

CS57300  
PURDUE UNIVERSITY  
AUGUST 30, 2021

---

# DATA MINING

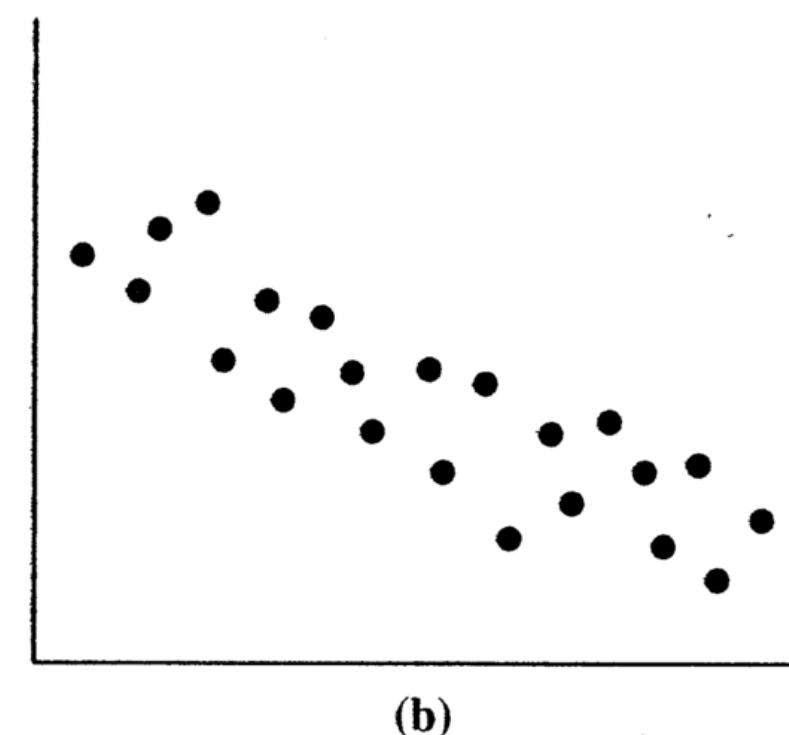
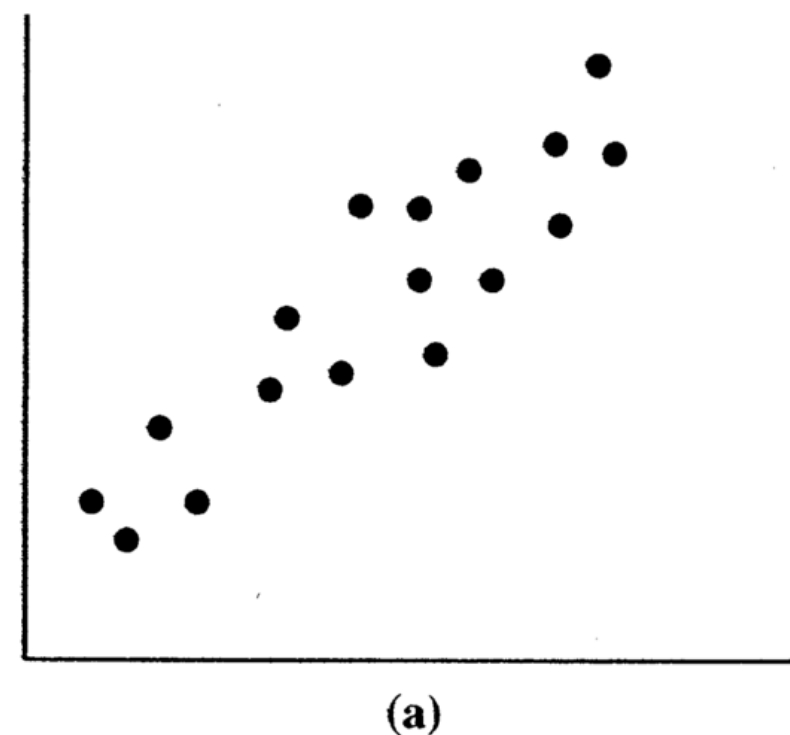
# COVARIANCE AND CORRELATION

# COVARIANCE

- ▶ Measures how variables  $X$  and  $Y$  vary together:

$$\begin{aligned} COV(X, Y) &= E[(X - E[X])(Y - E[Y])] \\ &= E[XY] - E[X]E[Y] \end{aligned}$$

- ▶ Positive if large values of  $X$  are associated with large values of  $Y$
- ▶ Negative if large values of  $X$  are associated with small values of  $Y$



Measures **linear** relationship

# COVARIANCE

- ▶ For discrete random variable pair  $(X, Y)$  that can take on the values of  $(x_i, y_i)$  for  $i=1, \dots, n$  with equal probabilities  $1/n$ :


$$COV(X, Y) = \frac{1}{n} \sum_{i=1}^n (x_i - E[X])(y_i - E[Y])$$

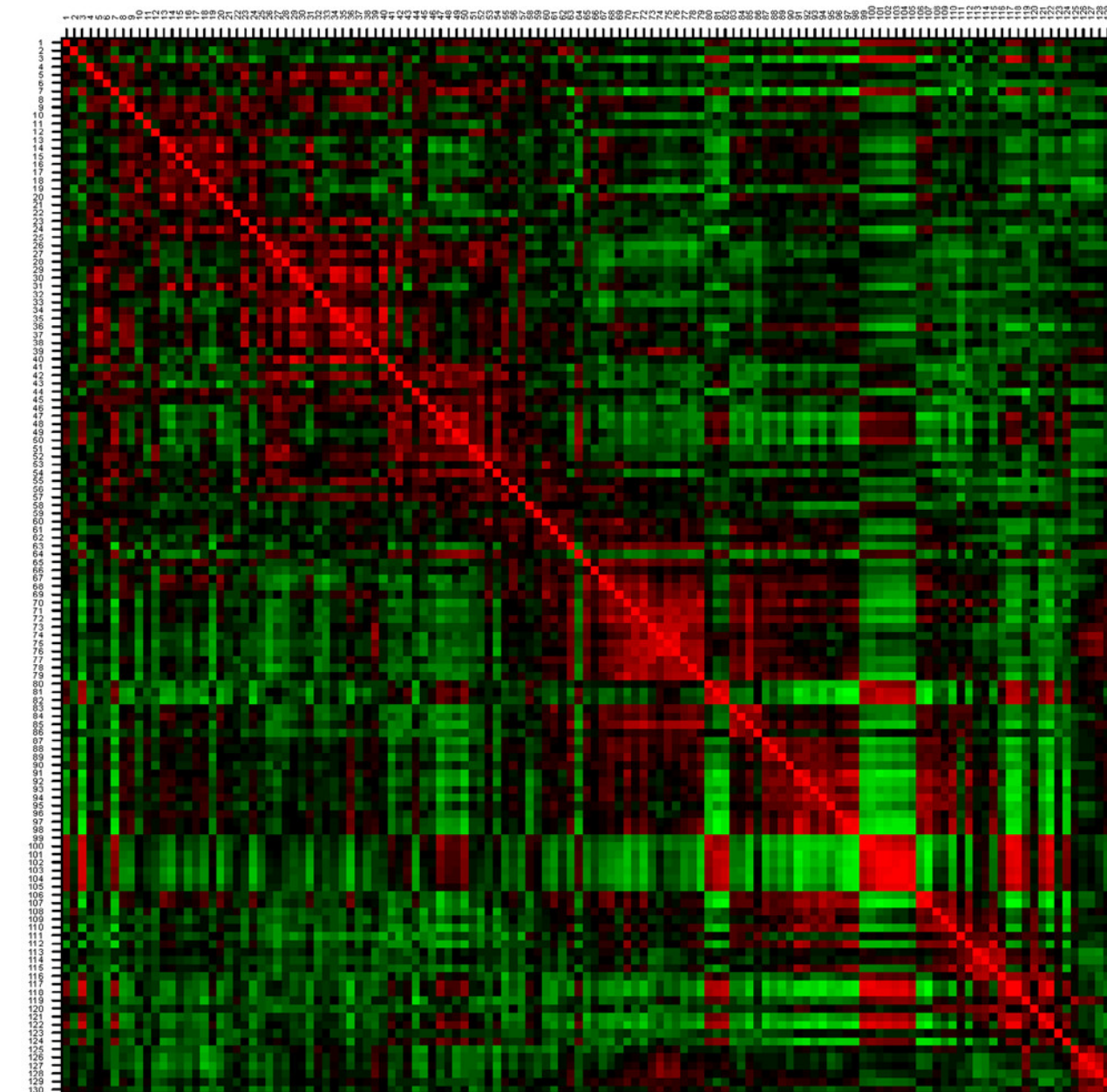
- ▶ Covariance matrix ( $\Sigma$ )
  - ▶ Symmetric matrix of covariances for  $p$  variables
  - ▶  $\Sigma_{ij} = COV(X_i, X_j)$

# CORRELATION COEFFICIENT

- ▶ Covariance depends on ranges of  $X_j$  and  $X_k$
- ▶ Correlation standardizes covariance by dividing through standard deviations

$$\rho(X_j, X_k) = \frac{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{X}_j)(x_{ik} - \bar{X}_k)}{\sigma_{X_j} \sigma_{X_k}}$$

- ▶ Correlation matrix 
  - ▶ Symmetric matrix of correlations for  $p$  variables
  - ▶ What values are on the diagonal?

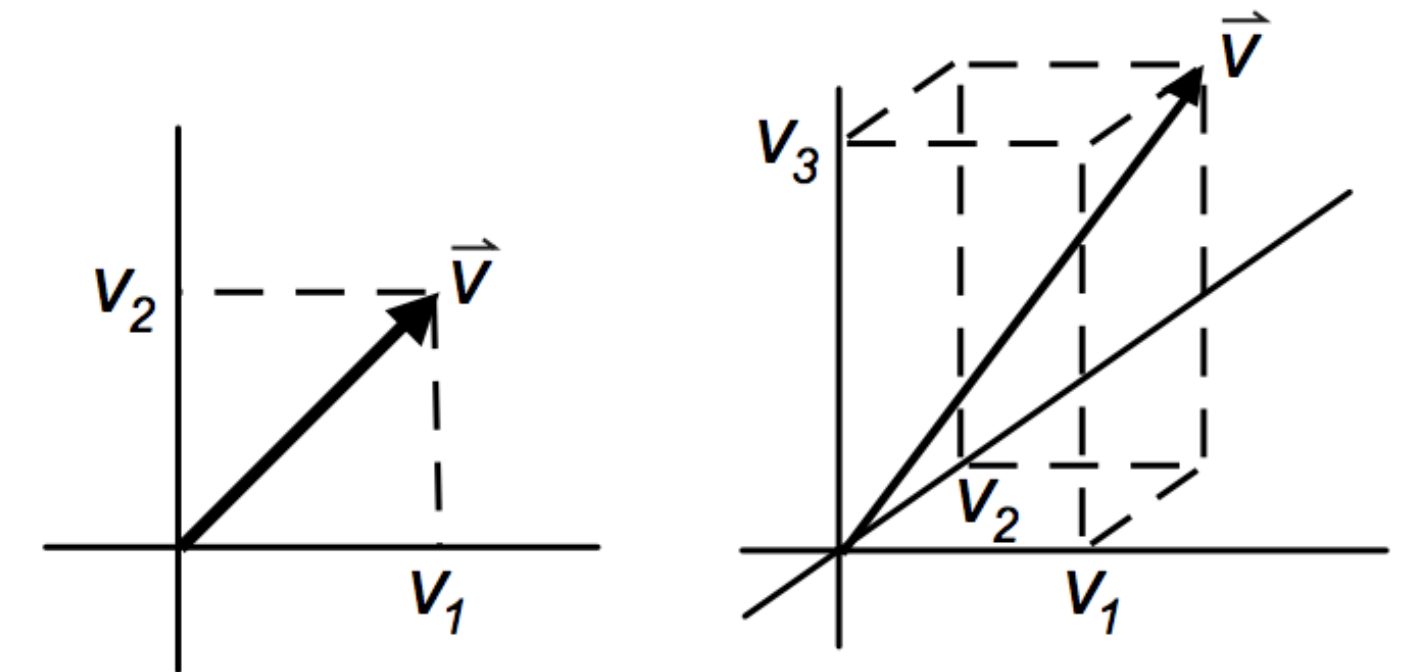


# LINEAR ALGEBRA

# VECTORS

- ▶ A **vector** is a 1D array of values
- ▶ We use the notation  $x_i$  to denote the  $i$ th entry of  $x$
- ▶ Vectors can be graphically depicted as arrows in  $n$ -dimensional space

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$



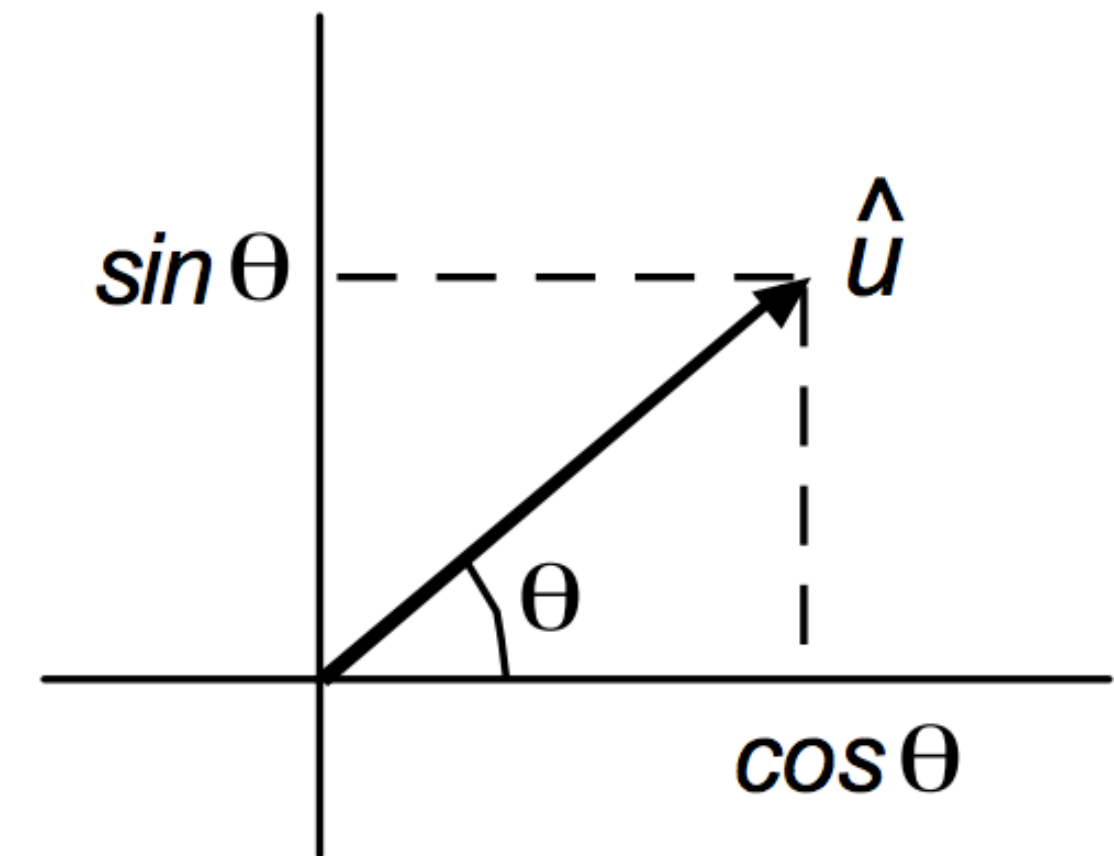
- ▶ The **norm** (length) of a vector is defined as  $\|x\| = \sqrt{\sum_{i=1}^n x_i^2}$



## MORE ON VECTORS

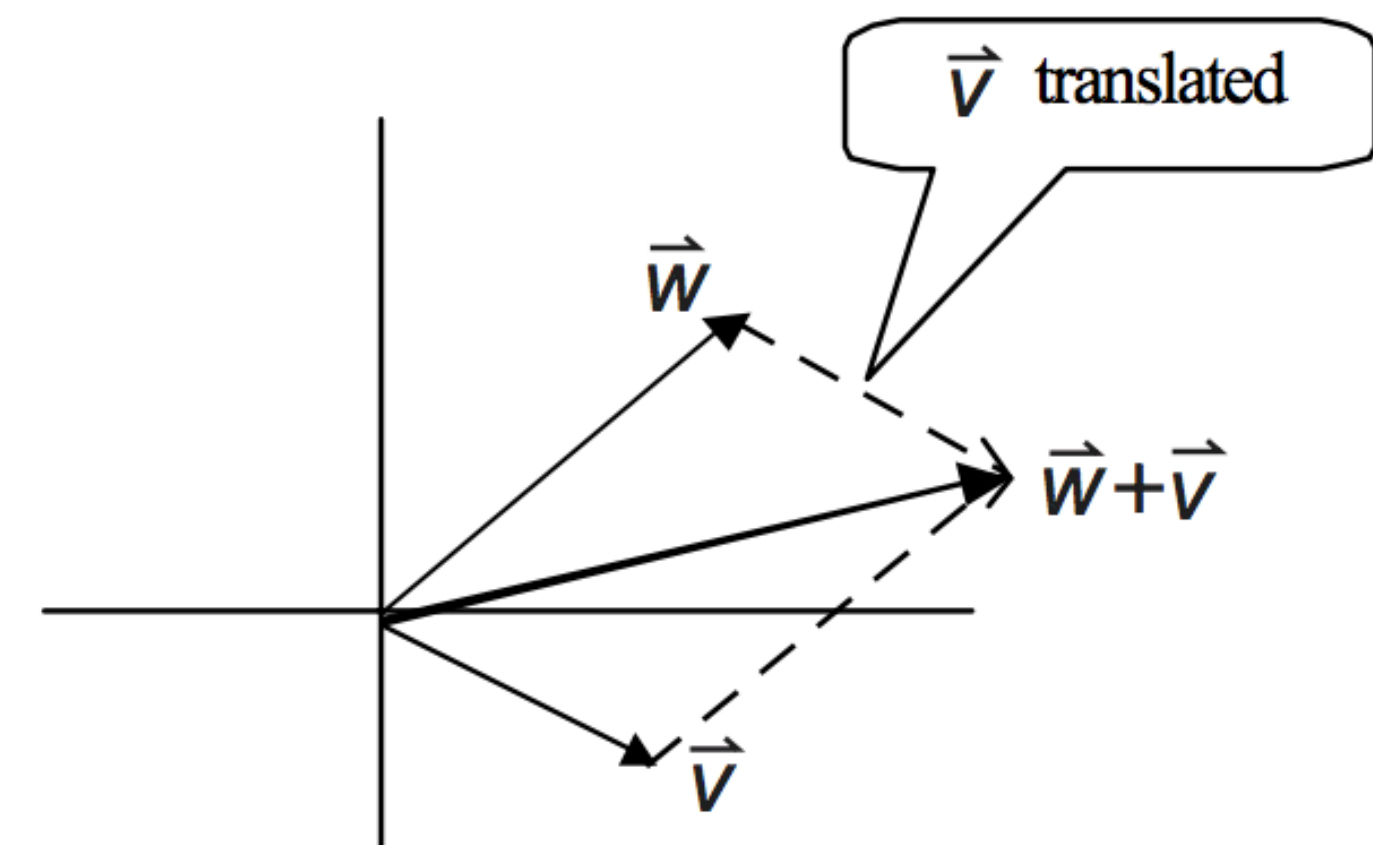
- ▶ A **unit vector** is a vector of length 1. A 2-D unit vector can be parameterized as:

$$\hat{u}(\theta) = \begin{pmatrix} \cos(\theta) \\ \sin(\theta) \end{pmatrix}$$



- ▶ Multiplying a vector by a scalar simply changes the length of the vector by that factor  $\|a\mathbf{x}\| = |a|\|\mathbf{x}\|$  (when  $a$  is negative, the direction of the vector is reversed)

- ▶ Vector addition:  $\mathbf{z} = \mathbf{w} + \mathbf{v} \Leftrightarrow z_i = w_i + v_i$





# INNER PRODUCT

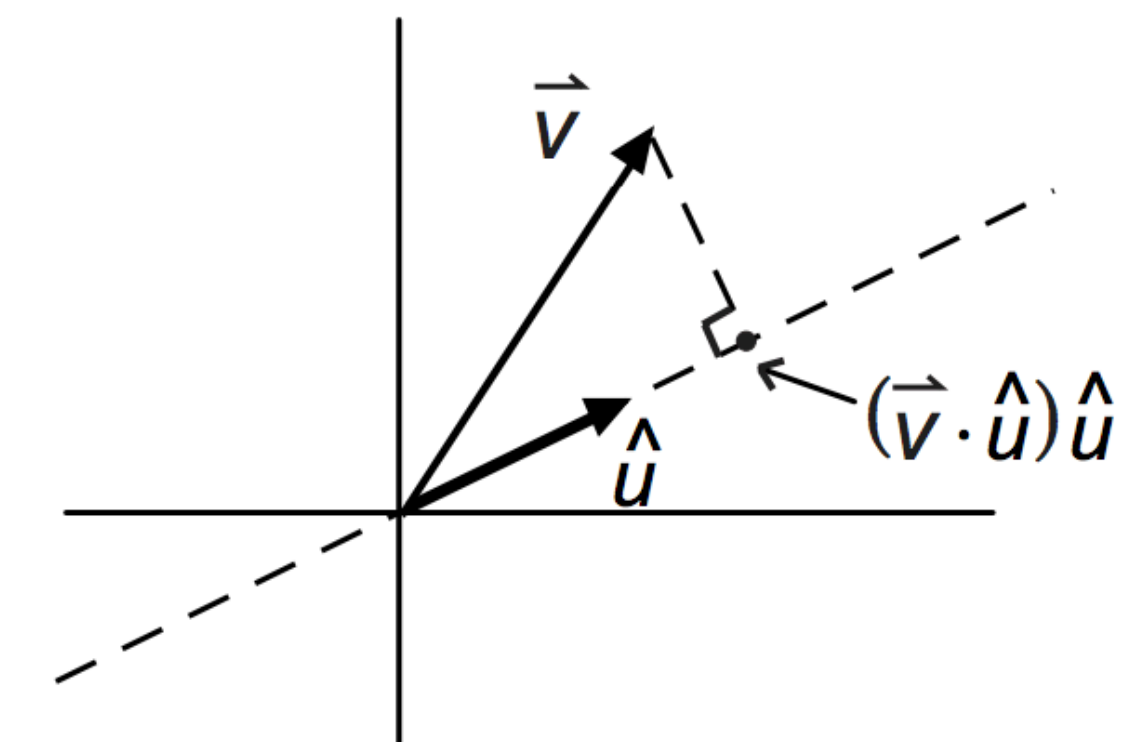
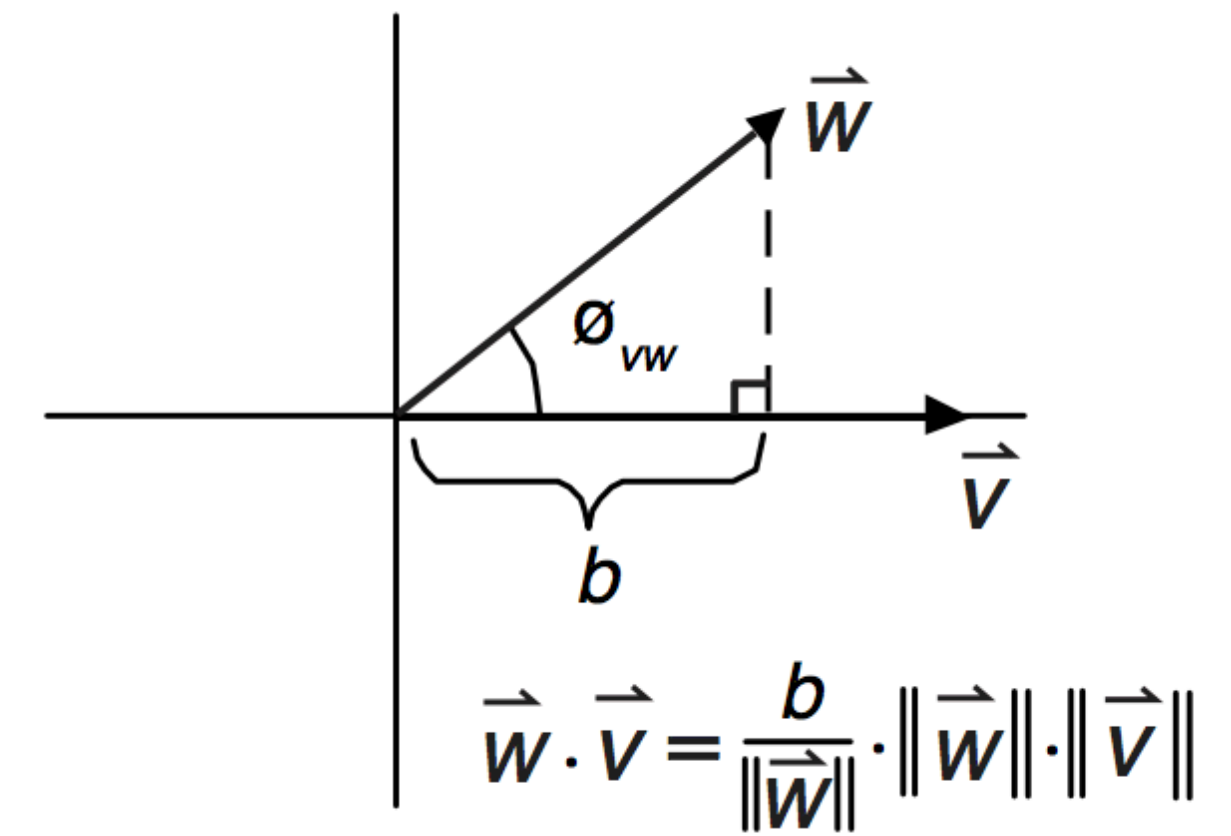
- ▶ The **inner product** of two vectors is the sum of pairwise product of components

$$w \cdot v = \sum_{i=1}^n w_i v_i$$

- ▶ Its equivalent geometric definition is:

$$w \cdot v = \|w\| \|v\| \cos(\phi_{vw})$$

- ▶ The inner product of a vector  $v$  with a unit vector  $u$  is the length of  $v$ 's projection on  $u$ .
- ▶ Two vectors are *orthogonal* to each other if their inner product is 0.



# VECTOR SPACE

- ▶ A vector space can be **spanned** by a set of vectors iff one can write any vector in the vector space as a linear combination of the set
  - ▶ Can the 3D vector space be spanned by  $(1, 1, 0)$  and  $(0, 2, 3)$ ?
- ▶ A set of vectors  $\{v_1, v_2, \dots, v_n\}$  is linearly independent iff the only solution to the following equation is  $\alpha_k = 0$  (for all  $k$ )

$$\sum_{k=1}^n \alpha_k v_k = 0$$

# BASIS

- ▶ A **basis** for a vector space is a linearly independent spanning set.
  - ▶ Is  $(1, 1, 0), (0, 2, 3), (0, 1, 0), (2, 5, 3)$  a basis for the 3D vector space?
- ▶ The **standard basis** of a vector space is the set of unit vectors that lie along the axes of the space
  - ▶  $e_1=(1, 0, \dots, 0), e_2=(0, 1, \dots, 0), \dots, e_n=(0, 0, \dots, 1)$

# MATRICES

- ▶ A **matrix** is a 2D array of values
- ▶ We use  $A_{ij}$  to denote the entry in row  $i$  and column  $j$
- ▶ Higher dimensional matrices are called tensors

$$A = \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1n} \\ A_{21} & A_{22} & \cdots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{m1} & A_{m2} & \cdots & A_{mn} \end{bmatrix}$$

## BASIC MATRIX OPERATIONS

For  $A, B \in \mathbb{R}^{m \times n}$ , matrix addition/subtraction is just the elementwise addition or subtraction of entries

$$C \in \mathbb{R}^{m \times n} = A + B \iff C_{ij} = A_{ij} + B_{ij}$$

For  $A \in \mathbb{R}^{m \times n}$ , transpose is an operator that “flips” rows and columns

$$C \in \mathbb{R}^{n \times m} = A^T \iff C_{ji} = A_{ij}$$

For  $A \in \mathbb{R}^{m \times n}$ ,  $B \in \mathbb{R}^{n \times p}$  matrix multiplication is defined as

$$C \in \mathbb{R}^{m \times p} = AB \iff C_{ij} = \sum_{k=1}^n A_{ik} B_{kj}$$

**Note:** Matrix multiplication is associative ( $A(BC) = (AB)C$ ), distributive ( $A(B + C) = AB + AC$ ), *not commutative* ( $AB \neq BA$ )

# SPECIAL TYPES OF MATRICES

- ▶ A **square matrix** is a matrix with the same number of rows and columns
- ▶ A **diagonal matrix** is a matrix for which all entries outside the main diagonal are zero

# IDENTITY AND INVERSE MATRIX

The identity matrix  $I \in \mathbb{R}^{n \times n}$  is a square matrix with ones on diagonal and zeros elsewhere, has property that for  $A \in \mathbb{R}^{m \times n}$

$$AI = IA = A \text{ (for different sized } I\text{)}$$



# ORTHOGONAL MATRIX

- ▶ An **orthogonal matrix** is a square matrix for which every column is a unit vector, and every pair of columns is orthogonal.
- ▶ If  $A$  is an orthogonal matrix, then

$$A^T A = I \quad \text{and} \quad A^{-1} = A^T \quad \text{and} \quad A A^T = I$$

# ORTHOGONAL MATRIX

- ▶ If  $A$  is an orthogonal matrix, then  $A^T A = I$  and  $A^{-1} = A^T$  and  $AA^T = I$
- ▶ So  $A^T$  is also an orthogonal matrix, which means that every row of  $A$  is a unit vector and every pair of rows of  $A$  is orthogonal.

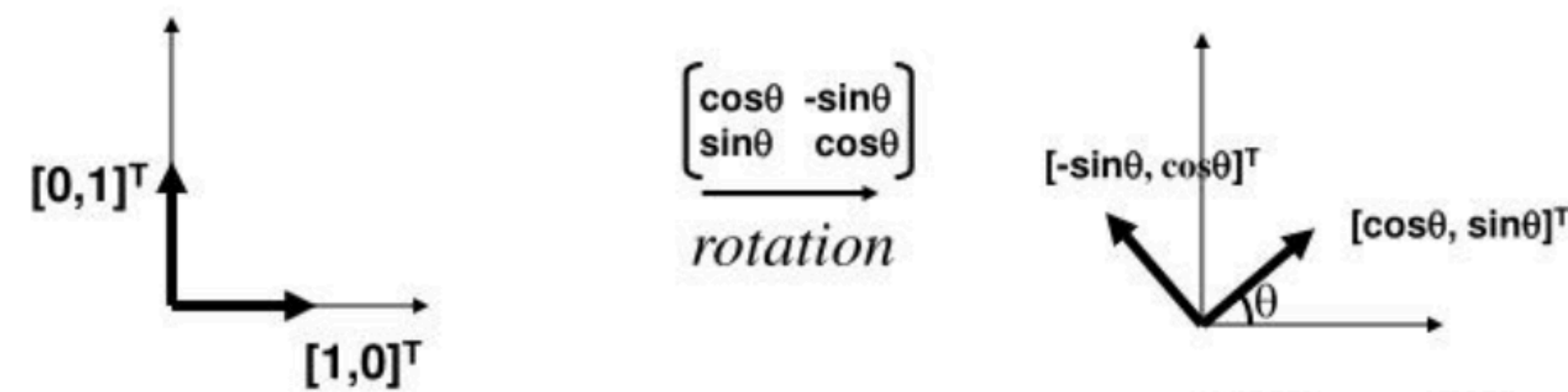
# OTHER DEFINITIONS/PROPERTIES

Transpose of matrix multiplication,  $A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{n \times p}$   
$$(AB)^T = B^T A^T$$

Inverse of product,  $A \in \mathbb{R}^{n \times n}, B \in \mathbb{R}^{n \times n}$  *both square and invertible*  
$$(AB)^{-1} = B^{-1} A^{-1}$$

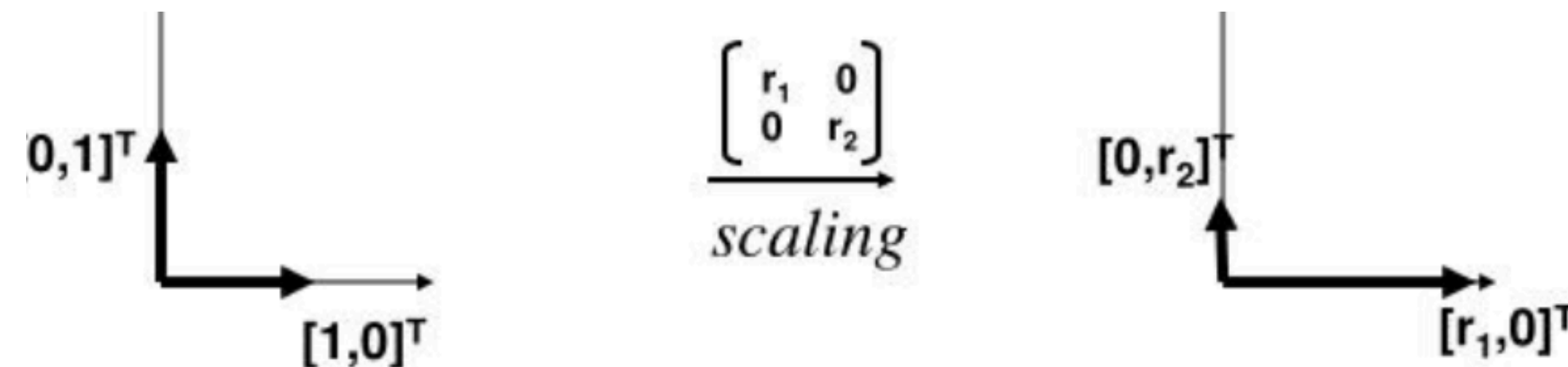
## REPRESENTING LINEAR TRANSFORMATION USING MATRICES

- ▶ When  $A$  is an orthogonal matrix,  $Ax$  rotates  $x$



Can also be  
interpreted as  
change of basis

- ▶ When  $A$  is a diagonal matrix,  $Ax$  stretch or squeeze the axes



- ▶ More general square matrix involves both rotation and scaling

# EIGENVALUES AND EIGENVECTORS

- ▶ An **eigenvector** is a non-zero vector that changes by only a scalar factor when a particular linear transformation is applied to it, and the scalar is **eigenvalue**.

$$Ax = \lambda x$$

- ▶ How to calculate eigenvalues and eigenvectors?
  - ▶  $(A - \lambda I)x = 0$ . Let the determinant of  $A - \lambda I$  be 0.

# EIGENDECOMPOSITION

- ▶ Let  $A$  be a square matrix with  $N$  linearly independent eigenvectors,  $q_i$  ( $i=1, \dots, N$ ). Then  $A$  can be factorized as:
  - ▶  $A = Q\Lambda Q^{-1}$
  - ▶  $Q$  is the square matrix whose  $i$ -th column is the eigenvector  $q_i$  of  $A$ ,  $\Lambda$  is the diagonal matrix whose diagonal elements are the corresponding eigenvalues, i.e.,  $\Lambda_{ii} = \lambda_i$

# DERIVATION OF EIGENDECOMPOSITION

- ▶ Let  $A$  be a square matrix with  $N$  linearly independent eigenvectors,  $q_i$  ( $i=1, \dots, N$ ). Then  $A$  can be factorized as:  $A = Q\Lambda Q^{-1}$



# EIGENDECOMPOSITION

- ▶ Let  $A$  be a square matrix with  $N$  linearly independent eigenvectors,  $q_i$  ( $i=1, \dots, N$ ). Then  $A$  can be factorized as:
  - ▶  $A = Q\Lambda Q^{-1}$
  - ▶  $Q$  is the square matrix whose  $i$ -th column is the eigenvector  $q_i$  of  $A$ ,  $\Lambda$  is the diagonal matrix whose diagonal elements are the corresponding eigenvalues, i.e.,  $\Lambda_{ii} = \lambda_i$
  - ▶ For a symmetric matrix  $A$ ,  $Q$  is an orthogonal matrix, that is,  $A = Q\Lambda Q^T$

# SINGULAR VALUE DECOMPOSITION (SVD)

- ▶ A rectangular matrix  $A$  can be broken down into the product of three matrices: an orthogonal matrix  $U$ , a diagonal matrix  $S$ , and the transpose of an orthogonal matrix  $V$ .

The diagram illustrates the SVD decomposition of a matrix  $A$  into three matrices:  $U$ ,  $S$ , and  $V^T$ . Matrix  $A$  is represented by a blue rectangle with dimensions  $m$  (height) and  $n$  (width). It is equal to the product of matrix  $U$  (dimensions  $m \times m$ ), matrix  $S$  (dimensions  $m \times n$ ), and matrix  $V^T$  (dimensions  $n \times n$ ). The matrices are represented by blue rectangles, and the dimensions are indicated by labels  $m$  and  $n$  around each rectangle. The decomposition is shown as  $A = U * S * V^T$ .

# DERIVATION OF SINGULAR VALUE DECOMPOSITION (SVD)

- ▶ Columns of  $U$  are eigenvectors of  $AA^T$
- ▶ Columns of  $V$  are eigenvectors of  $A^TA$
- ▶ Diagonal entries of  $S$  are the square roots of the non-zero eigenvalues of  $AA^T$  (as well as  $A^TA$ )