# Teaching
# Principal Components Analysis
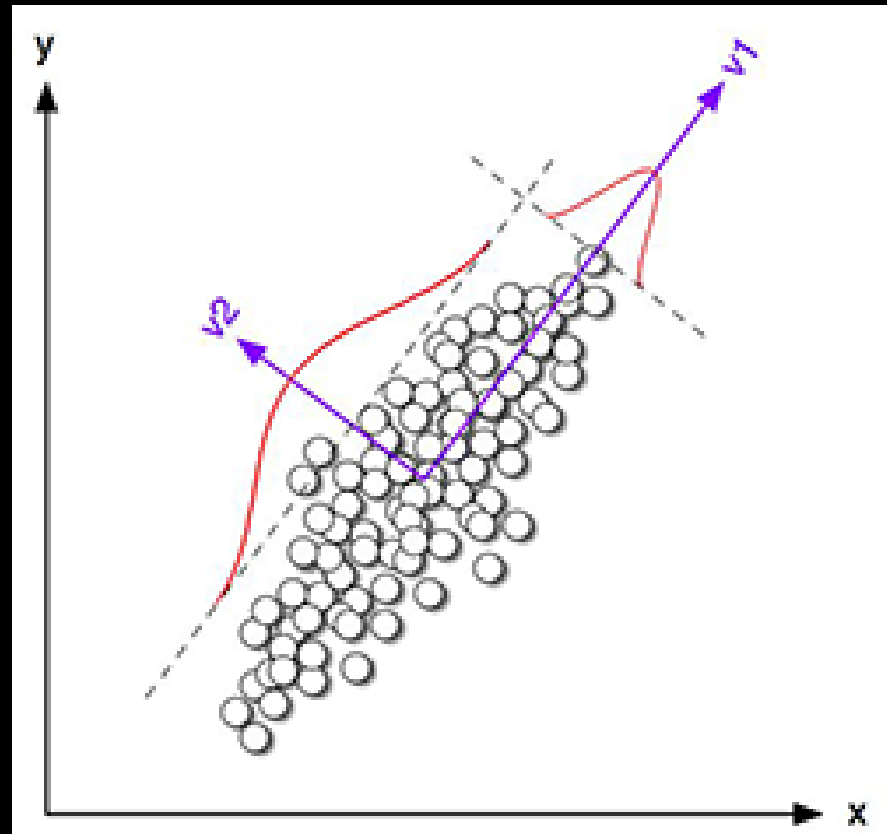# with
# Minitab

*Dr. Jaime Curts*

*The University of Texas Pan American*

ACA 2009 to be held June 25-28, 2009
at École de Technologie Supérieure (ETS),
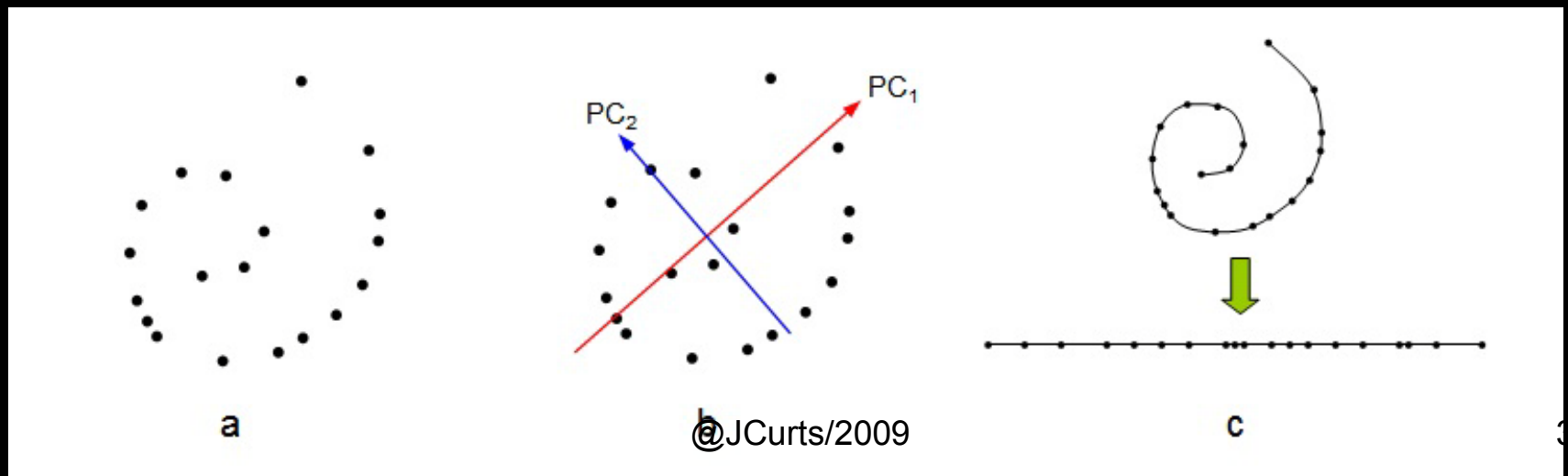Université du Québec, Montréal, Québec, Canada

# Introduction

**The purpose of this presentation is to introduce the logical and arithmetic operators and simple matrix functions of Minitab® –a well-known software package for teaching statistics- as a computer-aid to teach Principal Components Analysis (PCA) to graduate students in the field of Education.**
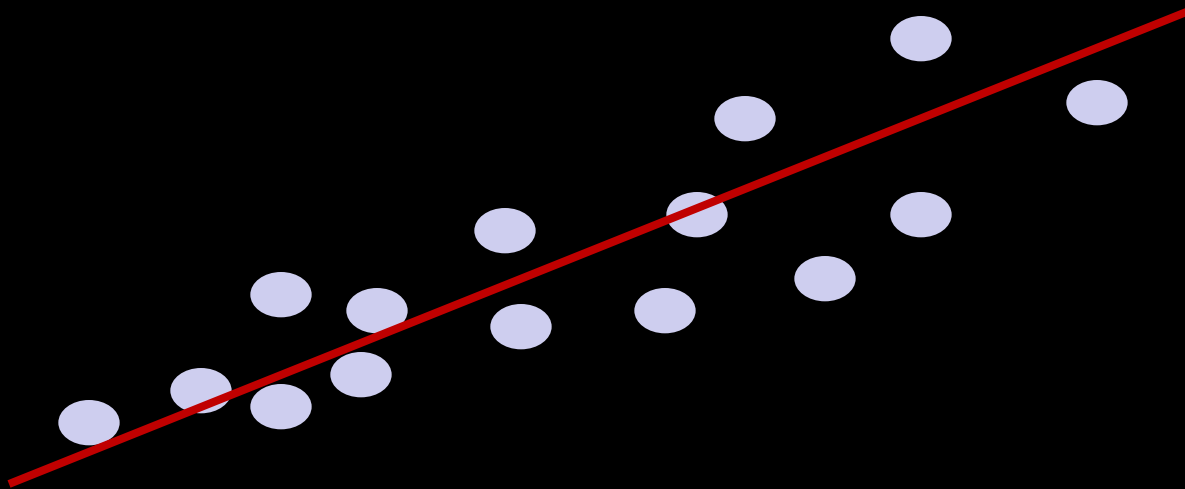
**PCA, originally proposed by Pearson (1901) is a mathematical technique –a vector space transform- that has its roots in linear algebra and in statistics.**

**Its main purpose is to reduce a correlated multidimensional data set to an uncorrelated lower dimensional space with maximum variance.**

**PCA concepts can be a roadblock for non-mathematical oriented students, since statistical definitions (i.e., variance-covariance, correlation) need to be connected to matrix algebra (eigenvectors of a variance-covariance matrix) and to graphical vector representation (including matrix rotation).**

- Given m points in a n dimensional space, for large n, how does one project on to a low dimensional space while preserving broad trends in the data and allowing it to be visualized?

- Choose a line that fits the data so the points are spread

@JCurts/2009

# A sample of *n* observations in the 2-D space



**Goal:** to account for the variation in a sample
in as few variables as possible, to some accuracy

Formally, minimize sum of squares of distances to the line.

Why sum of squares?
Because it allows fast minimization, assuming the line passes through 0

Minimizing sum of squares of distances to the line is the same as maximizing the sum of squares of the projections on that line, thanks to Pythagoras.

@JCurts/2009

Y1 and Y2 are new coordinates.
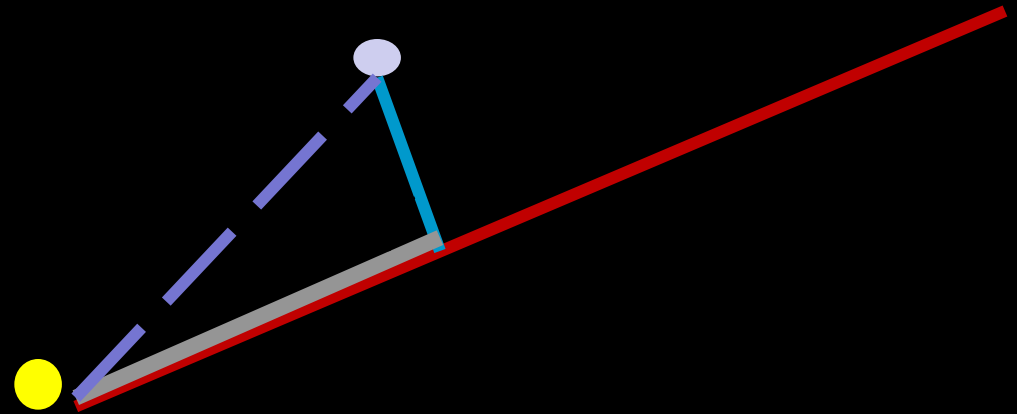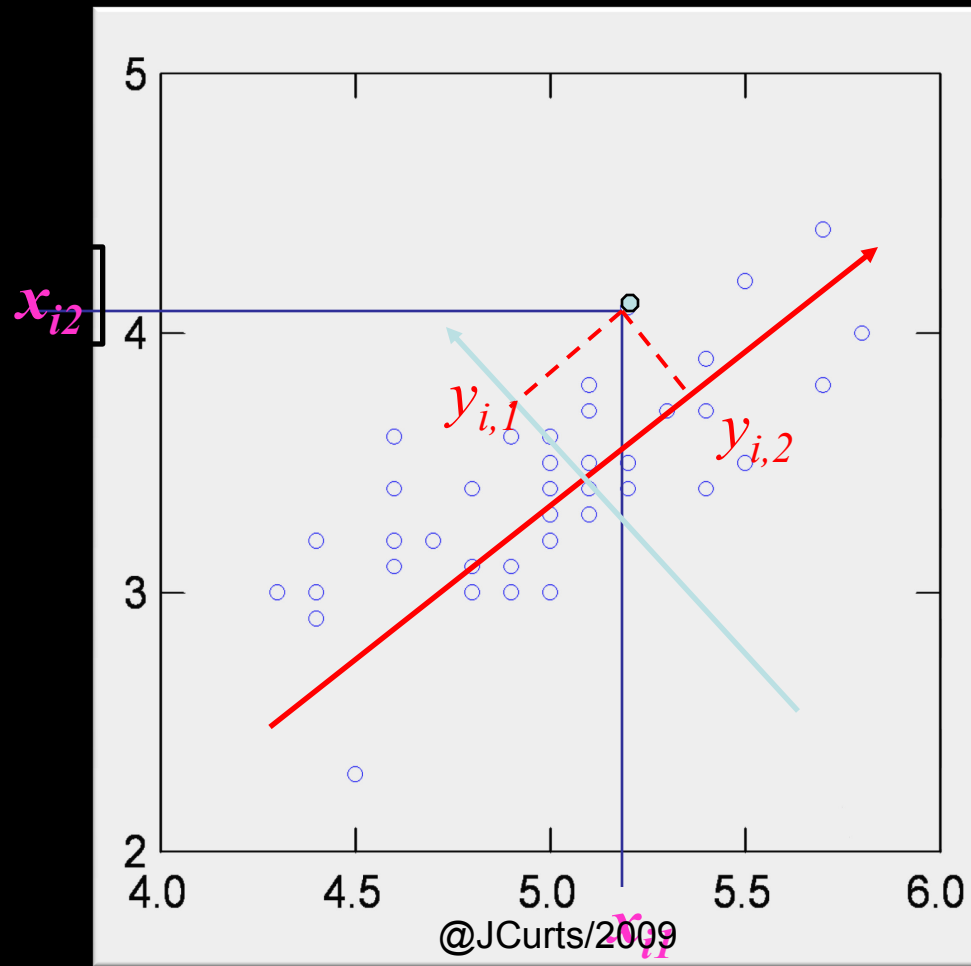  •Y1 represents the direction where the data values have the largest uncertainty.
  •Y2 is perpendicular to Y1.

To find Y1 and Y2, we need to make transformation from X1 and X2. To simplify the discussion, we move the origin to $(\bar{x}_1, \bar{x}_2)$ and redefine the (X1,X2) coordinate as

x1 = X1 - $\bar{x}_1$    , x2 = X2 - $\bar{x}_2$    , so that the origin is (0,0).

The relationship is illustrated in the following graph. We would like to present the data of a given lab, p = (x1,x2) in terms of p = (y1,y2).

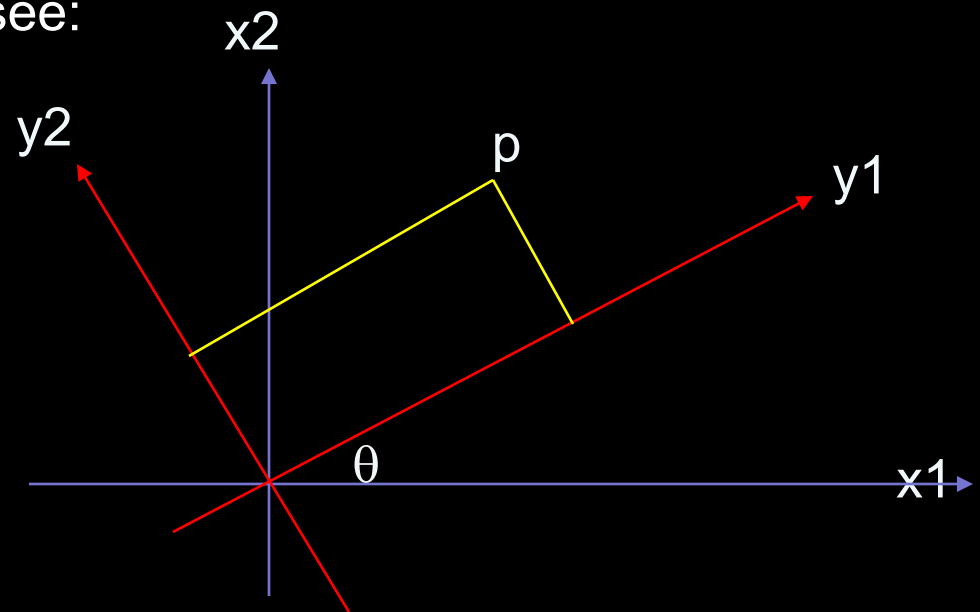From basic geometry relations, we see:
y1 = (cosθ) x1 + (sinθ) x2
y2 = (-sinθ) x1 + (cosθ) x2

The angle θ is determined so that the observations along the Y1 axis has the largest variability.
**But  HOW?**

@JCurts/2009

10

Boxplot of Stack data

**For any given value of theta, then, it is a simple matter to work out the values of Y1 for each of our twenty observations. When θ is 5 degrees, for example, the calculations are:**

| Z1 | Z2 | New Variable Y1 |
|---|---|---|
| 1.90 | 0.47 | **1.94** |
| 0.99 | 0.85 | 1.06 |
| 1.22 | 0.09 | 1.22 |
| 0.54 | -0.68 | 0.47 |
| 0.31 | 0.47 | 0.35 |
| 0.08 | -1.25 | -0.03 |
| -0.15 | -0.30 | -0.17 |
| -0.60 | 0.09 | -0.59 |
| -1.06 | -1.63 | -1.20 |
| -1.29 | -1.25 | -1.39 |
| -1.74 | -1.63 | -1.88 |
| -1.52 | -1.06 | -1.60 |
| 0.76 | 0.47 | 0.80 |
| 1.90 | 2.57 | 2.12 |
| 0.31 | 0.28 | 0.33 |
| -0.38 | -0.10 | -0.38 |
| -0.15 | 0.66 | -0.09 |
| -0.15 | 0.85 | -0.07 |
| -0.38 | 0.66 | -0.32 |
| -0.60 | 0.47 | -0.56 |

Mean (X1) = 7.65; VAR (X1) = 19.23

Mean (Z1) = Mean (Z2) = 0;

Variance (Z1) = Variance (Z2) = 1

**VAR (Y1) =** $(\cos\theta)\, x1 + (\sin\theta)\, x2 =$ **1.12**

Note that each of the original variables has a variance of 1.0, but the variance of the new axis is 1.12, which constitutes more than half of the total variance for the entire dataset (e.g., 1.12/2.00 or 56%).
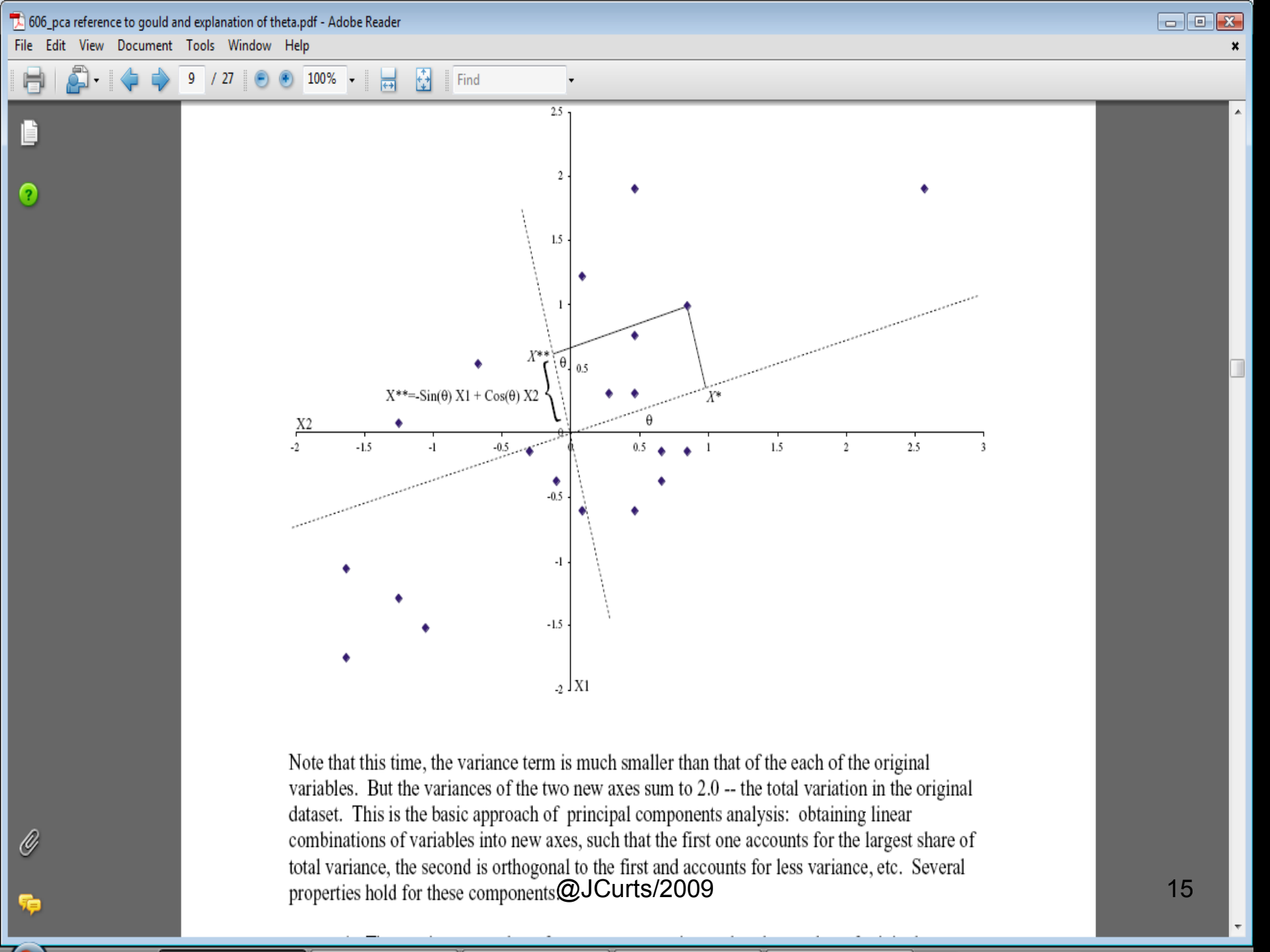
Each value of theta will yield a different set of scores on Y1, and will also result in distinct values for the variance term. If we calculate transformed values and variances for different values of theta, we can compare the variance of the new axis to the total for our dataset.

Note that as we increase the angle, the new variable accounts for an increasing fraction of total variance, until 45 degrees, and then declines; by the time theta is 90 degrees, the new axis is equivalent to X2, and, not surprisingly, its proportion of variance is back to 1.00 or 50.0%.

| Theta | Var (Y1) | Proportion |
|-------|----------|------------|
| 5 | 1.121 | 56.0% |
| 10 | 1.238 | 61.9% |
| 15 | 1.348 | 67.4% |
| 20 | 1.447 | 72.4% |
| 25 | 1.533 | 76.7% |
| 30 | 1.603 | 80.1% |
| 35 | 1.654 | 82.7% |
| 40 | 1.685 | 84.3% |
| **45** | **1.696** | **84.8%** |
| 50 | 1.685 | 84.3% |
| 55 | 1.654 | 82.7% |
| 60 | 1.603 | 80.1% |
| 65 | 1.533 | 76.7% |
| 70 | 1.447 | 72.4% |
| 75 | 1.348 | 67.4% |
| 80 | 1.238 | 61.9% |
| 85 | 1.121 | 56.0% |
| 90 | 1.000 | 50.0% |

$$X^* = \cos(\theta)\, X1 + \sin(\theta)\, X2$$

X2

X*

θ

X1

@JCurts/2009

14

$$X** = -Sin(\theta)\ X1 + Cos(\theta)\ X2$$

Note that this time, the variance term is much smaller than that of the each of the original variables. But the variances of the two new axes sum to 2.0 -- the total variation in the original dataset. This is the basic approach of principal components analysis: obtaining linear combinations of variables into new axes, such that the first one accounts for the largest share of total variance, the second is orthogonal to the first and accounts for less variance, etc. Several properties hold for these components @JCurts/2009

15

The transformation from (x1,x2) to (y1,y2) results several nice properties

1. The variability along y1 is largest.
2. Y1 and y2 are uncorrelated, that is, orthogonal.
3. The confidence region based on (y1,y2) is easy to construct, and provide useful interpretations of the two sample plots.

Questions remain unanswered are

1. How to determine the angle $\theta$ so that the variability of observations along the y1 axis is maximized?
2. How to construct the ellipse for confidence region with different levels of confidences?
3. How to interpret the two-sample plots?

✓ How to determine the Y1 and Y2 axis so that the variability of observations along the Y1 axis is maximized and Y2 is orthogonal to Y1?

✓ Rewrite the linear relation between (Y1,Y2) and (x1,x2) in matrix notation:
Y1 = (cosθ) x1 + (sinθ) x2
Y2 = (-sinθ) x1 + (cosθ) x2

$$Y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} (\cos\theta)x_1 + (\sin\theta)x_2 \\ (-\sin\theta)x_1 + (\cos\theta)x_2 \end{bmatrix} = \begin{bmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} A_1'X \\ A_2'X \end{bmatrix} = AX$$

NOTE: X is bivariate , so is Y, and

V(X) = $\begin{bmatrix} V(x_1) & Cov(x_1,x_2) \\ Cov(x_1,x_2) & V(x_2) \end{bmatrix}$ , V(Y) = A'V(X)A = $\begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$

$\lambda_1$ and $\lambda_2$ are called the eigen values. Which are the solutions of $\left| V(X) - \lambda I \right| = 0$
And, V(Y1) = $\lambda_1$ , V(Y2) = $\lambda_2$, Correlation between Y1 and Y2 = 0.

$\lambda_1$ and $\lambda_2$ are called the eigen values. Which are the solutions of
And, $V(Y1) = \lambda_1$ , $V(Y2) = \lambda_2$, Correlation between Y1 and Y2 = 0.

The angle $\theta$ = $$(.5)\arctan\left(\frac{2\rho\sigma_1\sigma_2}{\sigma_1^2 - \sigma_2^2}\right)$$ if $\sigma_1 \neq \sigma_2$

when $\sigma_1 = \sigma_2$ , $\theta$ = 45° The angle $\theta$ = $$= \arctan\left(\frac{\lambda_1 - \sigma_1^2}{\rho\sigma_1\sigma_2}\right)$$

Note the angle depends on the correlation between X1 and X2 , as well as, on the variances of X1 and X2, respectively.
• When $\rho$ is close to zero, the angle is also close to zero. If V(X1) and V(X2) are close, then, the scatter plots are scattered like a circle. That is, there is no clear major principal component.
•When $\rho$ is close to zero and V(X1) is much larger than V(X2), then, the angle will be close to zero, and the data points are likely to be parallel to the X-axis. On the other hand, if V(X1) is much smaller than V(X2), the angle will be close to 90°, and the data points will be more likely parallel to the Y-axis.

**Consider, now, we actually observe the following two sample data:**

$$\begin{bmatrix} x_{11} & x_{21} \\ x_{12} & x_{22} \\ x_{13} & x_{23} \\ \vdots & \vdots \\ x_{1n} & x_{2n} \end{bmatrix}$$

The sample means are given by $\begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \end{bmatrix}$

The sample variance-covariance matrix is given by

$$\hat{V}(X) = \begin{bmatrix} s_1^2 & rs_1 s_2 \\ rs_1 s_2 & s_2^2 \end{bmatrix}$$

r is the Pearson's correlation coefficient, and $S^2$ is the sample variance. S is the sample standard deviation.

V(Y) is the solution of $\left| \hat{V}(X) - \lambda I \right| = 0$

The solutions for $\lambda$ are given by $\dfrac{(s_1^2 + s_2^2) \pm \sqrt{(s_1^2 + s_2^2)^2 - 4\,1 - r^2)s_1^2 s_2^2}}{2}$

NOTE: V(Y1) + V(Y2) = $\lambda_1 + \lambda_2 = s_1^2 + s_2^2$ = V(X1) + V(X2)

Using the sample data, the angle is estimated by

$$\theta = (.5)\arctan\left(\frac{2rs_1s_2}{s_1^2 - s_2^2}\right) \qquad s_1 \neq s_2$$

$$= \arctan\left(\frac{\lambda_1 - s_1^2}{rs_1s_2}\right)$$

# Correlation and Covariance Matrix

File   Edit   Data   Calc   Stat   Graph   Editor   Tools   Window   Help

**Boxplot of Stack data**
print

**Data Display**

  Matrix CORR1

1.00000   0.70560
0.70560   1.00000
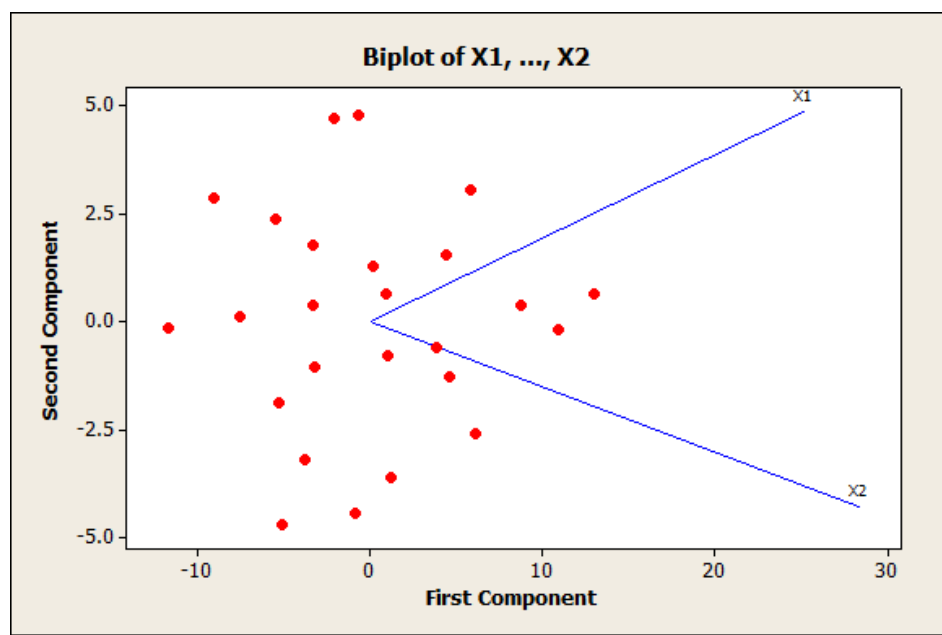
Matrix COVA1

20.2767   15.585
15.5850   24.060

  Matrix M3

0.663139    0.748496
0.748496   -0.663139

**Principal Component Analysis: X1, X2**

Eigenanalysis of the Covariance Matrix

Eigenvalue  37.868   6.469
Proportion   0.854   0.146
Cumulative   0.854   1.000

Variable     PC1      PC2
X1         0.663    0.748
X2         0.748   -0.663



Biplot of X1, ..., X2

@JCurts/2009

Welcome to Minitab, press F1 for help.

Editable

Scatterplot of C14 vs C15

**Pearson correlation of C14 and C15 = -0.000**

Thanks…………………….

jbcurts@utpa.edu