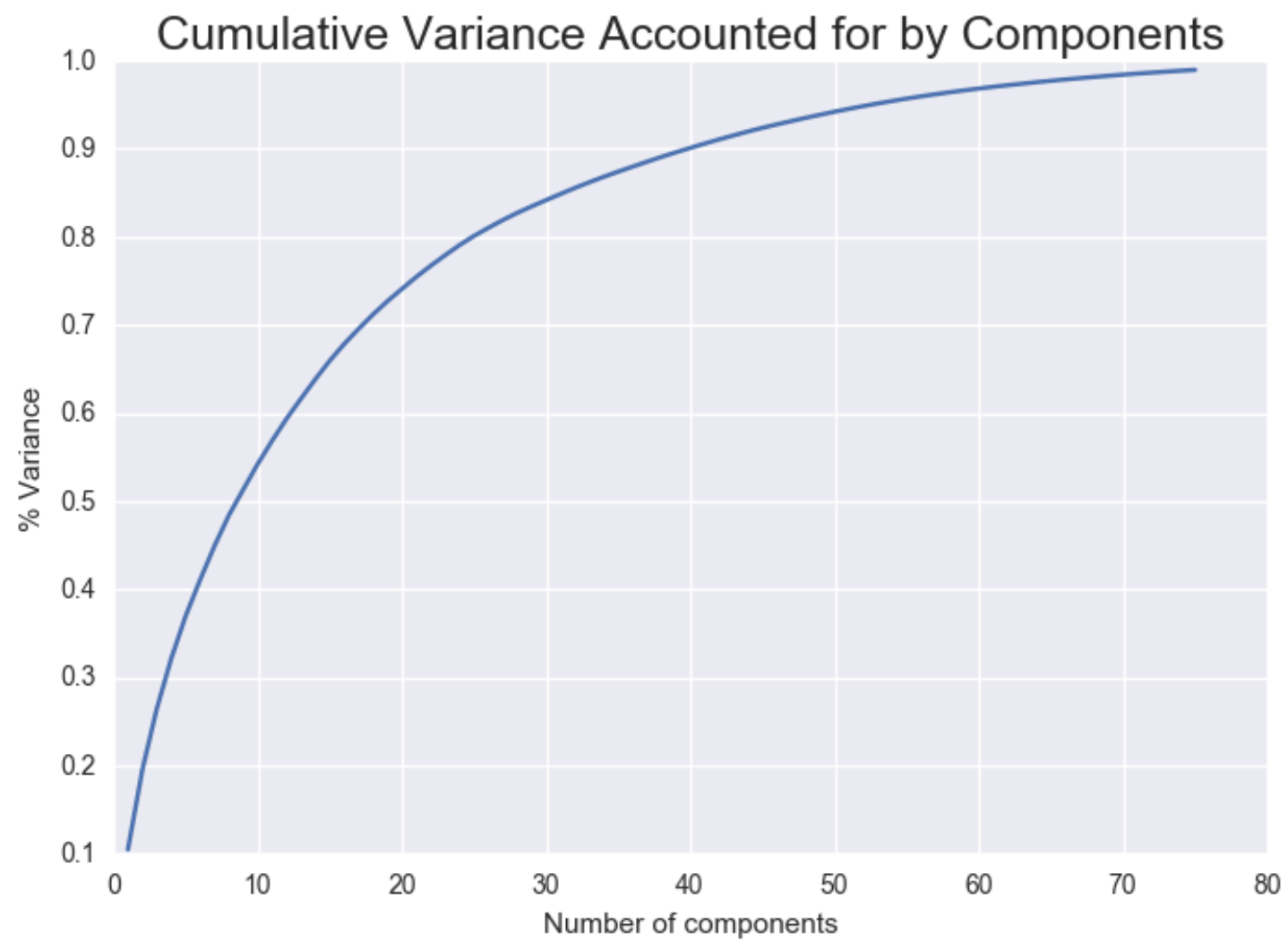
**High variance:** e.g., minutes spent inside cabin

# Eigenvalues

# Eigenvectors

# original data set

output from PCA

Cumulative Variance Accounted for by Components

# PCA is:

## Covariance Matrix

A measure of how each variable is associated with one another.

## Eigenvectors

The directions in which our data are dispersed.

## Eigenvalues

The relative importance (magnitude) of these different directions.