

# Scientific Computing for Biologists

## Data as Vectors: Geometry of Bivariate Relationships I

Instructor: Paul M. Magwene

# Overview of Lecture

- Variable space/Subject space representations
- Vector Geometry
  - Vectors are directed line segments
  - Vector length
- Vector Arithmetic
  - Addition, subtraction
  - Scalar multiplication
  - Linear combinations of vectors
  - Dot product and projection
- Vector representations of multivariate data
  - Mean as projection in subject space
  - Bivariate regression in geometric terms

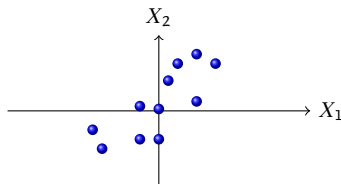
# Variable Space Representation of a Data Set

Consider a data set in which we've measured variables

$X = X_1, X_2, \dots, X_p$ , on a set of subjects (observations)  $s_1, \dots, s_n$ .

	$X_1$	$X_2$
$s_1$	0.9	1.4
$s_2$	1.1	1.7
$\vdots$	$\vdots$	$\vdots$
$s_n$	0.5	1.55

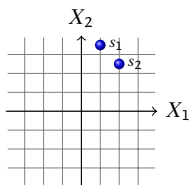
Such data is most often represented by drawing the observations as points in space of dimension  $p$ . This is the *variable space representation* of the data.



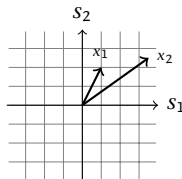
# Subject Space Representation of a Data Set

An alternate representation is to consider the variables in the space of the subjects. This is the *subject space* representation.

	$X_1$	$X_2$
$s_1$	1	3.5
$s_2$	2	2.5



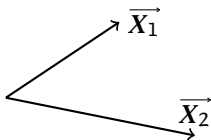
(a) Variable space representation



(b) Subject space representation

# How do we come up with a useful representation of variables in subject space?

- Let the variables be represented by centered vectors
  - lengths of vectors are proportional to standard deviation
  - angle between vectors represents association or similarity



This representation of variables as vectors in the space of the subjects is the view that we'll develop over the next few lectures.

# Variable Space vs Subject Space Representations

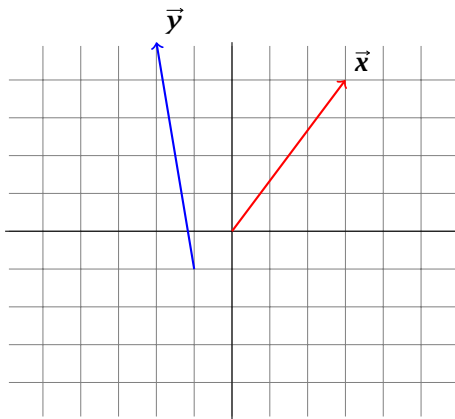
- Variable space – useful when multivariate analyses are focused on the relationship among observations in the data
- Subject space – useful when multivariate analyses are focused on the relationship among the variables of the data

We'll even see methods that simultaneously depict both subject space and variable space simultaneously.

Understanding these two different ways of thinking about complex data will help you to understand what various multivariate statistical methods do; and what their constraints or limitations are.

# Vector Geometry

Vectors are directed line segments.

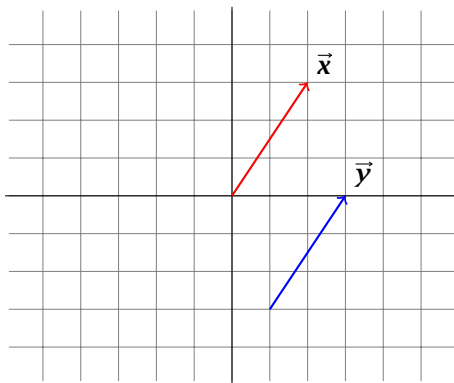


All of the figures and algebraic formulas I show you apply to  $n$ -dimensional vectors.

# Vector Geometry

Vectors have direction and length:

$$\vec{x} = [x_1, x_2]' = [2, 3]'; |\vec{x}| = \sqrt{x_1^2 + x_2^2 + \cdots + x_n^2}$$



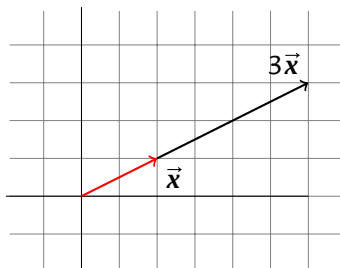
Often starting point is ignored, in which case  $\vec{x} = \vec{y}$ .



# Scalar Multiplication of a Vector

Let  $k$  be a scalar.

$$k\vec{x} = \begin{bmatrix} kx_1 \\ kx_2 \\ \vdots \\ kx_n \end{bmatrix}$$

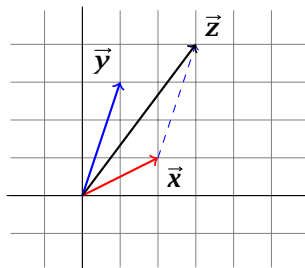


$$\vec{x} = [2, 1]'; \quad 3\vec{x} = [6, 3]'.$$

# Vector Addition

Let  $\vec{x} = [2, 1]'$ ;  $\vec{y} = [1, 3]'$

$$\vec{z} = \vec{x} + \vec{y} = \begin{bmatrix} x_1 + y_1 \\ x_2 + y_2 \\ \vdots \\ x_n + y_n \end{bmatrix}$$

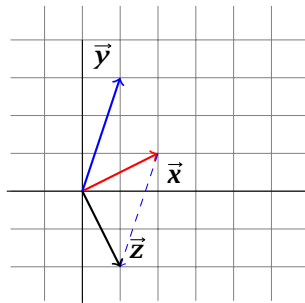


Addition follows the 'head-to-tail' rule.

# Vector Subtraction

Let  $\vec{x} = [2, 1]'$ ;  $\vec{y} = [1, 3]'$

$$\vec{z} = \vec{x} - \vec{y} = \begin{bmatrix} x_1 - y_1 \\ x_2 - y_2 \end{bmatrix} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$$

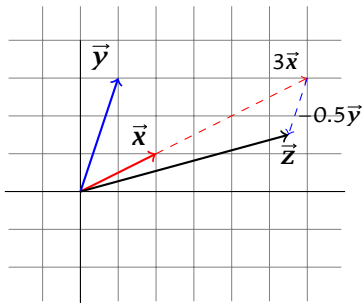


Follow the addition rule for  $-1\vec{y}$ .

# Linear Combinations of Vectors

A linear combination of vectors is of the form  $\vec{z} = b_1\vec{x} + b_2\vec{y}$

$$\vec{z} = 3\vec{x} - 0.5\vec{y} = 3 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - 0.5 \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

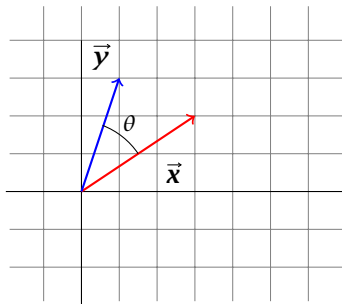


# Dot Product

The dot (inner) product of two vectors,  $\vec{x} \cdot \vec{y}$  is a scalar.

$$\begin{aligned}\vec{x} \cdot \vec{y} &= x_1 y_1 + x_2 y_2 + \cdots + x_n y_n \\ &= |\vec{x}| |\vec{y}| \cos \theta\end{aligned}$$

where  $\theta$  is the angle (in radians) between  $\vec{x}$  and  $\vec{y}$



$$\vec{x} = [3, 2]', \vec{y} = [1, 3]'; \vec{x} \cdot \vec{y} = \sqrt{13}\sqrt{10}\cos \theta = 9$$

# Useful Geometric Quantities as Dot Product

Length:

$$|\vec{x}|^2 = \vec{x} \cdot \vec{x} = x_1^2 + x_2^2 + \cdots + x_n^2$$

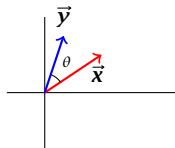
$$|\vec{y}|^2 = \vec{y} \cdot \vec{y}$$

Distance:

$$|\vec{x} - \vec{y}|^2 = \vec{x} \cdot \vec{x} + \vec{y} \cdot \vec{y} - 2\vec{x} \cdot \vec{y}$$

Angle:

$$\cos \theta = \vec{x} \cdot \vec{y} / (|\vec{x}| |\vec{y}|)$$



# Dot Product Properties

Some additional properties of the dot product that are useful to know:

$$\vec{x} \cdot \vec{y} = \vec{y} \cdot \vec{x} \text{ (commutative)}$$

$$\vec{x} \cdot (\vec{y} + \vec{z}) = \vec{x} \cdot \vec{y} + \vec{x} \cdot \vec{z} \text{ (distributive)}$$

$$(k\vec{x}) \cdot \vec{y} = \vec{x} \cdot (k\vec{y}) = k(\vec{x} \cdot \vec{y}) \text{ where } k \text{ is a scalar}$$

$$\vec{x} \cdot \vec{y} = 0 \text{ iff } \vec{x} \text{ and } \vec{y} \text{ are orthogonal}$$

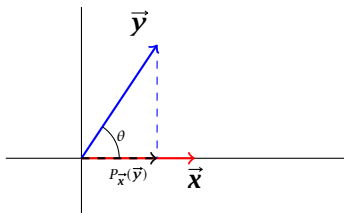
# Vector Projection

The *projection* of  $\vec{y}$  onto  $\vec{x}$ ,  $P_{\vec{x}}(\vec{y})$ , is the vector obtained by placing  $\vec{y}$  and  $\vec{x}$  tail to tail and dropping a line, perpendicular to  $\vec{x}$ , from the head of  $\vec{y}$  onto the line defined by  $\vec{x}$ .

$$P_{\vec{x}}(\vec{y}) = \left( \frac{\vec{x} \cdot \vec{y}}{\vec{x} \cdot \vec{x}} \right) \vec{x} = \left( \frac{\vec{x} \cdot \vec{y}}{|\vec{x}|^2} \right) \vec{x} = \left( \frac{\vec{x} \cdot \vec{y}}{|\vec{x}|} \right) \frac{\vec{x}}{|\vec{x}|}$$

The *component* of  $\vec{y}$  in  $\vec{x}$ ,  $C_{\vec{x}}(\vec{y})$ , is the length of  $P_{\vec{x}}(\vec{y})$ .

$$C_{\vec{x}}(\vec{y}) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}|} = |\vec{y}| \cos \theta$$



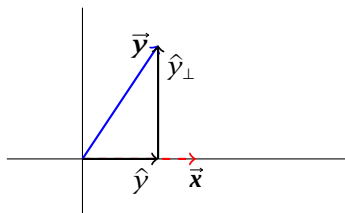


## Vector Projection II

$\vec{y}$  can be decomposed into two parts:

1. a vector parallel to  $\vec{x}$ ,  $\hat{y} = P_{\vec{x}}(\vec{y})$ ,
2. a vector perpendicular to  $\vec{x}$ ,  $\hat{y}_{\perp}$ .

$$\vec{y} = \hat{y} + \hat{y}_{\perp}$$



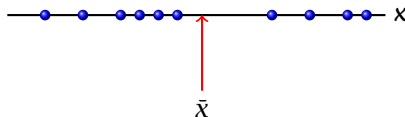
- $\hat{y}_{\perp}$  is *orthogonal* to  $\hat{y}$  and  $\vec{x}$ .
- $\hat{y}$  is the closest vector to  $\vec{y}$  in the subspace defined by  $\vec{x}$

# Vector Geometry of Simple Statistics

# Geometry of the Mean in Variables Space

The mean is a single number summary of a set (vector) of values,  $\vec{x}$ . Mathematically, the mean is the value that minimizes the following quantity:

$$\sum_{i=1}^n (x_i - \bar{x})^2$$



# Algebraic and Vector Formulas for the Mean

Let  $\vec{x} = [x_1, x_2, \dots, x_n]'$

Algebraic formula for the mean of  $\vec{x}$ :

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Vector formula for the mean:

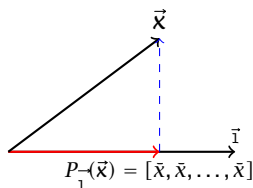
$$\begin{aligned}\bar{x} &= \frac{\vec{1} \cdot \vec{x}}{\vec{1} \cdot \vec{1}} \\ &= \frac{\vec{1} \cdot \vec{x}}{|\vec{1}|^2}\end{aligned}$$

Where  $\vec{1} = [1, 1, \dots, 1]'$  is the 1-vector of dimension  $n$ . (Note that the 1-vector is not the same as a unit vector!)

# Geometry of the Mean in Subject Space

- The mean can be interpreted in terms of the projection of  $\vec{x}$  onto the 1-vector:

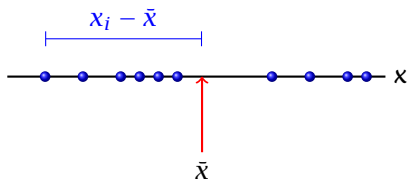
$$\begin{aligned}P_1(\vec{x}) &= \left( \frac{\vec{1} \cdot \vec{x}}{|\vec{1}|^2} \right) \vec{1} \\&= \bar{x} \vec{1} \\&= [\bar{x}, \bar{x}, \dots, \bar{x}]'\end{aligned}$$



# Variable Space Geometry of Sample Variance

Sample variance is proportional to the sum of squared deviates about the mean:

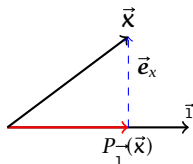
$$S_x^2 = \frac{1}{(n-1)} \sum (x_i - \bar{x})^2$$



# Vector Geometry of Sample Variance

- Let  $\vec{e}_x = \vec{x} - \bar{x}\vec{1}$
- The sample variance can be expressed in terms of dot products of  $\vec{e}_x$  with itself:

$$s_x^2 = \frac{\vec{e}_x \cdot \vec{e}_x}{n-1} = \frac{|\vec{e}_x|^2}{n-1}$$



# Mean centering

In the previous slide, we considered the vector:

$$\vec{e}_x = \vec{x} - \bar{x}\vec{1}$$

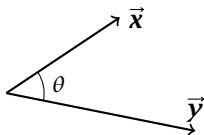
$\vec{e}_x$  is the “mean centered” version of  $\vec{x}$ , i.e. it’s the vector we get when we subtract the mean of  $\vec{x}$ ,  $\bar{x}$ , from every element of  $\vec{x}$ .

For convenience, I will usually state the variables of interest are mean centered and use the notation  $\vec{x}$  instead of  $\vec{e}_x$  so as to avoid a proliferation of subscripts.



# Covariance and Correlation in Vector Geometric Terms

Let  $X$  and  $Y$  be mean centered variables, and let  $\vec{x}$  and  $\vec{y}$  be their corresponding mean centered vector representations.



Vector formulas for covariance and correlation:

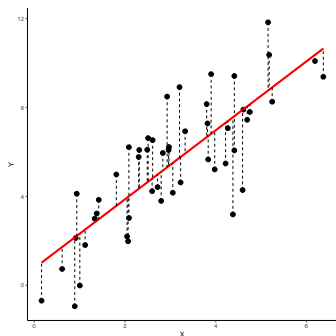
$$\text{Covariance: } \text{cov}(X, Y) = \frac{\vec{x} \cdot \vec{y}}{n - 1}$$

$$\text{Correlation: } \text{corr}(X, Y) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| |\vec{y}|} = \cos \theta$$

## Geometric interpretation of correlation

The correlation between two variables  $X$  and  $Y$  is equivalent to the cosine of the angle between their mean-centered vector representations!

# Bivariate Regression: Variable Space Representation



The standard bivariate regression equation relating one observed variable  $X$  (the predictor) to another observed variable of interest,  $Y$  (the outcome) is usually written as:

$$\hat{Y} = a + bX.$$

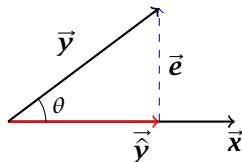
where  $\hat{Y}$  is the predicted value of  $Y$  and  $a$  (intercept) and  $b$  (slope) are chosen to minimize  $\sum (Y_i - \hat{Y}_i)^2$ .

# Geometry of Bivariate Regression

In vector geometric terms, *regression is projection*! Consider the regression formula for mean-centered vectors:

$$\hat{Y} = bX$$

In vector terms we can view this as:



$\hat{y}$  is the closest vector to  $\vec{y}$  in the subspace defined by  $\vec{x}$ , i.e. it is the scalar multiple of  $\vec{x}$  that minimizes  $|\vec{e}|$

$$\hat{\vec{y}} = P_{\vec{x}}(\vec{y}) = \left( \frac{\vec{x} \cdot \vec{y}}{\vec{x} \cdot \vec{x}} \right) \vec{x}$$

## Bivariate regression in vector terms

For mean centered vectors the regression equation is:

$$\vec{\hat{y}} = b\vec{x}$$

From the previous slide we see that we can solve  $b$  as:

$$b = \frac{\vec{x} \cdot \vec{y}}{\vec{x} \cdot \vec{x}}$$

## Bivariate regression: Alternate formulas for slope I

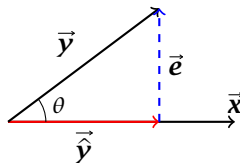
Regression equation for mean-centered vectors:  $\hat{\vec{y}} = b\vec{x}$

There are multiple, equivalent ways to write the solution for  $b$ :

$$\begin{aligned} b &= \frac{\vec{x} \cdot \vec{y}}{(\vec{x} \cdot \vec{x})} \\ &= \frac{\vec{x} \cdot \vec{y}}{|\vec{x}|^2} \\ &= \frac{|x||y| \cos \theta}{|x|^2} \\ &= \cos \theta \frac{|y|}{|x|} \\ &= r_{XY} \frac{|y|}{|x|} \end{aligned}$$

# Geometry of Goodness of Fit

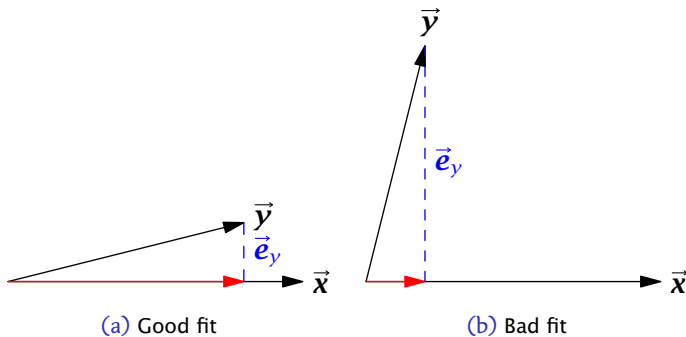
Geometric interpretation of regression goodness-of-fit:



$$|\hat{\vec{y}}|^2 + |\vec{e}|^2 = |\vec{y}|^2$$

The better the goodness-of-fit, the smaller the angle,  $\cos \theta$ , and the shorter residual vector,  $\vec{e}$ .

# Geometry of Goodness of Fit



# Bivariate Regression, Goodness of Fit

How well does our prediction agree with our outcome?

- Measure the angle between  $\vec{\hat{y}}$  and  $\vec{y}$ :

$$R = \cos \theta_{\vec{y}, \vec{\hat{y}}} = \frac{|\vec{\hat{y}}|}{|\vec{y}|}$$

- In the single-predictor case  $R = r_{XY}$ , but this is not generally true when we have multiple predictors.
- Note that  $|\vec{y}|$  can be expressed as follows:

$$\begin{aligned} |\vec{\hat{y}}|^2 + |\vec{e}|^2 &= |\vec{y}|^2 \\ SS_{\text{regression}} + SS_{\text{residual}} &= SS_{\text{total}} \end{aligned}$$

- With simple substitution we can show that:

$$\begin{aligned} SS_{\text{regression}} &= R^2 SS_{\text{total}} \\ SS_{\text{residual}} &= (1 - R^2) SS_{\text{total}} \end{aligned}$$



## Two-group ANOVA as Regression

We can also use a geometric perspective to test whether the mean of a variable differs between two groups of subjects.

- Setup a 'dummy variable' as the predictor  $X_g$ . We assign all subjects in group 1 the value 1 and all subjects in group 2 the value -1 on the dummy variable. We then regress the variable of interest,  $Y$ , on  $X_g$ .

$$y = X_g b + e$$

Group	Raw		Centered	
	$Y_i$	$X_i$	$y_i$	$x_i$
1	2	-1	-3	$-\frac{4}{3}$
	3	-1	-1	$-\frac{4}{3}$
2	5	1	0	$\frac{2}{3}$
	6	1	1	$\frac{2}{3}$
	6	1	1	$\frac{2}{3}$
	7	1	2	$\frac{2}{3}$
Mean	5	$\frac{1}{3}$	0	0

## Two-group ANOVA as Regression, cont

- When the means are different in the two groups,  $X_g$  will be a good predictor of the variable of interest, hence  $\vec{y}$  and  $\vec{x}_g$  will have a small angle between them.
- When the means in the two groups are similar, the dummy variable will not be a good predictor. Hence the angle between  $\vec{y}$  and  $\vec{x}_g$  will be large.

