

Covariance

- Variance and Covariance are a measure of the "spread" of a set of points around their center of mass (mean)
- Variance - measure of the deviation from the mean for points in one dimension e.g. heights
- Covariance as a measure of how much **each of the dimensions** vary from the mean **with respect to each other**.
- Covariance is measured between 2 dimensions to see if there is a relationship between the 2 dimensions e.g. number of hours studied & marks obtained.
- The covariance between one dimension and itself is the variance

Covariance

$$\text{covariance}(X,Y) = \frac{\sum_{i=1}^n (\bar{X}_i - \bar{X})(\bar{Y}_i - \bar{Y})}{(n-1)}$$

- So, if you had a 3-dimensional data set (x,y,z) , then you could measure the covariance between the x and y dimensions, the y and z dimensions, and the x and z dimensions. Measuring the covariance between x and x, or y and y, or z and z would give you the variance of the x, y and z dimensions respectively.

Covariance Matrix

- Representing Covariance between dimensions as a matrix e.g. for 3 dimensions:

$$C = \begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) & \text{cov}(x,z) \\ \text{cov}(y,x) & \text{cov}(y,y) & \text{cov}(y,z) \\ \text{cov}(z,x) & \text{cov}(z,y) & \text{cov}(z,z) \end{bmatrix}$$

Variances

- Diagonal is the **variances** of x, y and z
- $\text{cov}(x,y) = \text{cov}(y,x)$ hence matrix is **symmetrical** about the diagonal
- N-dimensional data will result in **NxN covariance** matrix

Covariance

- What is the interpretation of covariance calculations?

e.g.: 2 dimensional data set

x: number of hours studied for a subject

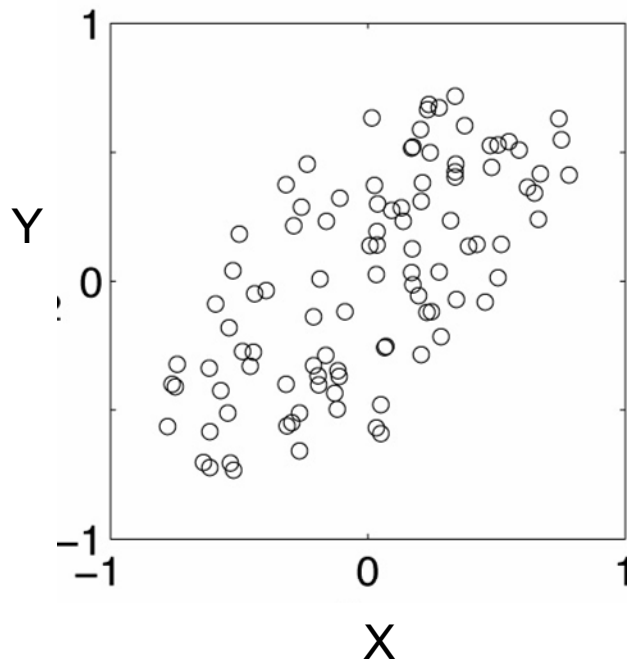
y: marks obtained in that subject

covariance value is say: 104.53

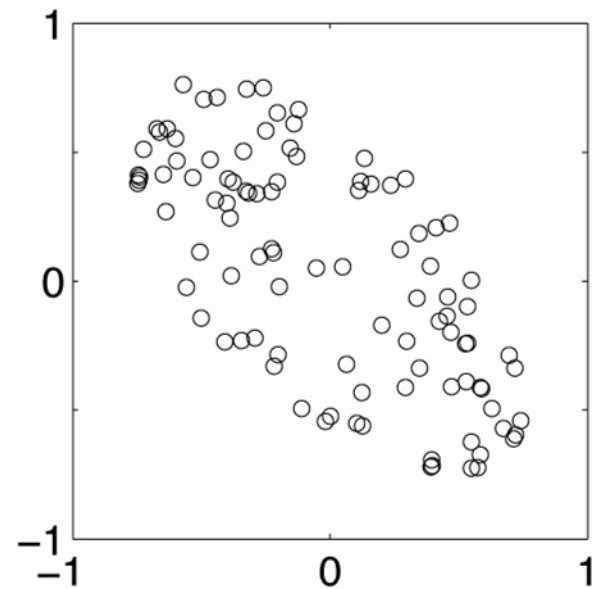
what does this value mean?

Covariance examples

positive covariance



negative covariance



Covariance

- Exact value is not as important as it's sign.
- A positive value of covariance indicates **both dimensions increase or decrease together** e.g. as the number of hours studied increases, the marks in that subject increase.
- A negative value indicates **while one increases the other decreases**, or vice-versa e.g. active social life at PSU vs performance in CS dept.
- If covariance is zero: the two dimensions are **independent** of each other e.g. heights of students vs the marks obtained in a subject

Covariance

- Why bother with calculating covariance when we could just plot the 2 values to see their relationship?

Covariance calculations are used to find relationships between dimensions in high dimensional data sets (usually greater than 3) where visualization is difficult.

PCA

- **principal components analysis (PCA)** is a technique that can be used to **simplify a dataset**
- It is a linear transformation that chooses a new coordinate system for the data set such that
 greatest variance by any projection of the data set comes to lie on the first axis (then called the **first principal component**),
 the second greatest variance on the second axis, and so on.
- PCA can be used for **reducing dimensionality** by eliminating the later principal components.

PCA Toy Example

Consider the following 3D points

1	2	4	3	5	6
2	4	8	6	10	12
3	6	12	9	15	18

If each component is stored in a byte,
we need $18 = 3 \times 6$ bytes

PCA Toy Example

Looking closer, we can see that all the points are related geometrically: they are all the same point, scaled by a factor:

<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3	= 1 *	<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3	<table><tr><td>2</td></tr><tr><td>4</td></tr><tr><td>6</td></tr></table>	2	4	6	= 2 *	<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3	<table><tr><td>4</td></tr><tr><td>8</td></tr><tr><td>12</td></tr></table>	4	8	12	= 4 *	<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3
1																										
2																										
3																										
1																										
2																										
3																										
2																										
4																										
6																										
1																										
2																										
3																										
4																										
8																										
12																										
1																										
2																										
3																										
<table><tr><td>3</td></tr><tr><td>6</td></tr><tr><td>9</td></tr></table>	3	6	9	= 3 *	<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3	<table><tr><td>5</td></tr><tr><td>10</td></tr><tr><td>15</td></tr></table>	5	10	15	= 5 *	<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3	<table><tr><td>6</td></tr><tr><td>12</td></tr><tr><td>18</td></tr></table>	6	12	18	= 6 *	<table><tr><td>1</td></tr><tr><td>2</td></tr><tr><td>3</td></tr></table>	1	2	3
3																										
6																										
9																										
1																										
2																										
3																										
5																										
10																										
15																										
1																										
2																										
3																										
6																										
12																										
18																										
1																										
2																										
3																										

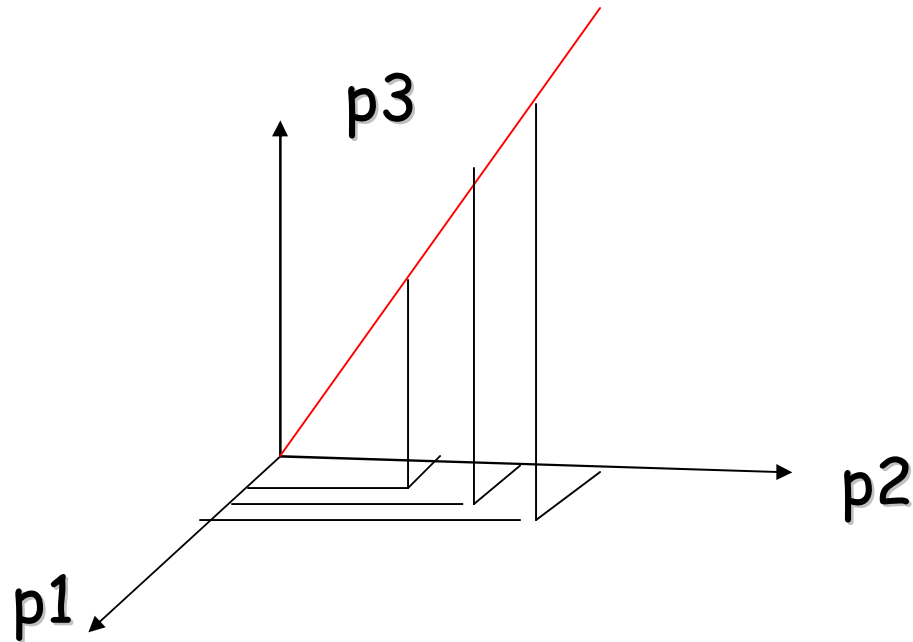
PCA Toy Example

1		1		2		1		4		1
2	= 1 *	2		4	= 2 *	2		8	= 4 *	2
3		3		6		3		12		3
3		1		5		1		6		1
6	= 3 *	2		10	= 5 *	2		12	= 6 *	2
9		3		15		3		18		3

They can be stored using only 9 bytes (50% savings!):
Store one point (3 bytes) + the multiplying constants (6 bytes)

Geometrical Interpretation:

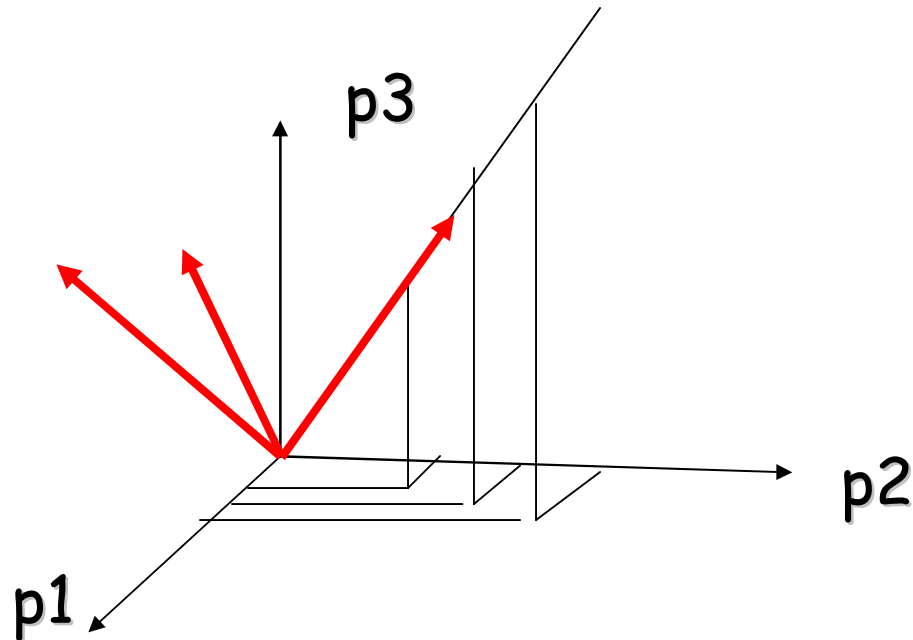
View each point in 3D space.



But in this example, all the points happen to belong to a line: a 1D subspace of the original 3D space.

Geometrical Interpretation:

Consider a new coordinate system where one of the axes is along the direction of the line:



In this coordinate system, every point has only one non-zero coordinate: we only need to store the direction of the line (a 3 bytes image) and the non-zero coordinate for each of the points (6 bytes).

Principal Component Analysis (PCA)

- Given a set of points, how do we know if they can be compressed like in the previous example?
 - The answer is to look into the correlation between the points
 - The tool for doing this is called PCA

PCA

- By finding the **eigenvalues and eigenvectors** of the covariance matrix, we find that the eigenvectors with the largest eigenvalues correspond to the dimensions that have the strongest correlation in the dataset.
- This is the principal component.
- PCA is a useful statistical technique that has found application in:
 - fields such as face recognition and image compression
 - finding patterns in data of high dimension.

PCA Theorem

Let $x_1 x_2 \dots x_n$ be a set of n $N \times 1$ vectors and let \bar{x} be their average:

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{iN} \end{bmatrix} \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{iN} \end{bmatrix}$$

PCA Theorem

Let X be the $N \times n$ matrix with columns

$x_1 - \bar{x}, x_2 - \bar{x}, \dots, x_n - \bar{x}$:

$$X = \begin{bmatrix} x_1 - \bar{x} & x_2 - \bar{x} & \cdots & x_n - \bar{x} \end{bmatrix}$$

Note: subtracting the mean is equivalent to translating the coordinate system to the location of the mean.

PCA Theorem

Let $Q = X X^T$ be the $N \times N$ matrix:

$$Q = X X^T = \begin{bmatrix} \mathbf{x}_1 - \bar{\mathbf{x}} & \mathbf{x}_2 - \bar{\mathbf{x}} & \cdots & \mathbf{x}_n - \bar{\mathbf{x}} \end{bmatrix} \begin{bmatrix} (\mathbf{x}_1 - \bar{\mathbf{x}})^T \\ (\mathbf{x}_2 - \bar{\mathbf{x}})^T \\ \vdots \\ (\mathbf{x}_n - \bar{\mathbf{x}})^T \end{bmatrix}$$

Notes:

1. Q is square
2. Q is symmetric
3. Q is the covariance matrix [aka scatter matrix]
4. Q can be very large (in vision, N is often the number of pixels in an image!)

PCA Theorem

Theorem:

Each x_j can be written as:
$$x_j = \bar{x} + \sum_{i=1}^{i=n} g_{ji} e_i$$

where e_i are the n eigenvectors of Q with non-zero eigenvalues.

Notes:

1. The eigenvectors $e_1 e_2 \dots e_n$ span an *eigenspace*
2. $e_1 e_2 \dots e_n$ are $N \times 1$ orthonormal vectors (directions in N -Dimensional space)
3. The scalars g_{ji} are the coordinates of x_j in the space.

$$g_{ji} = (x_j - \bar{x}) \cdot e_i$$

Using PCA to Compress Data

- Expressing x in terms of $e_1 \dots e_n$ has not changed the size of the data
- However, if the points are highly correlated many of the coordinates of x will be zero or closed to zero.

note: this means they lie in a lower-dimensional linear subspace

Using PCA to Compress Data

- Sort the eigenvectors e_i according to their eigenvalue:

$$\lambda_1 \geq \lambda_2 \geq \dots \lambda_n$$

- Assuming that $\lambda_i \approx 0$ if $i > k$

- Then

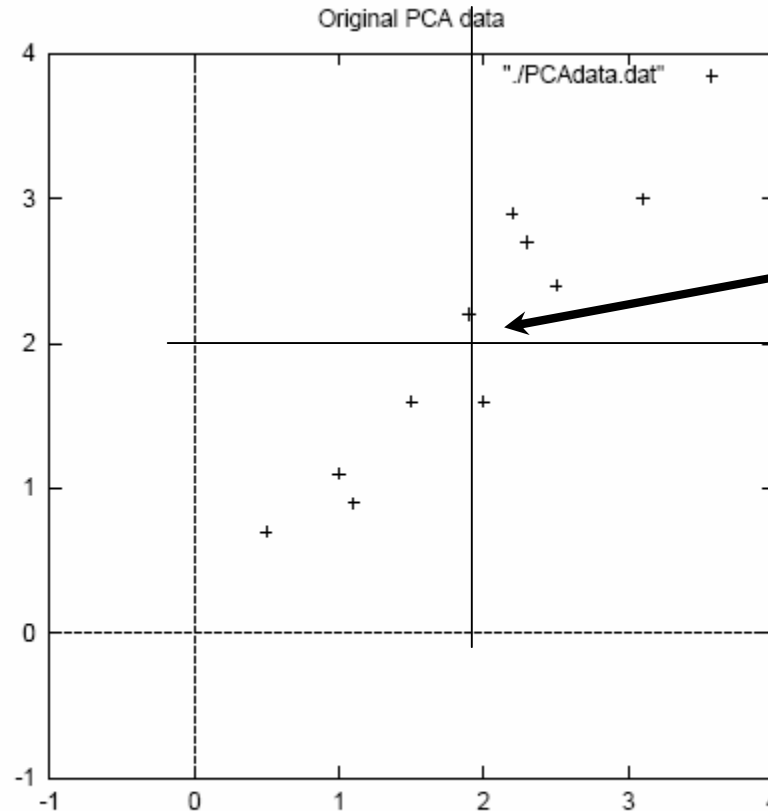
$$\mathbf{x}_j \approx \bar{\mathbf{x}} + \sum_{i=1}^{i=k} g_{ji} \mathbf{e}_i$$

PCA Example -STEP 1

<http://kybele.psych.cornell.edu/~edelman/Psych-465-Spring-2003/PCA-tutorial.pdf>

- DATA:

x	y
2.5	2.4
0.5	0.7
2.2	2.9
1.9	2.2
3.1	3.0
2.3	2.7
2	1.6
1	1.1
1.5	1.6
1.1	0.9



mean

this becomes the
new origin of the
data from now on

Figure 3.1: PCA example data, original data on the left, data with the means subtracted on the right, and a plot of the data

PCA Example -STEP 2

- Calculate the covariance matrix

$$\text{cov} = \begin{pmatrix} .616555556 & .615444444 \\ .615444444 & .716555556 \end{pmatrix}$$

- since the non-diagonal elements in this covariance matrix are positive, we should expect that both the x and y variable increase together.

PCA Example -STEP 3

- Calculate the eigenvectors and eigenvalues of the covariance matrix

$$\text{eigenvalues} = \begin{pmatrix} .0490833989 \\ 1.28402771 \end{pmatrix}$$

$$\text{eigenvectors} = \begin{pmatrix} -.735178656 & -.677873399 \\ .677873399 & -.735178656 \end{pmatrix}$$

PCA Example -STEP 3

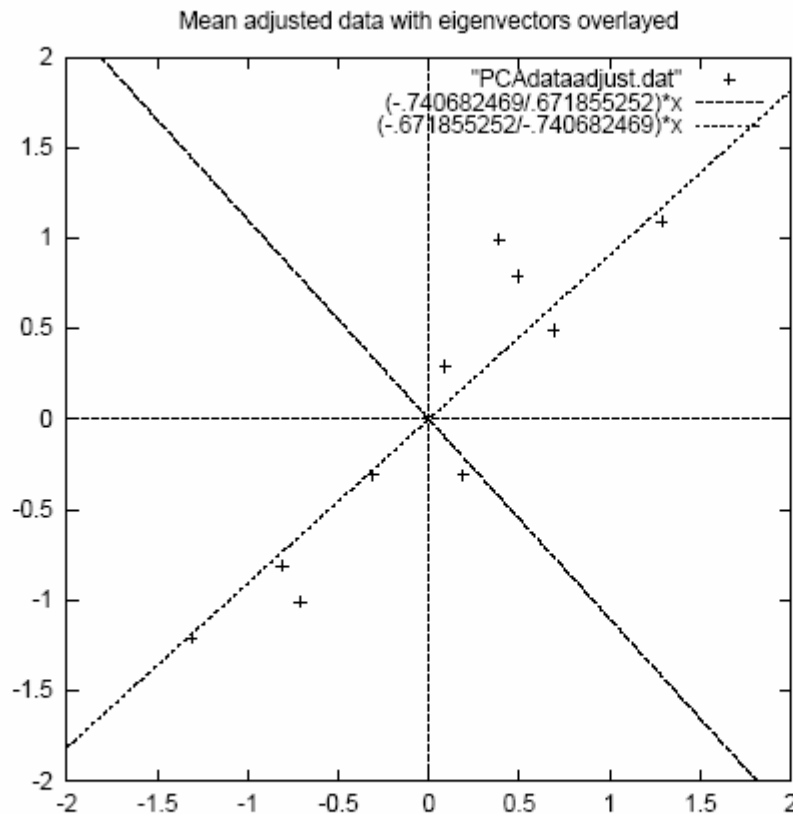


Figure 3.2: A plot of the normalised data (mean subtracted) with the eigenvectors of the covariance matrix overlayed on top.

- eigenvectors are plotted as diagonal dotted lines on the plot.
- Note they are perpendicular to each other.
- Note one of the eigenvectors goes through the middle of the points, like drawing a line of best fit.
- The second eigenvector gives us the other, less important, pattern in the data, that all the points follow the main line, but are off to the side of the main line by some amount.

PCA Example -STEP 4

- Feature Vector

$$\text{FeatureVector} = (\text{eig}_1 \text{ eig}_2 \text{ eig}_3 \dots \text{eig}_n)$$

We can either form a feature vector with both of the eigenvectors:

$$\begin{pmatrix} -.677873399 & -.735178656 \\ -.735178656 & .677873399 \end{pmatrix}$$

or, we can choose to leave out the smaller, less significant component and only have a single column:

$$\begin{pmatrix} -.677873399 \\ -.735178656 \end{pmatrix}$$

PCA Example -STEP 5

- Deriving new data coordinates

$\text{FinalData} = \text{RowFeatureVector} \times \text{RowZeroMeanData}$

RowFeatureVector is the matrix with the eigenvectors in the columns *transposed* so that the eigenvectors are now in the rows, with the most significant eigenvector at the top

RowZeroMeanData is the mean-adjusted data *transposed*, ie. the data items are in each column, with each row holding a separate dimension.

Note: this is essential Rotating the coordinate axes so higher-variance axes come first.

PCA Example -STEP 5

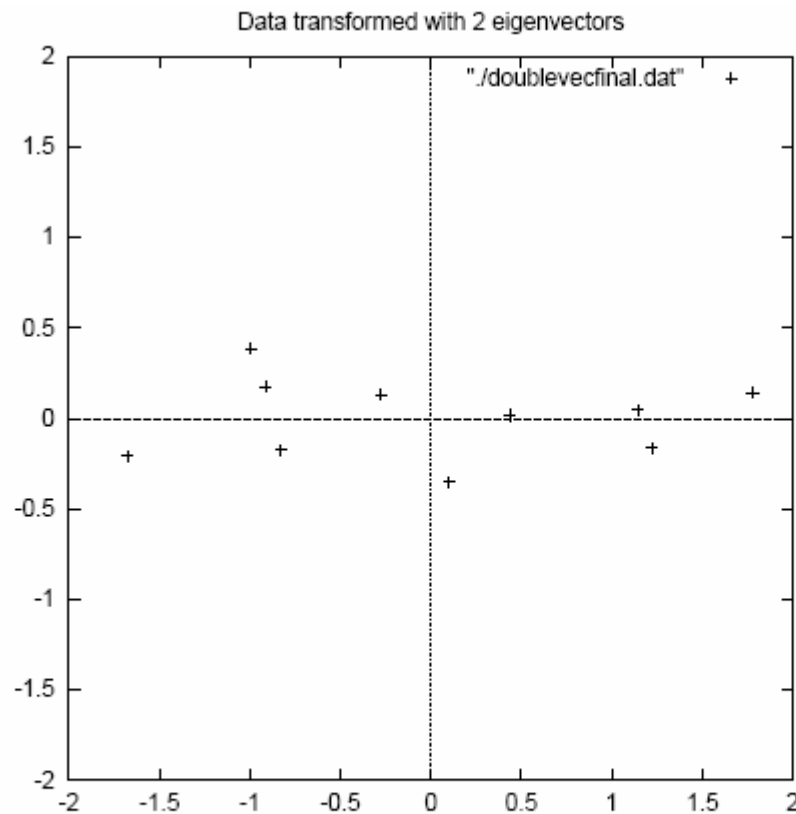
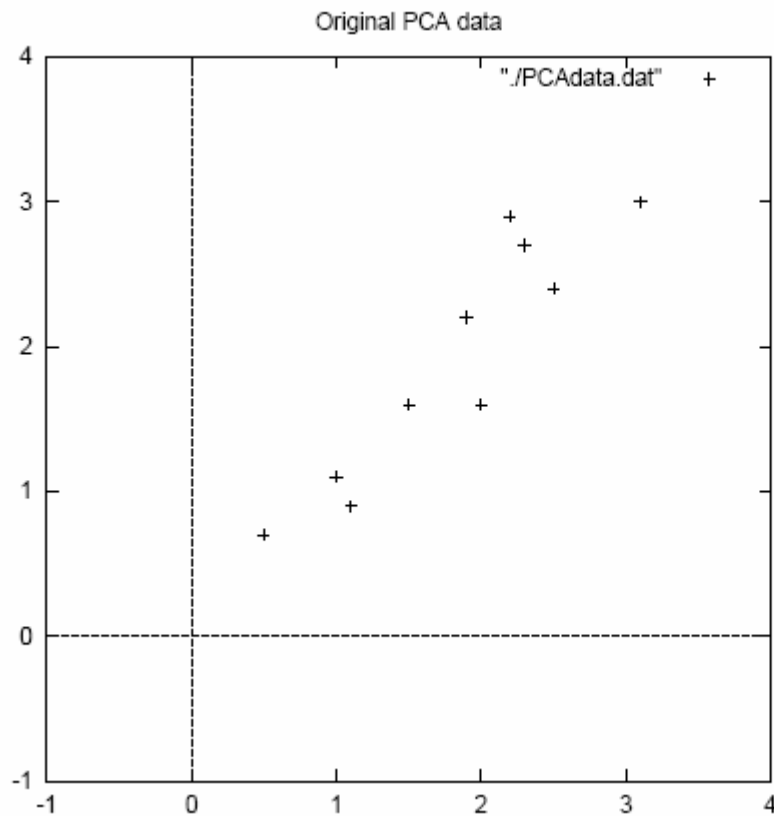


Figure 3.3: The table of data by applying the PCA analysis using both eigenvectors, and a plot of the new data points.

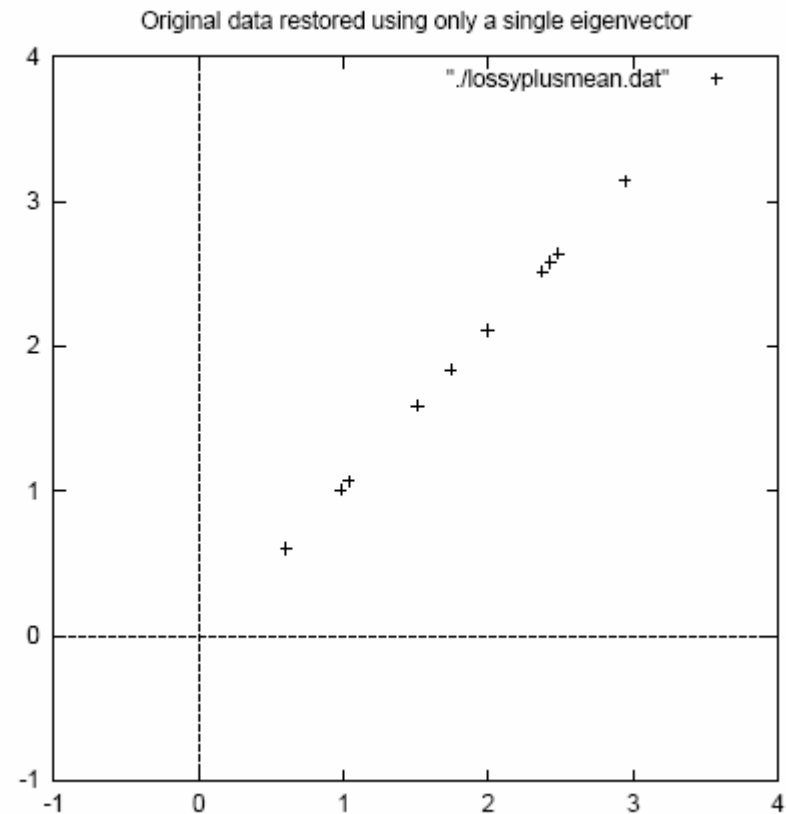
PCA Example : Approximation

- If we reduced the dimensionality, obviously, when reconstructing the data we would lose those dimensions we chose to discard. In our example let us assume that we considered only the x dimension...

PCA Example : Final Approximation



2D point cloud

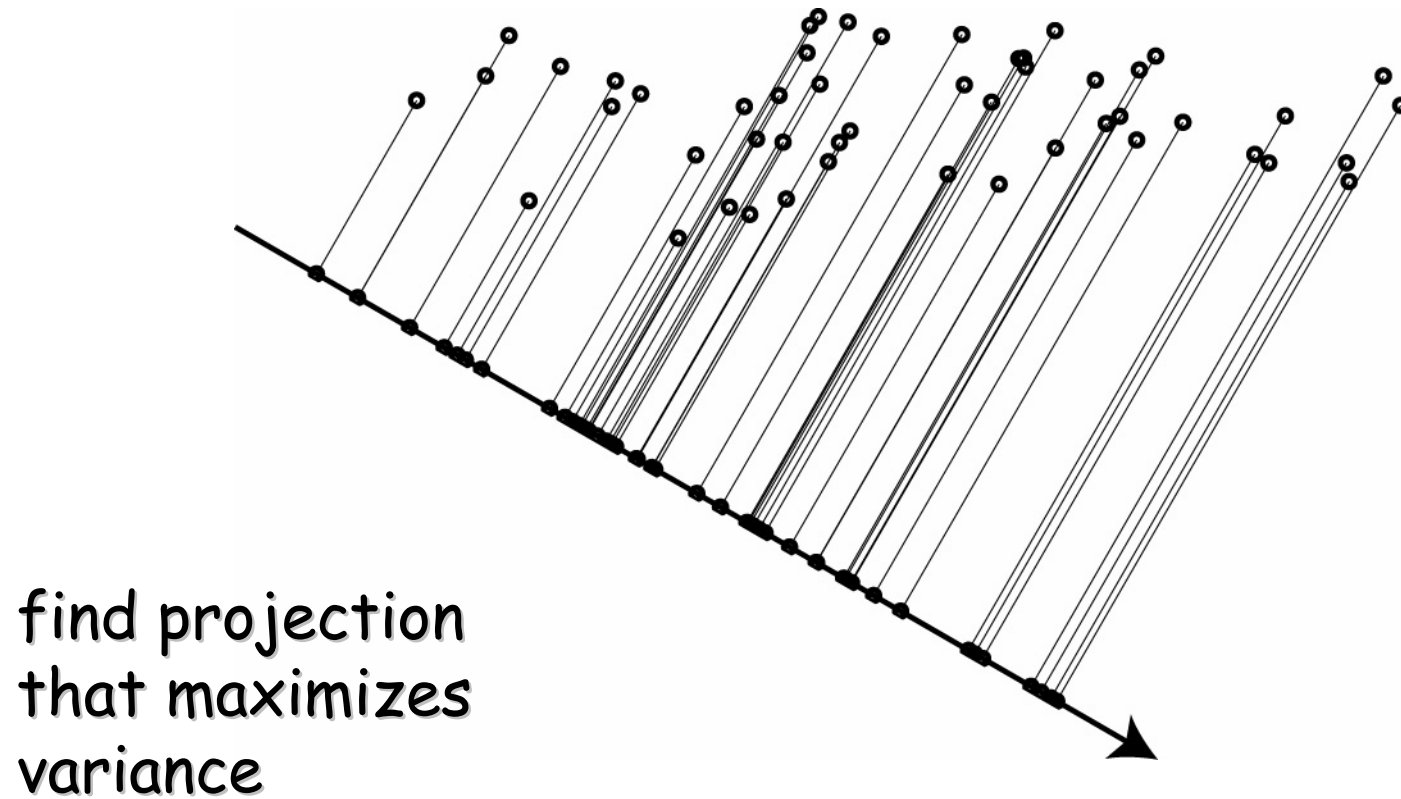


Approximation using
one eigenvector basis

Another way of thinking about Principal component

- direction of maximum variance in the input space
- happens to be same as the principal eigenvector of the covariance matrix

One-dimensional projection



Covariance to variance

- From the covariance, the variance of any projection can be calculated.
- Let w be a unit vector

$$\begin{aligned}\left\langle \left(w^T x\right)^2 \right\rangle - \left\langle w^T x \right\rangle^2 &= w^T C w \\ &= \sum_{ij} w_i C_{ij} w_j\end{aligned}$$

Maximizing variance

- Principal eigenvector of C
 - the one with the largest eigenvalue.

$$w^* = \arg \max_{w: |w|=1} w^T C w$$

$$\begin{aligned} \lambda_{\max}(C) &= \max_{w: |w|=1} w^T C w \\ &= w^{*T} C w^* \end{aligned}$$

Implementing PCA

- Need to find "first" k eigenvectors of Q:

$$Q = XX^T = \begin{bmatrix} \mathbf{x}_1 - \bar{\mathbf{x}} & \mathbf{x}_2 - \bar{\mathbf{x}} & \cdots & \mathbf{x}_n - \bar{\mathbf{x}} \end{bmatrix} \begin{bmatrix} (\mathbf{x}_1 - \bar{\mathbf{x}})^T \\ (\mathbf{x}_2 - \bar{\mathbf{x}})^T \\ \vdots \\ (\mathbf{x}_n - \bar{\mathbf{x}})^T \end{bmatrix}$$

Q is N x N (Again, N could be the number of pixels in an image. For a 256 x 256 image, N = 65536 !!)
Don't want to explicitly compute Q!!!!

Singular Value Decomposition (SVD)

Any $m \times n$ matrix X can be written as the product of 3 matrices:

$$X = UDV^T$$

Where:

- U is $m \times m$ and its columns are orthonormal vectors
- V is $n \times n$ and its columns are orthonormal vectors
- D is $m \times n$ diagonal and its diagonal elements are called the singular values of X , and are such that:

$$\sigma_1, \sigma_2, \dots, \sigma_n, 0$$

SVD Properties

$$X = UDV^T$$

- The columns of U are the eigenvectors of XX^T
- The columns of V are the eigenvectors of X^TX
- The squares of the diagonal elements of D are the eigenvalues of XX^T and X^TX