# Scientific Computing for Biologists

Data as Vectors: Introduction to Vector Geometry

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
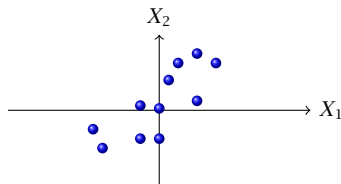    - Bivariate regression in geometric terms

# Variable Space Representation of a Data Set

Consider a data set in which we've measured variables
$\mathbf{X} = X_1, X_2, \ldots, X_p$, on a set of subjects (objects) $a_1, \ldots, a_n$.

|       | $X_1$ | $X_2$ |
|-------|-------|-------|
| $a_1$ | 0.9   | 1.4   |
| $a_2$ | 1.1   | 1.7   |
| $\vdots$ | $\vdots$ | $\vdots$ |
| $a_n$ | 0.5   | 1.55  |

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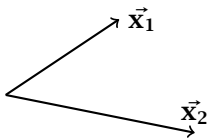representation* of the data.

# 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.

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
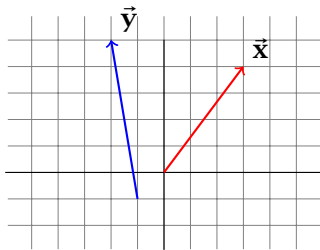


This representation of variables as vectors in the space of the subjects is the view that we'll develop over the next few lectures.

# Vector Geometry

Vectors are directed line segments.

$$\vec{\mathbf{x}} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = [x_1, x_2, \cdots, x_n]'$$
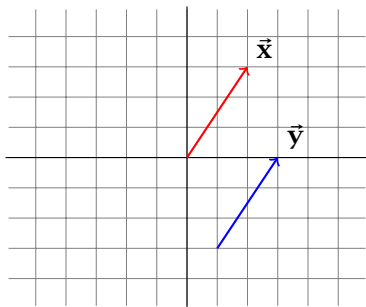


All of the figures and algebraic formulas I show you apply to $n$-dimensional vectors.

# Vector Geometry

Vectors have direction and length:

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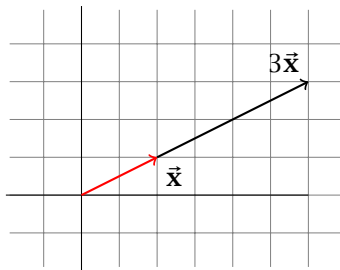


Often starting point is ignored, in which case $\vec{\mathbf{x}} = \vec{\mathbf{y}}$.

# Scalar Multiplication of a Vector

Let $k$ be a scalar.

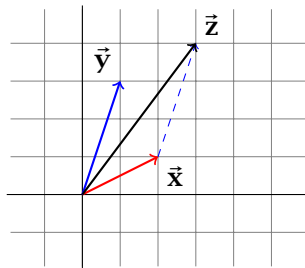$$k\vec{\mathbf{x}} = \left[ \begin{array}{c} kx_1 \\ kx_2 \\ \vdots \\ kx_n \end{array} \right]$$



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

# Vector Addition

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

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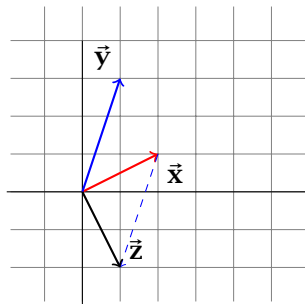


Addition follows the 'head-to-tail' rule.

# Vector Subtraction

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

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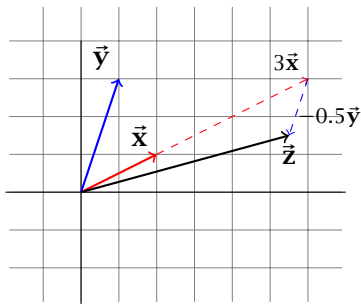


Follow the addition rule for $-1\vec{y}$.

# Linear Combinations of Vectors

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

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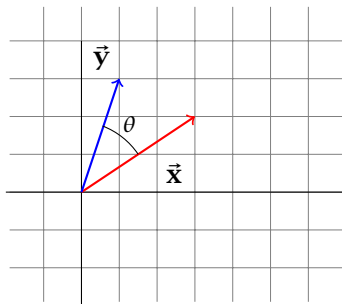
# Dot Product

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:

$$
\begin{aligned}
|\vec{\mathbf{x}}|^2 &= \vec{\mathbf{x}} \cdot \vec{\mathbf{x}} = x_1^2 + x_2^2 + \cdots + x_n^2 \\
|\vec{\mathbf{y}}|^2 &= \vec{\mathbf{y}} \cdot \vec{\mathbf{y}}
\end{aligned}
$$

Distance:

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

Angle:

$$
\cos \theta = \frac{\vec{\mathbf{x}} \cdot \vec{\mathbf{y}}}{|x||y|}
$$

# Dot Product Properties

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

$$
\begin{aligned}
\vec{x} \cdot \vec{y} &= \vec{y} \cdot \vec{x} \text{ (commutative)} \\
\vec{x} \cdot (\vec{y} + \vec{z}) &= \vec{x} \cdot \vec{y} + \vec{x} \cdot \vec{z} \text{ (distributive)} \\
(k\vec{x}) \cdot \vec{y} &= \vec{x} \cdot (k\vec{y}) = k(\vec{x} \cdot \vec{y}) \text{ where } k \text{ is a scalar} \\
\vec{x} \cdot \vec{y} &= 0 \text{ iff } \vec{x} \text{ and } \vec{y} \text{ are orthogonal}
\end{aligned}
$$

# Useful vectors

- Unit vector in the direction of $\vec{x}$ – a vector of length one parallel to $\vec{\mathbf{x}}$. Can be calculated as:

$$\text{Unit vector in the direction of } \vec{\mathbf{x}} = \frac{\vec{\mathbf{x}}}{|\vec{\mathbf{x}}|}$$

- One-vector – a $n$-dimensional vector of where every element is the number 1.

$$\vec{\mathbf{1}}_n = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

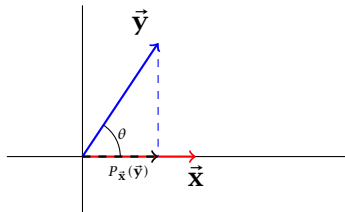Note that 1-vectors are not, in general, unit vectors!

# 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 $\vec{x}$.

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

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

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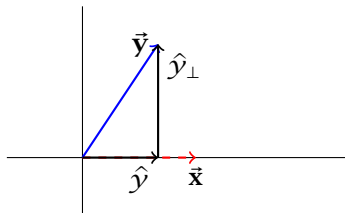
# Vector Projection II

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

1. a vector parallel to $\vec{\mathbf{x}}$, $\hat{y} = P_{\vec{\mathbf{x}}}(\vec{\mathbf{y}})$,

2. a vector perpendicular to $\vec{\mathbf{x}}$, $\hat{y}_\perp$.

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



- $\hat{y}$ is the closest vector to $\vec{\mathbf{y}}$ in the subspace defined by $\vec{\mathbf{x}}$, i.e. $|\hat{y}_\perp|$ is as small as possible
- $\hat{y}_\perp$ is *orthogonal* to $\hat{y}$ and $\vec{\mathbf{x}}$.
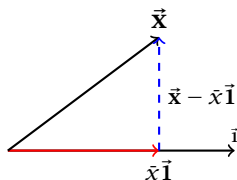
# Vector Geometry of Simple Statistics

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

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

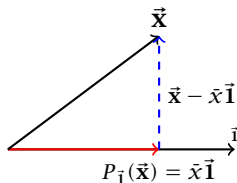# Sketch of Proof: Deriving the mean in vector geometric terms, part I

- The mean, $\bar{x}$, minimizes the quantity $\sum_{i=1}^{n}(x_i - \bar{x})^2$.
- The above can be written as $|\vec{\mathbf{x}} - \vec{\mathbf{1}}\bar{x}|^2$ where $\vec{\mathbf{x}} = [x_1, x_2, \ldots, x_n]'$ and $\vec{\mathbf{1}} = [1, 1, \ldots, 1]'$
- We are look for the scalar multiple, $\bar{x}$, of the one vector that minimizes $|\vec{\mathbf{x}} - \vec{\mathbf{1}}\bar{x}|^2$
- What does the geometry of $\vec{\mathbf{x}}$, $\vec{\mathbf{1}}$, and $\vec{\mathbf{x}} - \vec{\mathbf{1}}\bar{x}$ look like?



This picture looks familiar! Where did we see it before?

# Sketch of Proof: Deriving the mean in vector geometric terms, part II

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



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

# Vector and Algebraic Formulas for the Mean

Vector formula for the mean:

$$\bar{x} \quad = \quad \frac{\vec{\mathbf{1}} \cdot \vec{\mathbf{x}}}{\vec{\mathbf{1}} \cdot \vec{\mathbf{1}}}$$
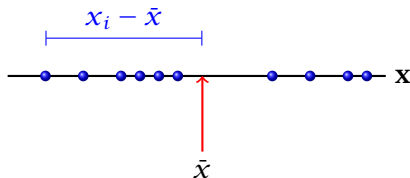
Algebraic formula for the mean of $\vec{\mathbf{x}}$:

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n} = \frac{1}{n} \sum_{i=1}^{n} x_i$$
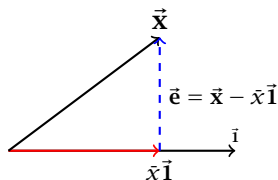
# Variable Space Geometry of Sample Variance

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 $\vec{\mathbf{e}}_{\mathbf{x}} = \vec{\mathbf{x}} - \bar{x}\vec{\mathbf{1}}$



The sample variance can be expressed in terms of dots products of $\vec{\mathbf{e}}_{\mathbf{x}}$ with itself:

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

# Mean centering

In the previous slide, we considered the vector:

$$\vec{e_x} = \vec{x} - \bar{x}\vec{1}$$

We can think of $\vec{e}_x$ as a "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 $\vec{x}$.
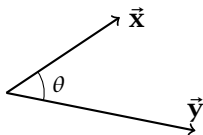
Important relationships for mean-centered vectors:

- The variance of $X$ is proportional to $|\vec{e}_x|^2$
- The standard deviation of $X$ is proportional to $|\vec{e}_x|$

For convenience, I will sometimes 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 variables of interest, and let $\vec{\mathbf{x}}$ and $\vec{\mathbf{y}}$ be their corresponding mean centered vector representations.



Vector formulas for covariance and correlation:

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

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

## 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!