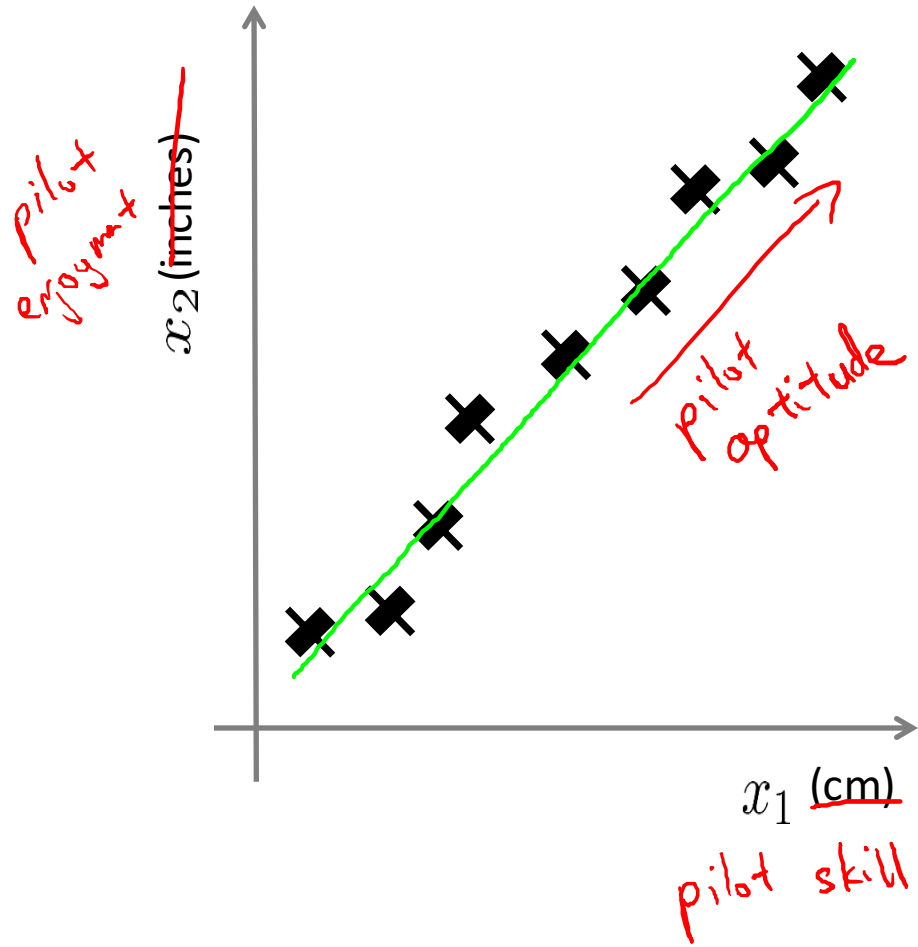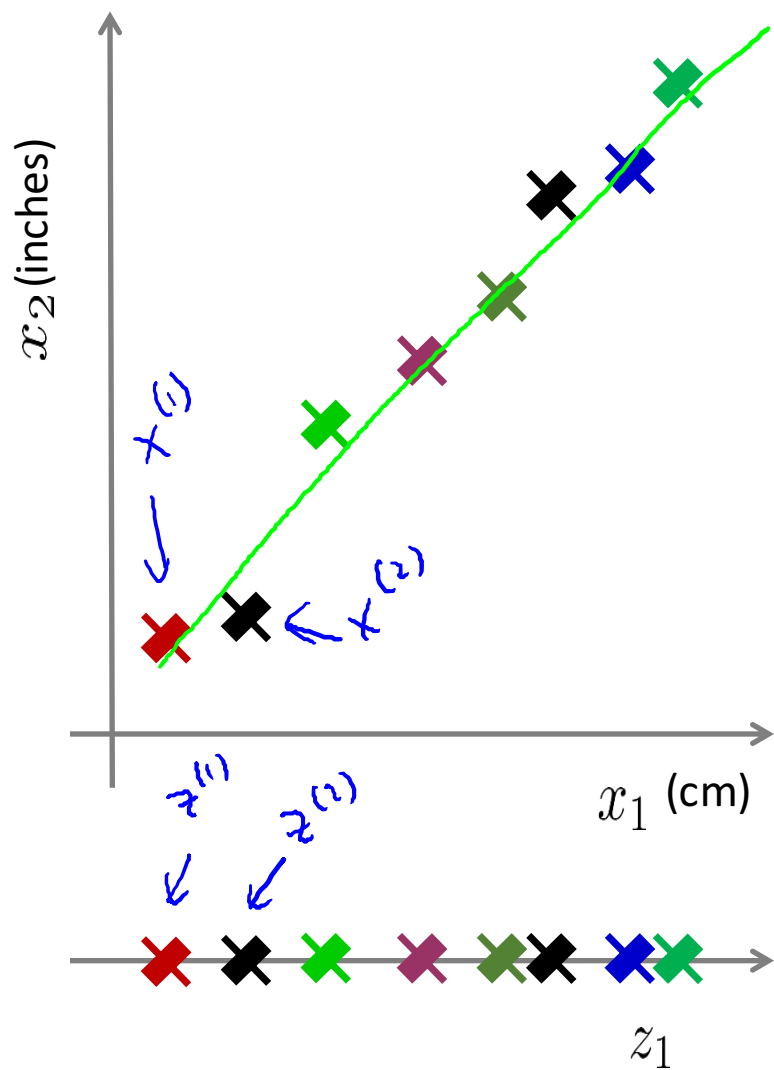# Principal Components Analysis (PCA)

# PCA

- Not all features (attributes) collected are independent.

- In fact, many attributes are highly correlated
  e.g. height and weight, total height and head height, attendance and GPA

- We can compact the data by choosing a new attribute that is a linear combination of the correlated factors

# Data Compression



Reduce data from 2D to 1D

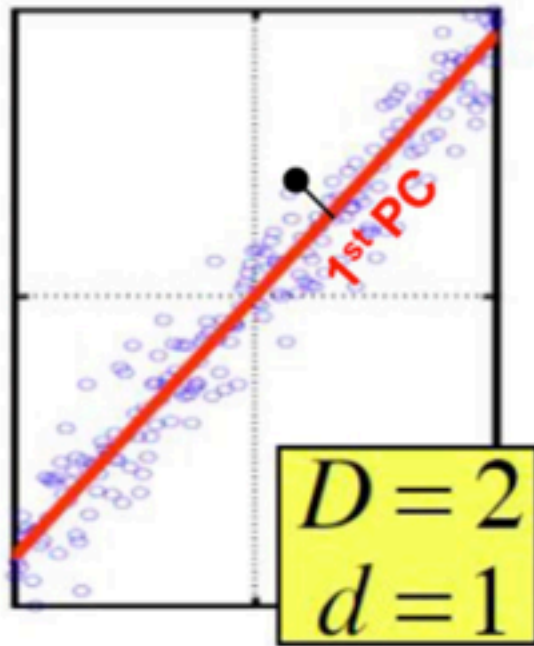# Data Compression



Reduce data from 2D to 1D

$$x^{(1)} \in \mathbb{R}^2 \rightarrow z^{(1)} \in \mathbb{R}$$

$$x^{(2)} \in \mathbb{R}^2 \rightarrow z^{(2)} \in \mathbb{R}$$

$$\vdots$$

$$x^{(m)} \in \mathbb{R}^2 \rightarrow z^{(m)} \in \mathbb{R}$$

# Principal components analysis



Principal Components (PC) are orthogonal directions that capture most of the variance in the data

1st PC – direction of greatest variability in data

$$D = 2$$
$$d = 1$$

# Principal components analysis

Principal Components (PC) are orthogonal directions that capture most of the variance in the data

1st PC – direction of greatest variability in data

$D = 2$
$d = 1$

Take a data point $x_i$ (D-dimensional vector)
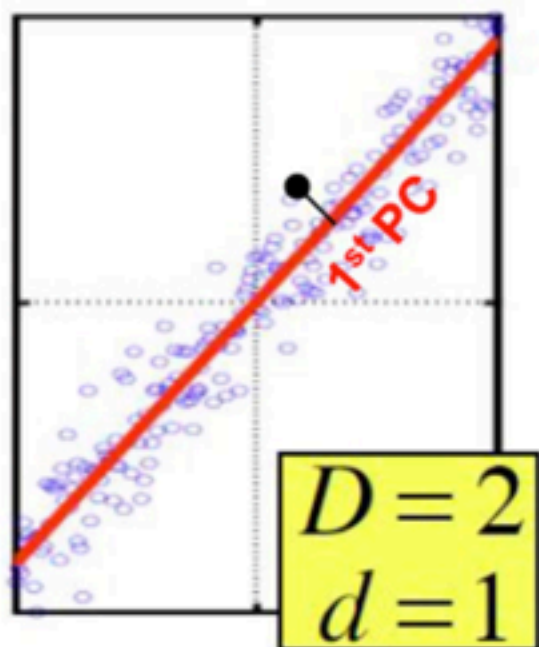
Projection of $x_i$ onto the 1st PC $v$ is $v^T x_i$

# Principal components analysis



Principal Components (PC) are orthogonal directions that capture most of the variance in the data

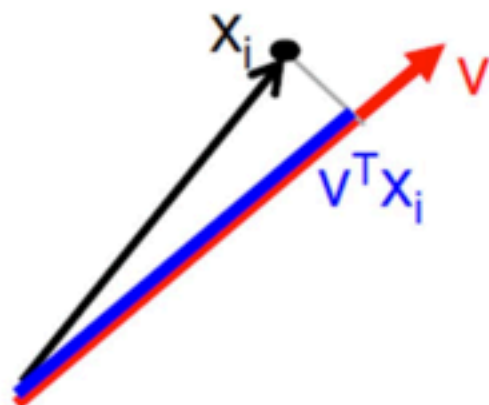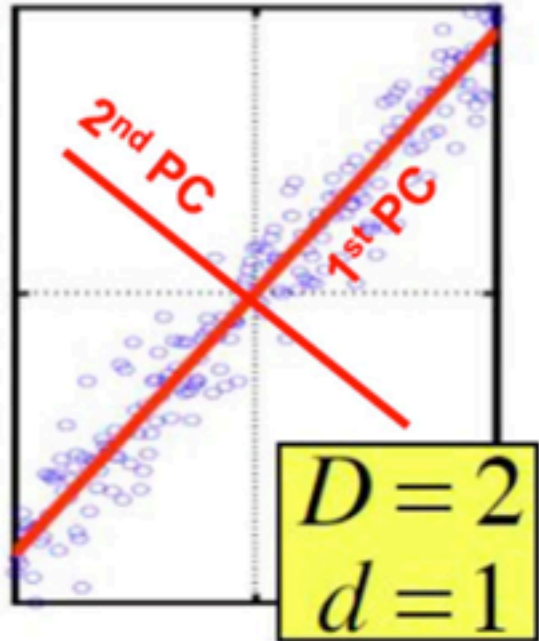1st PC – direction of greatest variability in data

2nd PC – Next orthogonal (uncorrelated) direction of greatest variability

$D = 2$
$d = 1$

Take a data point $x_i$ (D-dimensional vector)
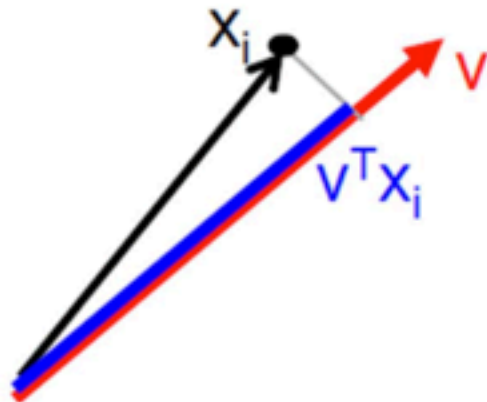
Projection of $x_i$ onto the 1st PC $v$ is $v^T x_i$

# PCA is …

- A backbone of modern data analysis.

- A black box that is widely used but poorly understood.

*OK … let's dispel the magic behind this black box*

# PCA - Overview

- It is a mathematical tool from applied linear algebra.

- It is a simple, non-parametric method of <span style="color:red">extracting</span> relevant information from confusing data sets.

- It provides a roadmap for how to <span style="color:red">reduce</span> a complex data set to a lower dimension.

# What do we need under our BELT?!!!

- Basics of statistical measures, e.g. variance and covariance.

- Basics of linear algebra:
  - Matrices
  - Vector space
  - Basis
  - Eigen vectors and eigen values

# Variance

- A measure of the spread of the data in a data set with mean $\overline{X}$

$$\sigma^2 = \frac{\sum_{i=1}^{n} \left( X_i - \overline{X} \right)^2}{\left( n - 1 \right)}$$

- Variance is claimed to be the original statistical measure of spread of data.

# Covariance

- Variance – measure of the deviation from the mean for points in one dimension, e.g., heights

- Covariance – a measure of how much each of the dimensions varies 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 and grade obtained.

- The covariance between one dimension and itself is the variance

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

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

# Covariance

- What is the interpretation of covariance calculations?

- Say you have a 2-dimensional data set

  - *X*: number of hours studied for a subject

  - *Y*: marks obtained in that subject

- And assume the covariance value (between X and Y) is: 104.53

- *What does this value mean*?

# Covariance

- Exact value is not as important as its sign.

- A <u>positive value</u> of covariance indicates that both dimensions increase or decrease together, e.g., as the number of hours studied increases, the grades in that subject also increase.

- A <u>negative value</u> indicates while one increases the other decreases, or vice-versa, e.g., active social life vs. performance in ECE Dept.

- If <u>covariance is zero</u>: the two dimensions are independent of each other, e.g., heights of students vs. grades obtained in a subject.

# Covariance

- Why bother with calculating (expensive) 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.

# Covariance Matrix

- Representing covariance among dimensions as a matrix, e.g., for 3 dimensions:

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

- Properties:

  - Diagonal: variances of the variables

  - cov($X$,$Y$)=cov($Y$,$X$), hence matrix is symmetrical about the diagonal (upper triangular)

  - $m$-dimensional data will result in $m$x$m$ covariance matrix

# Covariance Matrix

- Covariance matrix C is square and symmetric.
- Our aim is to convert C to a diagonal matrix.
- If we convert given matrix X to be a zero mean one, then we can define C as

$$C = \frac{1}{(n-1)} X X^T$$

where $X^T$ is the transpose of X

# Covariance Matrix

Theorem from Linear Algebra:

In linear algebra, a square matrix A is called **diagonalizable** if it is similar to a diagonal matrix, i.e., if there exists an invertible matrix P such that $P^{-1}AP$ is a diagonal matrix.
That is, we can discover a diagonal matrix D such that:

$$D = P^{-1}AP$$

# Covariance Matrix

$$P^{-1}AP = \begin{pmatrix} \lambda_1 & 0 & \ldots & 0 \\ 0 & \lambda_2 & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \ldots & \lambda_n \end{pmatrix},$$

then:

$$AP = P \begin{pmatrix} \lambda_1 & 0 & \ldots & 0 \\ 0 & \lambda_2 & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \ldots & \lambda_n \end{pmatrix}.$$

# Covariance Matrix

Writing $P$ as a block matrix of its column vectors $\vec{\alpha}_i$

$$P = (\ \vec{\alpha}_1 \quad \vec{\alpha}_2 \quad \cdots \quad \vec{\alpha}_n\ ),$$

the above equation can be rewritten as

$$A\vec{\alpha}_i = \lambda_i \vec{\alpha}_i \qquad (i = 1, 2, \cdots, n).$$

What does this remind us of?
Eigenvectors $\vec{\alpha_i}$ and eigenvalues $\lambda_i$

So the column vectors of P are eigenvectors of A, and the corresponding diagonal entry is the corresponding eigenvalue.

# Now … Let's go back

# PCA

# Example of a problem

- We collected *m* parameters about 100 students:
  - Height
  - Weight
  - Hair color
  - Average grade
  - …
- We want to find the most important parameters that best describe a student.

# Example of a problem

- Each student has a vector of data which describes him of length $m$:

  - (180,70,'purple',84,…)

- We have $n = 100$ such vectors. Let's put them in one matrix, where each column is one student vector.

- So we have a $m$x$n$ matrix. This will be the input of our problem.

# Which parameters can we ignore?

- <span style="color:red">Constant</span> parameter  (number of heads)

  – 1,1,…,1.

- Constant parameter with some <span style="color:red">noise</span> - (thickness of hair)

  – 0.003, 0.005,0.002,….,0.0008 ➔ low variance

- Parameter that is <span style="color:red">linearly dependent</span> on other parameters (head size and height)

  – Z= aX + bY

# Which parameters do we want to keep?

- Parameter that doesn't depend on others (e.g. eye color), i.e. uncorrelated ➜ low covariance.

- Parameter that changes a lot (grades)
  - The opposite of noise
  - High variance

# Covariance Matrix

- Assuming zero mean data (subtract the mean), consider the indexed vectors $\{\mathbf{x_1}, \mathbf{x_2}, ..., \mathbf{x_m}\}$ which are the *rows* of an *mxn* matrix $\mathbf{X}$.

- Each row corresponds to all measurements of a particular measurement type or attribute($\mathbf{x_i}$).

- Each column of $\mathbf{X}$ corresponds to a set of measurements from particular instance or example.

- We now arrive at a definition for the covariance matrix $\mathbf{S_X}$.

$$\mathbf{S_X} \equiv \frac{1}{n-1}\mathbf{XX}^T \quad \text{where} \quad X = \begin{bmatrix} \mathbf{x_1} \\ \vdots \\ \mathbf{x_m} \end{bmatrix}$$

# Covariance Matrix

$$\mathbf{S_X} \equiv \frac{1}{n-1}\mathbf{XX}^T \quad \text{where} \quad X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$$

- The $ij^{th}$ element of the variance is the dot product between the vector of the $i^{th}$ measurement type with the vector of the $j^{th}$ measurement type.

  - $\mathbf{S_X}$ is a square symmetric $m{\times}m$ matrix.

  - The diagonal terms of $\mathbf{S_X}$ are the variance of particular measurement types.

  - The off-diagonal terms of $\mathbf{S_X}$ are the covariance between measurement types.

# Diagonalize the Covariance Matrix

Our goals are to find the covariance matrix that:

1.    Minimizes redundancy, measured by covariance. (off-diagonal), i.e. we would like each variable to co-vary as little as possible with other variables.

2.    Maximizes the signal, measured by variance. (the diagonal)

Since covariance is non-negative, the optimized covariance matrix will be a diagonal matrix.

# Diagonalize the Covariance Matrix
## PCA Assumptions

- PCA assumes that all basis vectors $\{\mathbf{p_1}, \ldots, \mathbf{p_m}\}$ are orthonormal (i.e. $\mathbf{p_i} \cdot \mathbf{p_j} = \delta_{ij}$).

- Hence, in the language of linear algebra, PCA assumes $\mathbf{P}$ is an orthonormal matrix.

- Secondly, PCA assumes the directions with the largest variances are the most "important" or in other words, most principal.

- *Why are these assumptions easiest?*

# Solving PCA: Eigen Vectors of Covariance Matrix

- We will derive our first algebraic solution to PCA using linear algebra. This solution is based on an important property of *eigenvector decomposition*.

- Once again, the data set is **X**, an $m \times n$ matrix, where $m$ is the number of measurement types and $n$ is the number of data trials.

- The goal is summarized as follows:

  - Find some orthonormal matrix **P** where **Y = PX** such that is $S_Y \equiv \frac{1}{n-1} YY^T$ diagonalized. The rows of **P** are the *principal components* of **X**.
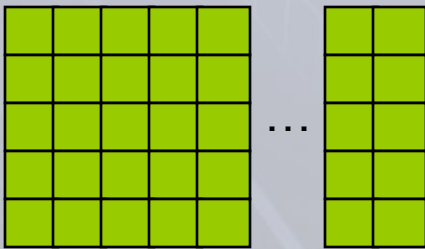
# PCA Process – STEP 1

- Subtract the mean from each of the dimensions

- This produces a data set whose mean is zero.

- Subtracting the mean makes variance and covariance calculation easier by simplifying their equations.

- The variance and co-variance values are not affected by the mean value.

- Suppose we have two measurement types $X_1$ and $X_2$, hence $m = 2$, and ten samples each, hence $n = 10$.

# Principal Component Analysis (PCA)
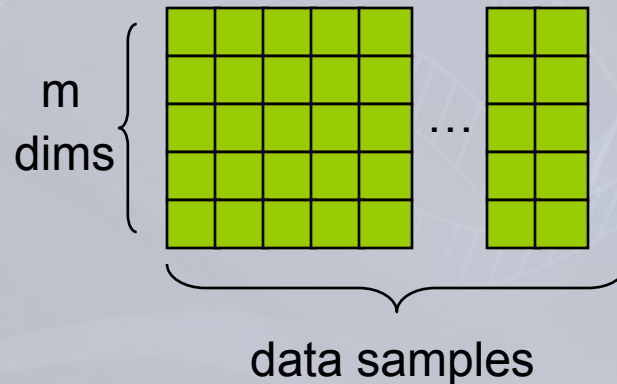
Steps in PCA: #1 Calculate Adjusted Data Set

**Adjusted Data Set: A**　　　　　**Data Set: D**　　**Mean values: M**



... = m dims { ... } -

data samples

$M_i$ is calculated by taking the mean of the values in dimension i

# PCA Process – STEP 1

| $X_1$ | $X_2$ | | | | $X'_1$ | $X'_2$ |
|-------|-------|---|---|---|--------|--------|
| 2.5 | 2.4 | | | | 0.69 | 0.49 |
| 0.5 | 0.7 | | | | $-1.31$ | $-1.21$ |
| 2.2 | 2.9 | | | | 0.39 | 0.99 |
| 1.9 | 2.2 | | | | 0.09 | 0.29 |
| 3.1 | 3.0 | $\Rightarrow$ | $\overline{X_1} = 1.81$ | $\Rightarrow$ | 1.29 | 1.09 |
| 2.3 | 2.7 | | $\overline{X_2} = 1.91$ | | 0.49 | 0.79 |
| 2.0 | 1.6 | | | | 0.19 | $-0.31$ |
| 1.0 | 1.1 | | | | $-0.81$ | $-0.81$ |
| 1.5 | 1.6 | | | | $-0.31$ | $-0.31$ |
| 1.2 | 0.9 | | | | $-0.71$ | $-1.01$ |

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

# PCA Process – STEP 2
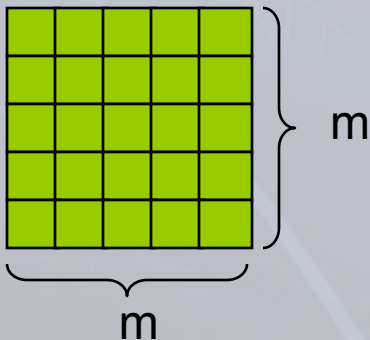
- Calculate the covariance matrix

$$\text{cov} = \begin{bmatrix} 0.616555556 & 0.615444444 \\ 0.615444444 & 0.716555556 \end{bmatrix}$$

- Since the non-diagonal elements in this covariance matrix are positive, we should expect that both the $X_1$ and $X_2$ variables increase together.

- Since it is symmetric, we expect the eigenvectors to be orthogonal.

# Principal Component Analysis (PCA)

Steps in PCA: #2 Calculate Co-variance matrix, C, from Adjusted Data Set, A

Co-variance Matrix: C

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

m

m

$C_{ij} = cov(i,j)$

Note: Since the means of the dimensions in the adjusted data set, A, are 0, the covariance matrix can simply be written as:

C = (A A$^T$) / (n-1)

# PCA Process – STEP 3
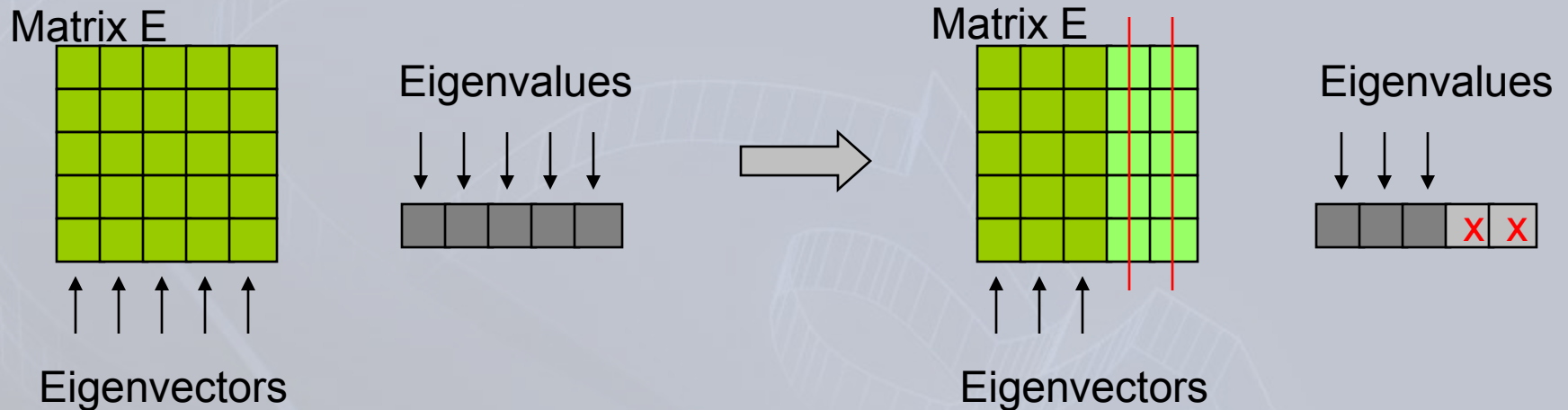
- Calculate the eigenvectors and eigenvalues of the covariance matrix

$$eigenvalues = \begin{bmatrix} 0.490833989 \\ 1.28402771 \end{bmatrix}$$

$$eigenvectors = \begin{bmatrix} -0.735178656 & -0.677873399 \\ 0.677873399 & -0.735178656 \end{bmatrix}$$

# Principal Component Analysis (PCA)

Steps in PCA: #3 Calculate eigenvectors and eigenvalues of C

Matrix E

Eigenvalues

Matrix E

Eigenvalues

Eigenvectors

Eigenvectors

If some eigenvalues are 0 or very small, we can essentially discard those eigenvalues and the corresponding eigenvectors, hence reducing the dimensionality of the new basis.

# PCA Process – STEP 3



Mean adjusted data with eigenvectors overlayed

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).

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 Process – STEP 4

- Reduce dimensionality and form *feature vector*

The eigenvector with the *highest* eigenvalue is the *principal component* of the data set.

In our example, the eigenvector with the largest eigenvalue is the one that points down the middle of the data.

Once eigenvectors are found from the covariance matrix, the next step is to order them by eigenvalue, highest to lowest. This gives the components in order of significance.

# PCA Process – STEP 4

Now, if you'd like, you can decide to *ignore* the components of lesser significance.

You do lose some information, but if the eigenvalues are small, you don't lose much

- *m* dimensions in your data

- calculate *m* eigenvectors and eigenvalues

- choose only the first *r* eigenvectors

- final data set has only *r* dimensions.

# PCA Process – STEP 4

- When the $\lambda_i$'s are sorted in descending order, the proportion of variance explained by the $r$ principal components is:

$$\frac{\displaystyle\sum_{i=1}^{r} \lambda_i}{\displaystyle\sum_{i=1}^{m} \lambda_i} = \frac{\lambda_1 + \lambda_2 + \ldots + \lambda_r}{\lambda_1 + \lambda_2 + \ldots + \lambda_p + \ldots + \lambda_m}$$

- If the dimensions are highly correlated, there will be a small number of eigenvectors with large eigenvalues and $r$ will be much smaller than $m$.

- If the dimensions are not correlated, $r$ will be as large as $m$ and PCA does not help.

# PCA Process – STEP 4

- Feature Vector

$$\text{FeatureVector} = (\lambda_1 \; \lambda_2 \; \lambda_3 \; \ldots \; \lambda_r)$$

(take the eigenvectors to keep from the ordered list of eigenvectors, and form a matrix with these eigenvectors in the columns)

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

$$\begin{bmatrix} -0.677873399 & -0.735178656 \\ -0.735178656 & 0.677873399 \end{bmatrix}$$

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

$$\begin{bmatrix} -0.677873399 \\ -0.735178656 \end{bmatrix}$$

# PCA Process – STEP 5

- Derive the new data

**FinalData = RowFeatureVector x 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*, i.e., the data items are in each column, with each row holding a separate dimension.

# PCA Process – STEP 5

- FinalData is the final data set, <u>with data items in columns, and dimensions along rows.</u>

- <u>What does this give us?</u>

  The original data *solely in terms of the vectors we chose*.
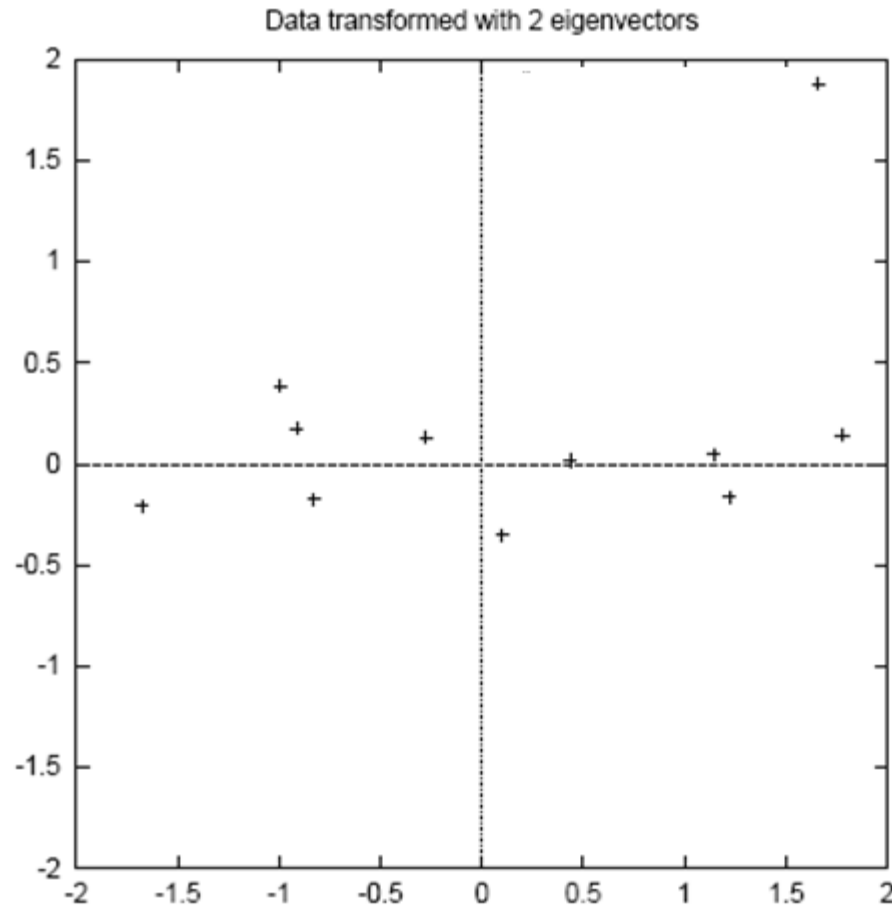
- We have changed our data from being in terms of the axes $X_1$ and $X_2$, to now be in terms of our 2 eigenvectors.

# PCA Process – STEP 5

FinalData (transpose: dimensions along columns)

| $newX_1$ | $newX_2$ |
|---|---|
| $-0.827870186$ | $-0.175115307$ |
| $1.77758033$ | $0.142857227$ |
| $-0.992197494$ | $0.384374989$ |
| $-0.274210416$ | $0.130417207$ |
| $-1.67580142$ | $-0.209498461$ |
| $-0.912949103$ | $0.175282444$ |
| $0.0991094375$ | $-0.349824698$ |
| $1.14457216$ | $0.0464172582$ |
| $0.438046137$ | $0.017646297$ |
| $1.22382956$ | $-0.162675287$ |

# PCA Process – STEP 5



Data transformed with 2 eigenvectors

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

# Applications of PCA

- Exploratory data analysis
- Data preprocessing, dimensionality reduction
- Data compression, data reconstruction
  - (lossy) data compression technique
  - The table describing the data with first k-principal components is smaller than original data table

# References

- J. SHLENS, "TUTORIAL ON PRINCIPAL COMPONENT ANALYSIS," 2005
  http://www.snl.salk.edu/~shlens/pub/notes/pca.pdf

- Introduction to PCA,
  http://dml.cs.byu.edu/~cgc/docs/dm/Slides/