



---

# Machine Learning: Elements of **Linear Algebra**

---

Alexander G. Ororbia II  
Introduction to Machine Learning  
CSCI-635  
8/30/2023

# This is a Crash Course Review...

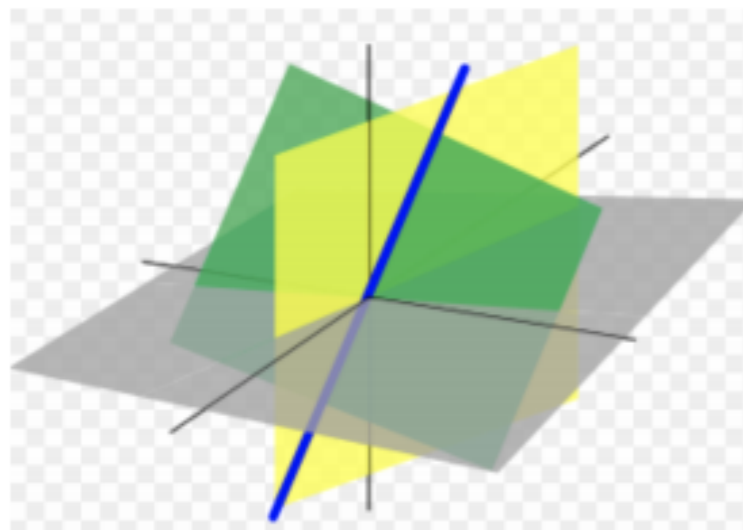
- Not a comprehensive survey of all of linear algebra
- Focused on subset most relevant to machine learning
- For a larger subset/treatment:  
***Linear Algebra* by Georgi E. Shilov**
- Also read: “**Linear Algebra for Dummies**”

# What is Linear Algebra?

- Linear algebra is the branch of mathematics concerning linear equations such as

$$a_1x_1 + \dots + a_nx_n = b$$

- In vector notation we say  $a^T x = b$
  - Called a linear transformation of  $x$
- Linear algebra is fundamental to geometry, for defining objects such as lines, planes, rotations



Linear equation  $a_1x_1 + \dots + a_nx_n = b$   
defines a plane in  $(x_1, \dots, x_n)$  space  
Straight lines define common solutions  
to equations

Linear algebra is used  
throughout engineering  
(based on continuous math)

$\subset$  = proper subset (not the whole thing)  $\subseteq$  = subset

$\exists$  = there exists

$\forall$  = for every

$\in$  = element of

$\cup$  = union (or)

$\cap$  = intersection (and)

s.t. = such that

w.r.t. = with respect to

$\implies$  implies

$\iff$  if and only if

$\sum$  = sum

$\prod$  = multiplication

$\setminus$  = set minus

$\therefore$  = therefore

## Mathematical Notation in CSCI 635

# Scalars

- A scalar is a single number
- Integers, real numbers, rational numbers, etc.
- Denoted with italic font:

*a, n, x*

## On number spaces/domains:

$\mathbb{R}$   $\mathcal{R}$  All real (continuous) numbers

$\mathbb{Z}$   $\mathcal{Z}$  All integers

$\mathbb{N}$  All natural (counting) numbers

# Vectors

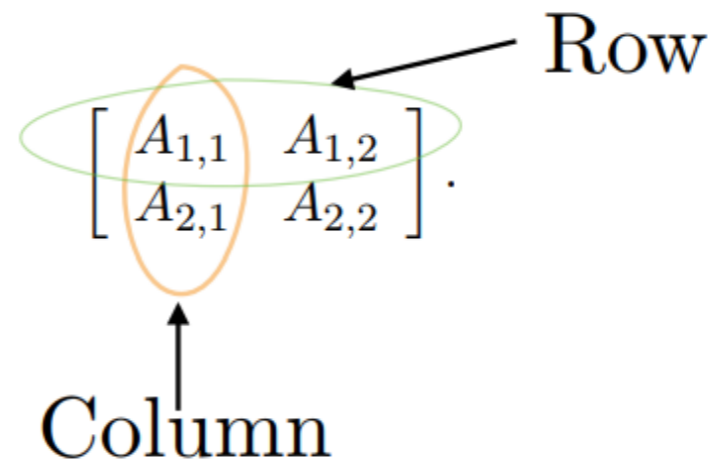
- An array of numbers arranged in order
- Each no. identified by an index
- Written in lower-case bold such as  $\mathbf{x}$ 
  - its elements are in italics lower case, subscripted

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}, \quad \mathcal{R}^n = \mathcal{R}^{n \times 1} \quad \text{“All real numbers”}$$

- If each element is in  $R$  then  $\mathbf{x}$  is in  $R^n$
- We can think of vectors as points in space
  - Each element gives coordinate along an axis

# Matrices

- 2-D array of numbers
  - So each element identified by two indices
- Denoted by bold typeface  $\mathbf{A}$ 
  - Elements indicated by name in italic but not bold
    - $A_{1,1}$  is the top left entry and  $A_{m,n}$  is the bottom right entry
    - We can identify nos in vertical column  $j$  by writing  $:$  for the horizontal coordinate



The "slice" operator

- $A_{i:}$  is  $i^{\text{th}}$  row of  $A$ ,  $A_{:j}$  is  $j^{\text{th}}$  column of  $A$
- If  $A$  has shape of height  $m$  and width  $n$  with real-values then  $\mathbf{A} \in \mathbb{R}^{m \times n}$

# Tensor

- Sometimes need an array with more than two axes
  - E.g., an RGB color image has three axes
- A tensor is an array of numbers arranged on a regular grid with variable number of axes
- Denote a tensor with this bold typeface: **A**
- Element  $(i, j, k)$  of tensor denoted by  $A_{i, j, k}$

Types of notation accepted (just be consistent & mean what you write):

$$\mathcal{R}^{1 \times 1 \times 1} = \mathcal{R}^{1 \times 1} = \mathcal{R}^1 = \mathcal{R}$$



# Shapes of Tensors



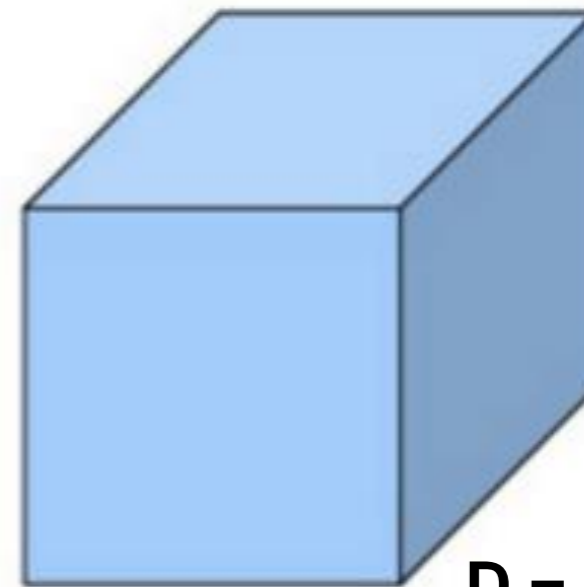
**D = 1**

1d-tensor



**D = 2**

2d-tensor



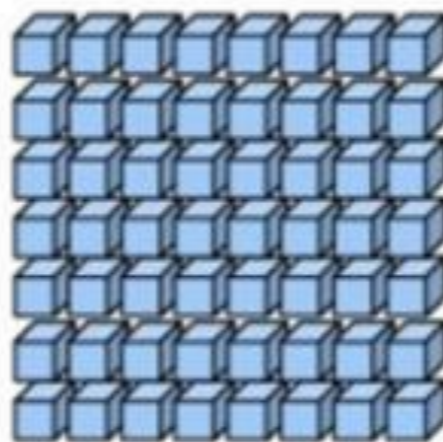
**D = 3**

3d-tensor



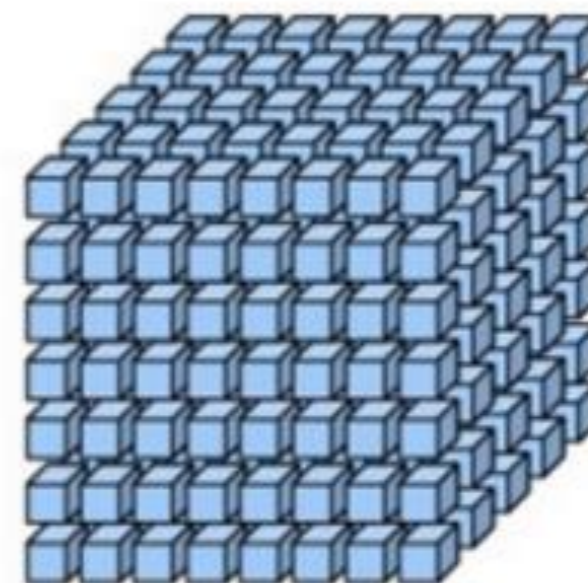
**D = 4**

4d-tensor



**D = 5**

5d-tensor



**D = 5**

6d-tensor

# Transpose of a Matrix

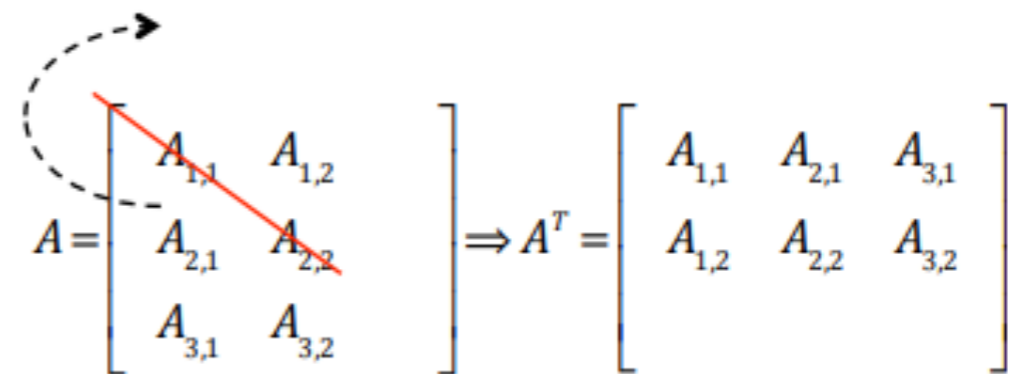
- An important operation on matrices
- The transpose of a matrix  $\mathbf{A}$  is denoted as  $\mathbf{A}^T$
- Defined as

$$(\mathbf{A}^T)_{i,j} = A_{j,i}$$

– The mirror image across a diagonal line

- Called the main diagonal, running down to the right starting from upper left corner

$$A = \begin{bmatrix} A_{1,1} & A_{1,2} & A_{1,3} \\ A_{2,1} & A_{2,2} & A_{2,3} \\ A_{3,1} & A_{3,2} & A_{3,3} \end{bmatrix} \Rightarrow A^T = \begin{bmatrix} A_{1,1} & A_{2,1} & A_{3,1} \\ A_{1,2} & A_{2,2} & A_{3,2} \\ A_{1,3} & A_{2,3} & A_{3,3} \end{bmatrix}$$


$$A = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \\ A_{3,1} & A_{3,2} \end{bmatrix} \Rightarrow A^T = \begin{bmatrix} A_{1,1} & A_{2,1} & A_{3,1} \\ A_{1,2} & A_{2,2} & A_{3,2} \end{bmatrix}$$

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T \quad \leftarrow \text{Useful property of transpose}$$

# Vector as Special Case of Matrix

- Vectors are matrices with a single column
- Often written in-line using transpose

$$\mathbf{x} = [x_1, \dots, x_n]^T$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \Rightarrow \mathbf{x}^T = [x_1, x_2, \dots, x_n]$$

- A scalar is a matrix with one element

$$a = a^T \longleftarrow \mathcal{R}^{1 \times 1}$$

# Matrix Addition/Subtraction

- Assume column-major matrices (for efficiency)
- Add/subtract operators follow basic properties of normal add/subtract
  - Matrix A + Matrix B is computed element-wise

0.5	-0.7
-0.69	1.8

 + 

0.5	-0.7
-0.69	1.8

 = 

.5 + .5 = 1.0	-.7 - .7 = -1.4
-.69 - .69 = -1.38	1.8 + 1.8 = 3.6

# Matrix Addition

- We can add matrices to each other if they have the same shape, by adding corresponding elements

– If  $A$  and  $B$  have same shape (height  $m$ , width  $n$ )

$$C = A + B \Rightarrow C_{i,j} = A_{i,j} + B_{i,j}$$

- A scalar can be added to a matrix or multiplied by a scalar

$$D = aB + c \Rightarrow D_{i,j} = aB_{i,j} + c$$

- Less conventional notation used in ML:

– Vector added to matrix  $C = A + \mathbf{b} \Rightarrow C_{i,j} = A_{i,j} + b_j$

- Called broadcasting since vector  $\mathbf{b}$  added to each row of  $A$

# Multiplying Matrices

- For product  $C=AB$  to be defined,  $A$  has to have the same no. of columns as the no. of rows of  $B$ 
  - If  $A$  is of shape  $m \times n$  and  $B$  is of shape  $n \times p$  then *matrix product*  $C$  is of shape  $m \times p$

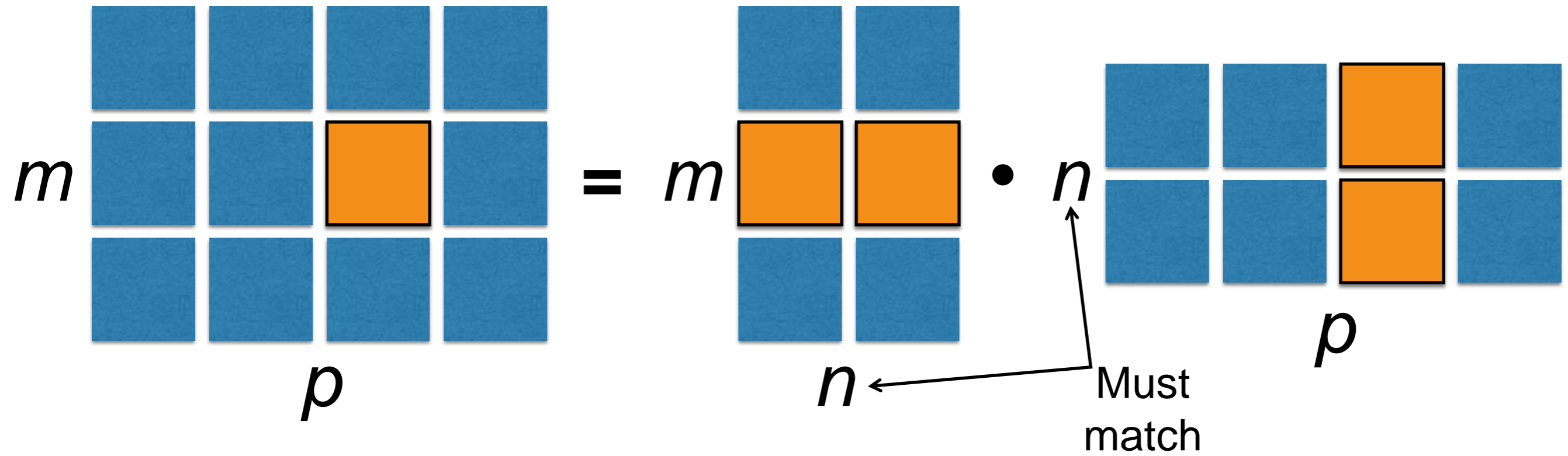
$$C = AB \Rightarrow C_{i,j} = \sum_k A_{i,k} B_{k,j}$$

- Dot product between two vectors  $\mathbf{x}$  and  $\mathbf{y}$  of same dimensionality is the matrix product  $\mathbf{x}^T \mathbf{y}$
- We can think of matrix product  $C=AB$  as computing  $C_{ij}$  the dot product of row  $i$  of  $A$  and column  $j$  of  $B$

# Matrix-Matrix Multiply

- Matrix-Matrix multiply (outer product)
  - Vector-Vector multiply (dot product)
- The usual workhorse of machine learning
- Vectorizes sums of products (builds on dot product)

0.5	-0.7	*	0.5	-0.7	=	$(.5 * .5) + (-.7 * -.69)$	$(.5 * -.7) + (-.7 * 1.8)$
-0.69	1.8		-0.69	1.8		$(-.69 * .5) + (1.8 * -.69)$	$(-.69 * -.7) + (1.8 * 1.8)$



**Referred to sometimes as matrix-matrix product or matrix-vector product (or multiply)**



# Matrix Product Properties

- Distributivity over addition:  $A(B + C) = AB + AC$
- Associativity:  $A(BC) = (AB)C$
- Not commutative:  $AB = BA$  is not always true
- Dot product between vectors is commutative:  
 $\mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x}$
- Transpose of a matrix product has a simple form:  $(AB)^T = B^T A^T$

# Hadamard Product

- Multiply each  $A(i, j)$  to each corresponding  $B(i, j)$
- Element-wise multiplication

0.5	-0.7
-0.69	1.8

 $\odot$ 

0.5	-0.7
-0.69	1.8

 = 

$.5 * .5 = .25$	$-.7 * .7 = .49$
$-.69 * -.69 = .4761$	$1.8 * 1.8 = 3.24$

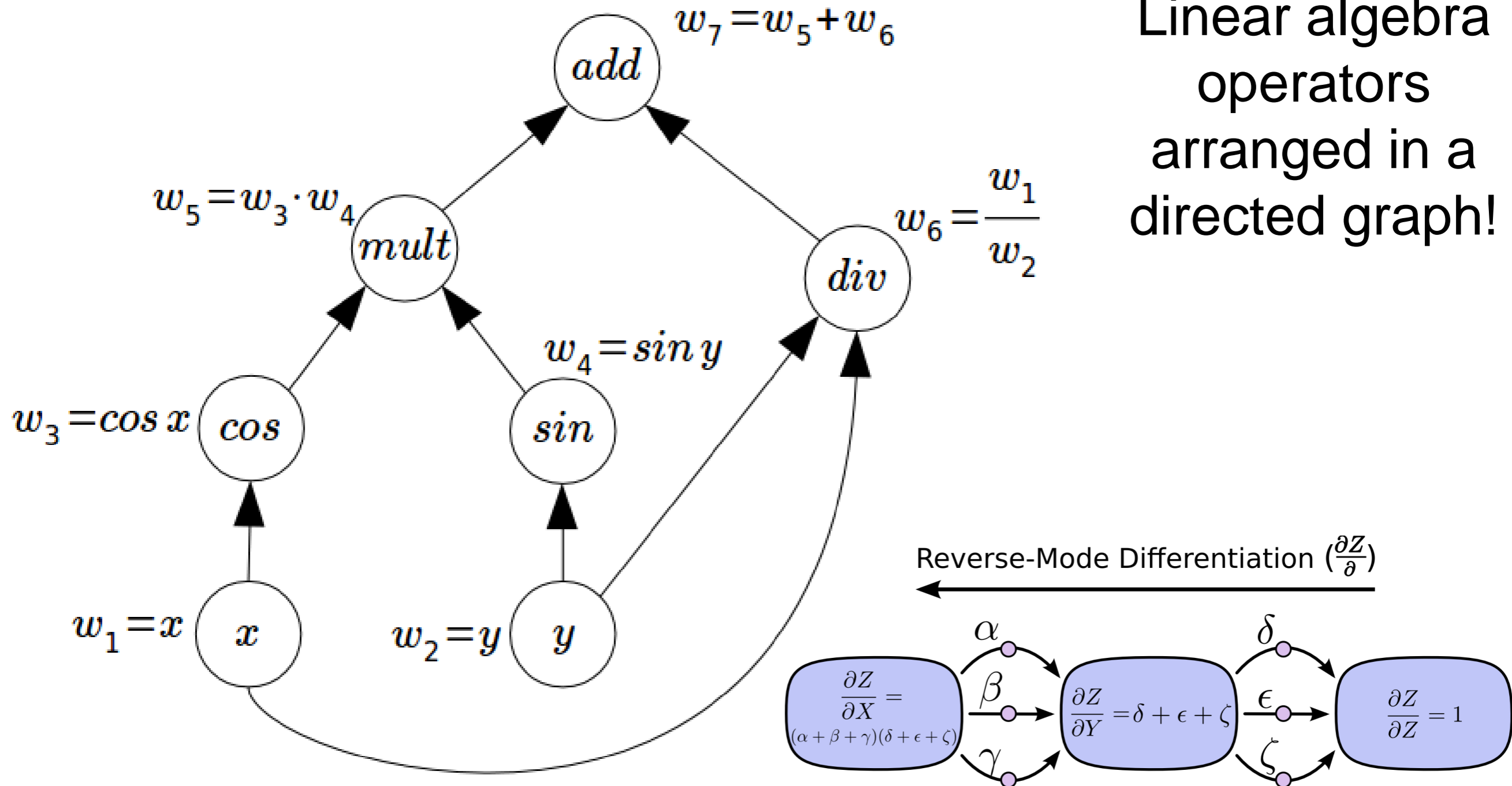
# Elementwise Functions

- Applied to each element (i, j) of matrix/vector argument
  - Could be  $\cos(\cdot)$ ,  $\sin(\cdot)$ ,  $\tanh(\cdot)$ , etc. (the “.” means argument)
- *Identity*:  $\phi(\mathbf{v}) = \mathbf{v}$
- *Logistic Sigmoid*:  $\phi(\mathbf{v}) = \sigma(\mathbf{v}) = \frac{1}{1+e^{-\mathbf{v}}}$
- *Softmax*:  $\phi(\mathbf{v}) = \frac{\exp(\mathbf{v})}{\sum_{c=1}^C \exp(\mathbf{v}_c)}$   $\mathbf{v} \in \mathbb{R}^C$
- *Linear Rectifier*:  $\phi(\mathbf{v}) = \max(0, \mathbf{v})$

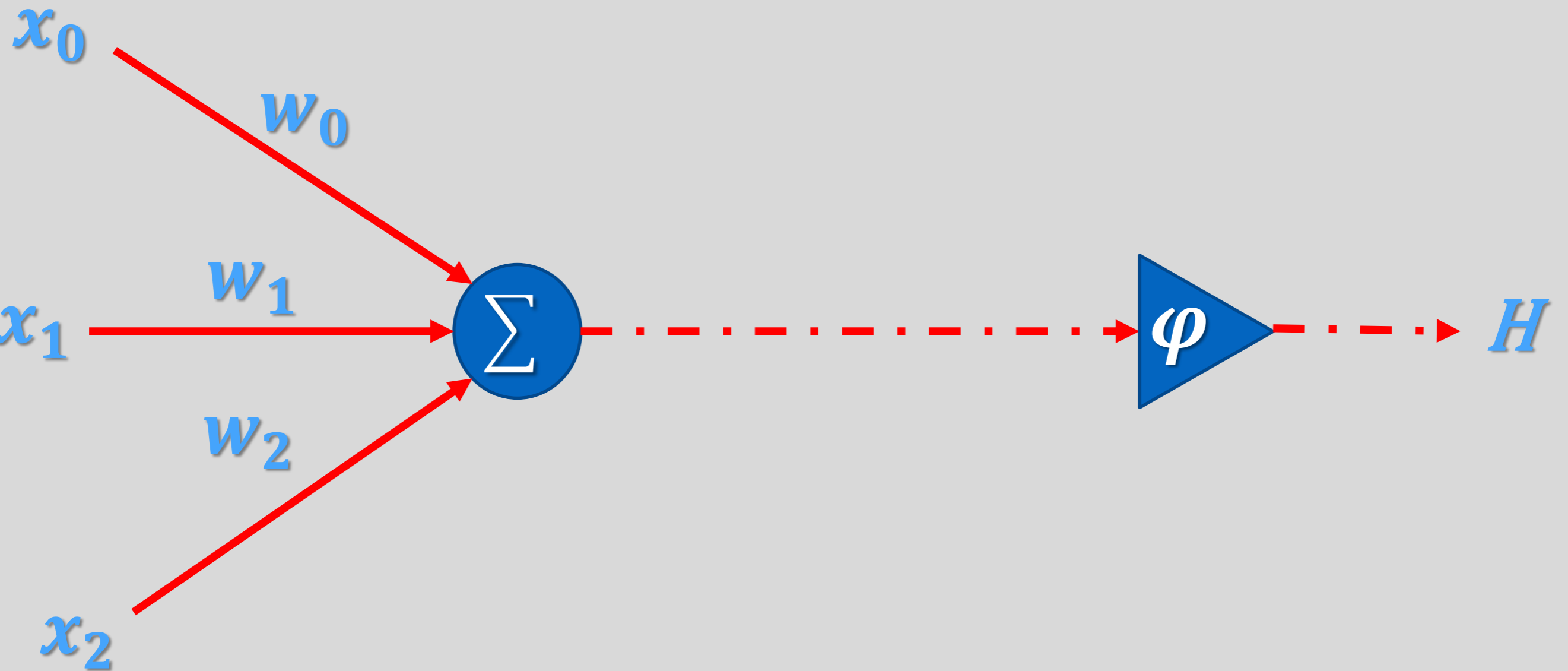
$$\varphi \left( \begin{array}{|c|c|} \hline 1.0 & -1.4 \\ \hline -0.69 & 1.8 \\ \hline \end{array} \right) = \begin{array}{|c|c|} \hline \varphi(1.0) = 1 & \varphi(-1.4) = 0 \\ \hline \varphi(-0.68) = 0 & \varphi(1.8) = 1.8 \\ \hline \end{array}$$

# Why Do We Care? Computation Graphs

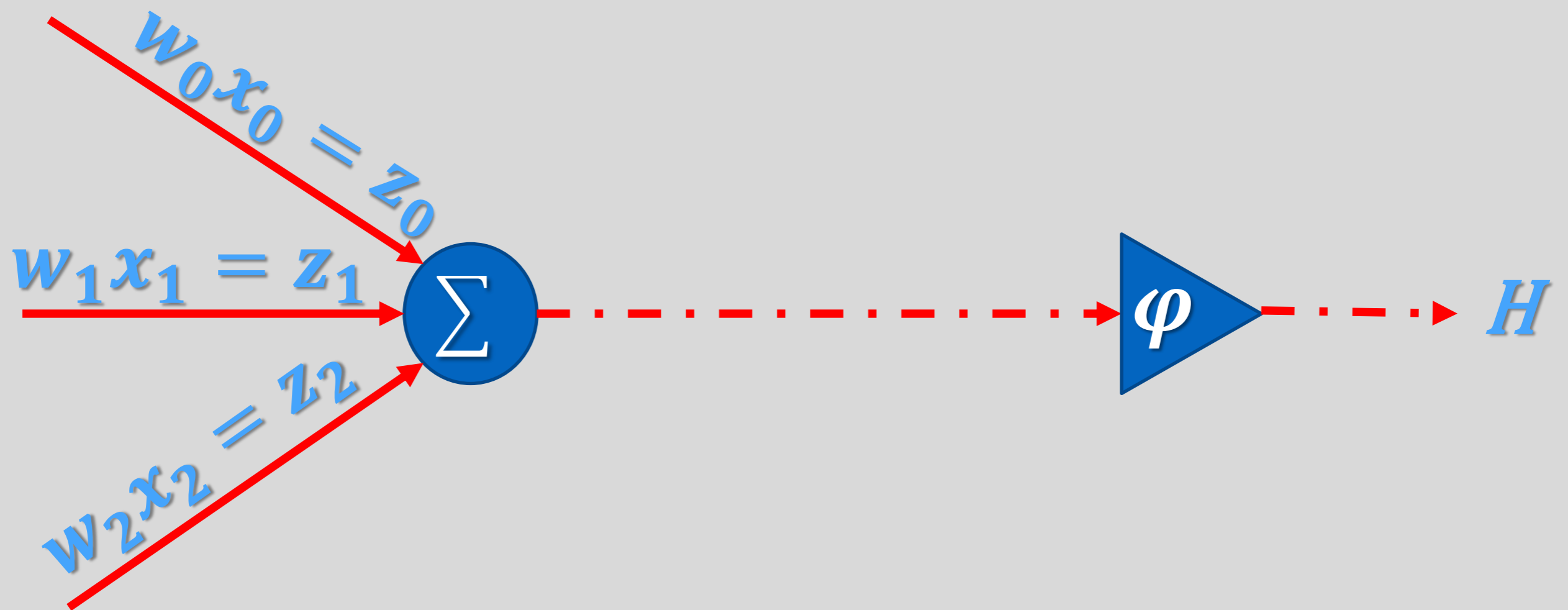
Linear algebra operators arranged in a directed graph!



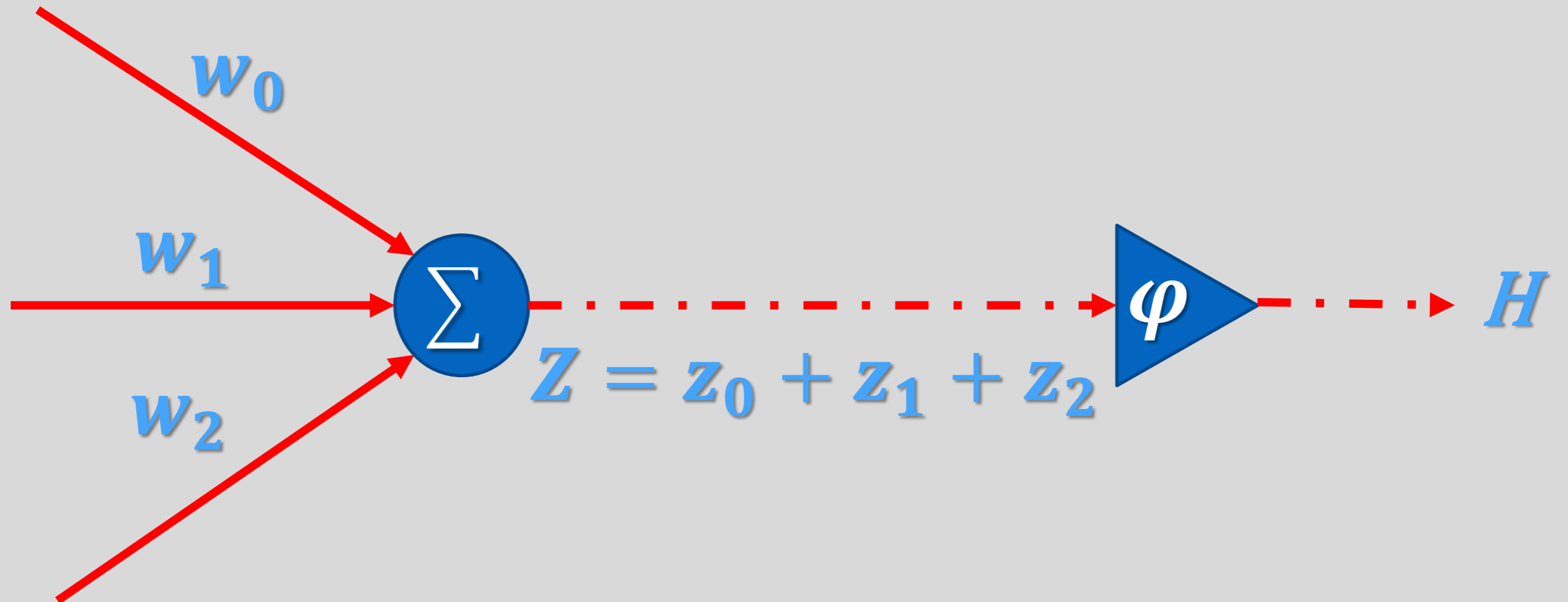
# A Simple Linear Predictor



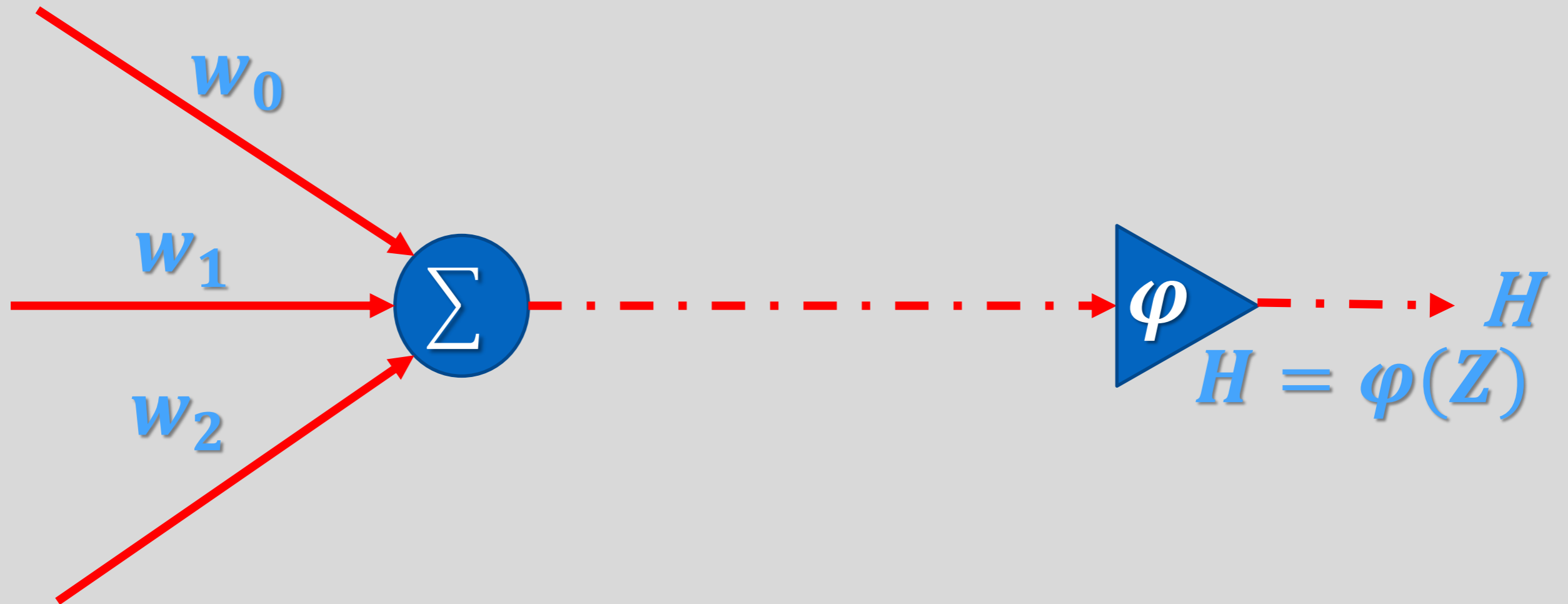
# A Simple Linear Predictor



# A Simple Linear Predictor



# A Simple Linear Predictor





# Vector Form (One Unit)

This calculates activation value of single (output) unit that is connected to 3 (input) sensors.

$$h_0: \begin{array}{|c|c|c|} \hline w_0 & w_1 & w_2 \\ \hline \end{array} * \begin{array}{|c|} \hline x_0 \\ \hline x_1 \\ \hline x_2 \\ \hline \end{array} = \boxed{\varphi(w_0 * x_0 + w_1 * x_1 + w_2 * x_2)}$$

# Vector Form (Two Units)

This vectorization easily generalizes to multiple (3) sensors feeding into multiple (2) units.

$$\begin{array}{l} h_0: \\ h_1: \end{array} \begin{array}{|c|c|c|} \hline w_0 & w_1 & w_2 \\ \hline w_3 & w_4 & w_5 \\ \hline \end{array} * \begin{array}{|c|} \hline x_0 \\ \hline x_1 \\ \hline x_2 \\ \hline \end{array} = \begin{array}{|c|} \hline \varphi(w_0 * x_0 + w_1 * x_1 + w_2 * x_2) \\ \hline \varphi(w_3 * x_0 + w_4 * x_1 + w_5 * x_2) \\ \hline \end{array}$$

***Known as vectorization!***

# Now Let Us Fully Vectorize This!

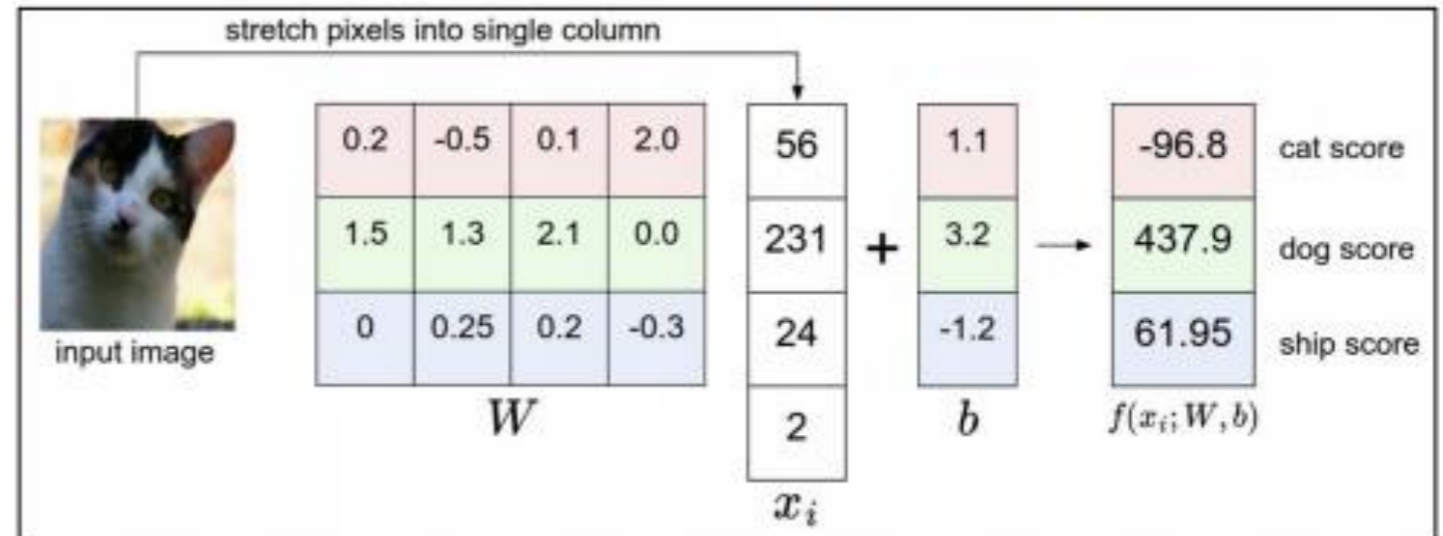
This vectorization is also important for formulating **mini-batches**.  
(Good for GPU-based processing.)

$$\begin{array}{l} h_0: \\ h_1: \end{array} \begin{array}{|c|c|c|} \hline w_0 & w_1 & w_2 \\ \hline w_3 & w_4 & w_5 \\ \hline \end{array} * \begin{array}{|c|c|} \hline x_0 & x_3 \\ \hline x_1 & x_4 \\ \hline x_2 & x_5 \\ \hline \end{array} = \begin{array}{|c|c|} \hline \varphi(w_0 * x_0 + \dots) & \varphi(w_0 * x_3 + \dots) \\ \hline \varphi(w_3 * x_0 + \dots) & \varphi(w_3 * x_3 + \dots) \\ \hline \end{array}$$

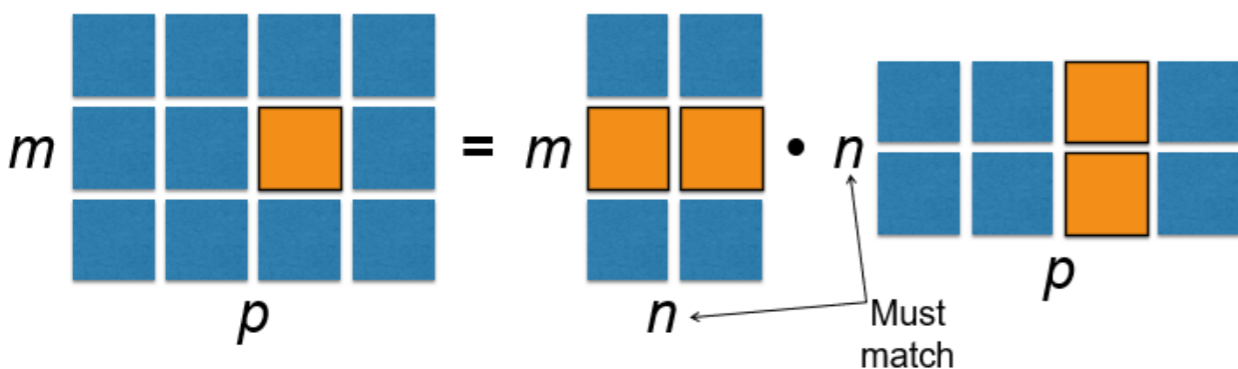
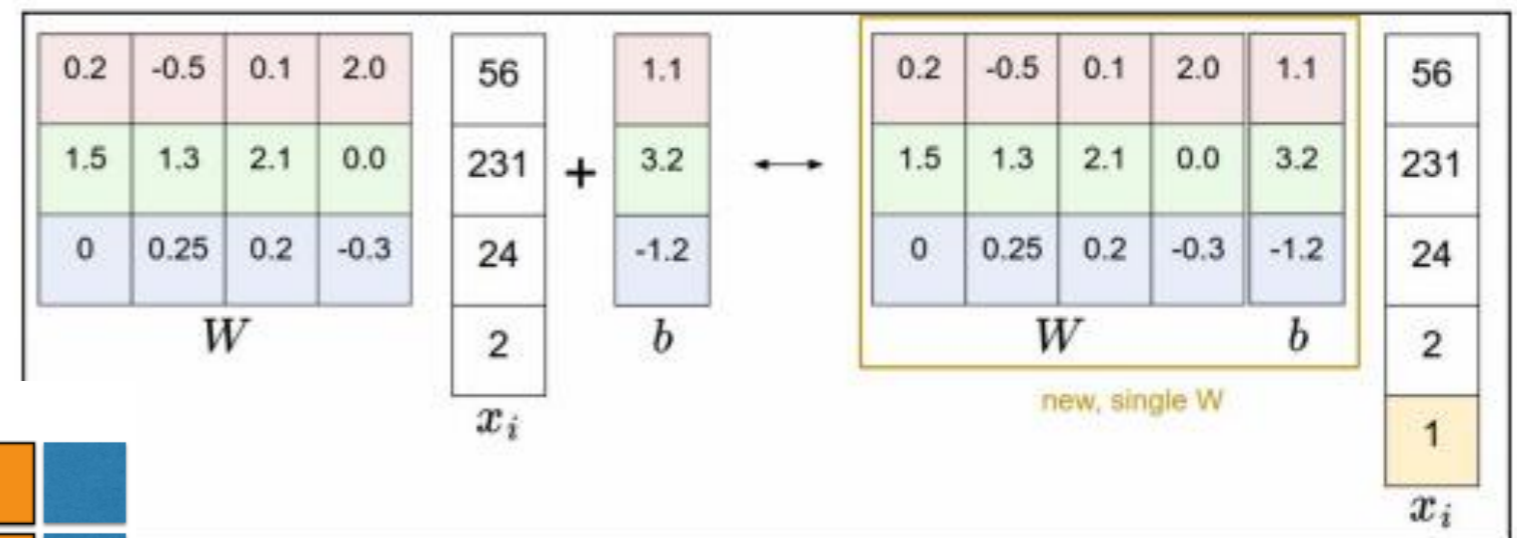
# Tensors in Machine Learning

Vector  $x$  is converted into vector  $y$  by multiplying  $x$  by a matrix  $W$

A linear classifier  $y = Wx + b$



A linear classifier with bias eliminated  $y = Wx$



# Identity and Inverse Matrices

- Matrix inversion is a powerful tool to analytically solve  $A\mathbf{x}=\mathbf{b}$
- Needs concept of Identity matrix
- Identity matrix does not change value of vector when we multiply the vector by identity matrix
  - Denote identity matrix that preserves n-dimensional vectors as  $I_n$
  - Formally  $I_n \in \mathbb{R}^{n \times n}$  and  $\forall \mathbf{x} \in \mathbb{R}^n, I_n \mathbf{x} = \mathbf{x}$
  - Example of  $I_3$   $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

# Norms

- Used for measuring the size of a vector
- Norms map vectors to non-negative values
- Norm of vector  $\mathbf{x} = [x_1, \dots, x_n]^T$  is distance from origin to  $\mathbf{x}$ 
  - It is any function  $f$  that satisfies:

$$f(\mathbf{x}) = 0 \Rightarrow \mathbf{x} = \mathbf{0}$$

$$f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y}) \quad \text{Triangle Inequality}$$

$$\forall \alpha \in \mathbb{R} \quad f(\alpha \mathbf{x}) = |\alpha| f(\mathbf{x})$$

# $L^p$ Norm

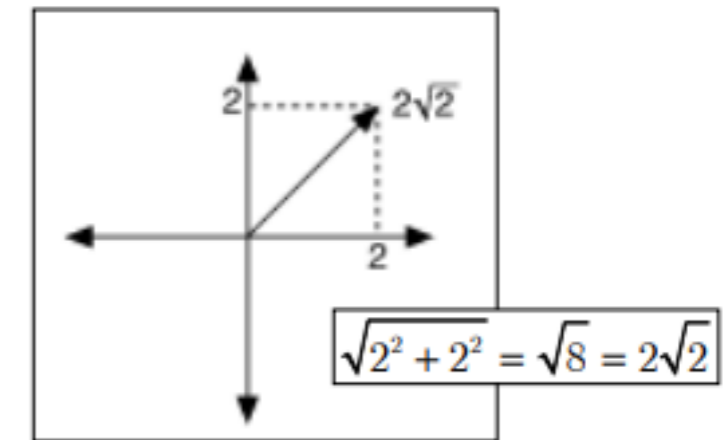
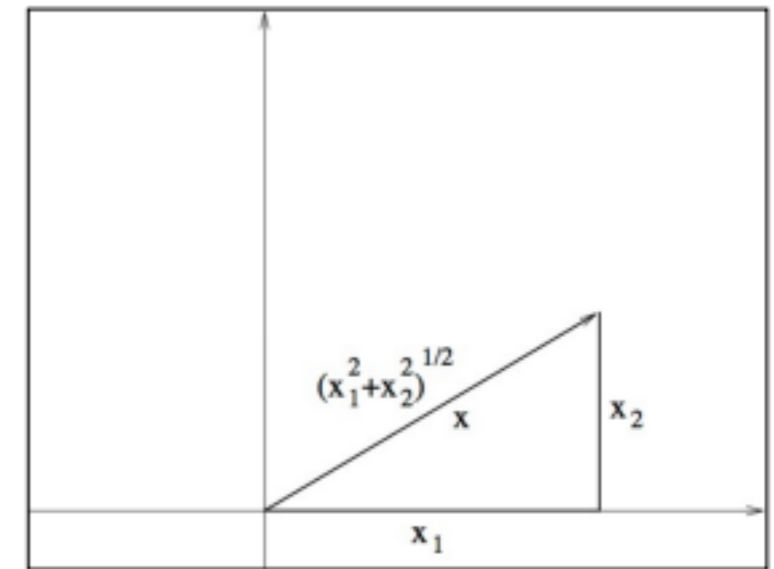
- **Definition:**

$$\|x\|_p = \left( \sum_i |x_i|^p \right)^{\frac{1}{p}}$$

- $L^2$  Norm

- Called Euclidean norm

- Simply the Euclidean distance between the origin and the point  $x$
- written simply as  $\|x\|$
- Squared Euclidean norm is same as  $x^T x$



- $L^1$  Norm

- Useful when 0 and non-zero have to be distinguished

- Note that  $L^2$  increases slowly near origin, e.g.,  $0.1^2=0.01$ )

- $L^\infty$  Norm  $\|x\|_\infty = \max_i |x_i|$

- Called max norm

# The Frobenius Norm

Similar to  $L^2$  norm

$$\|A\|_F = \left( \sum_{ij} A_{ij}^2 \right)^{\frac{1}{2}}$$

$$A = \begin{bmatrix} 2 & -1 & 5 \\ 0 & 2 & 1 \\ 3 & 1 & 1 \end{bmatrix} \quad \|A\| = \sqrt{4 + 1 + 25 + \dots + 1} = \sqrt{46}$$

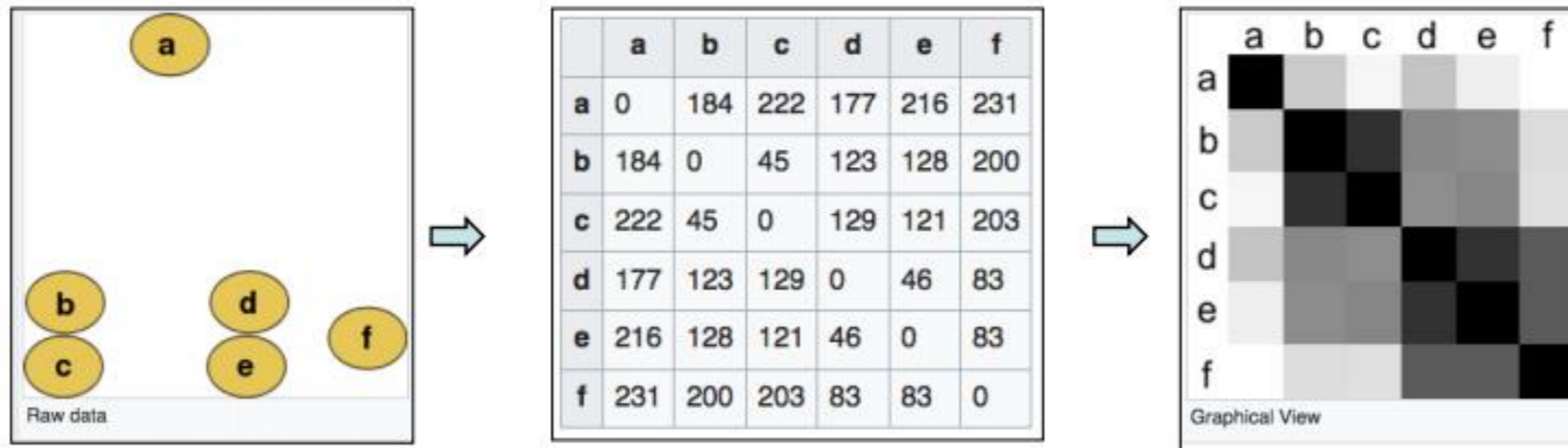


# Norms Can Serve as Building Blocks for Distance Measurements!

- Distance between two vectors  $(v, w)$ 
  - $\text{dist}(v, w) = \|v - w\|$ 
$$= \sqrt{(v_1 - w_1)^2 + \dots + (v_n - w_n)^2}$$

# Special kind of Matrix: Symmetric

- A symmetric matrix equals its transpose:  $A = A^T$ 
  - E.g., a distance matrix is symmetric with  $A_{ij} = A_{ji}$



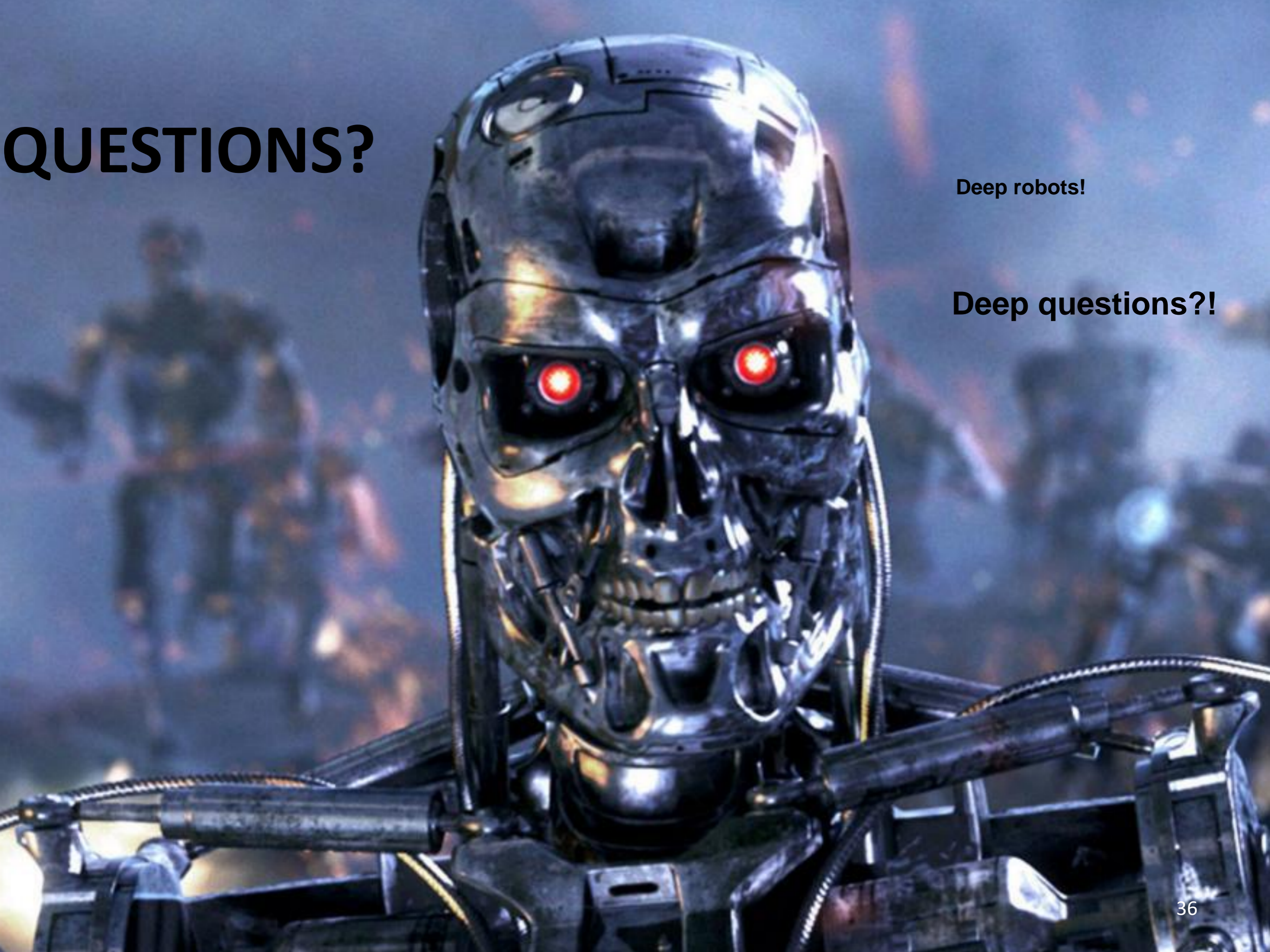
- E.g., covariance matrices are symmetric

$$\Sigma = \begin{pmatrix} 1 & .5 & .15 & .15 & 0 & 0 \\ .5 & 1 & .15 & .15 & 0 & 0 \\ .15 & .15 & 1 & .25 & 0 & 0 \\ .15 & .15 & .25 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & .10 \\ 0 & 0 & 0 & 0 & .10 & 1 \end{pmatrix}$$

# QUESTIONS?

Deep robots!

Deep questions?!



# Linear Transformation

- $Ax=b$

- where  $A \in \mathbb{R}^{n \times n}$  and  $b \in \mathbb{R}^n$

- More explicitly

$$\begin{array}{l} A_{11}x_1 + A_{12}x_2 + \dots + A_{1n}x_n = b_1 \\ A_{21}x_1 + A_{22}x_2 + \dots + A_{2n}x_n = b_2 \\ \dots \\ A_{n1}x_1 + A_{n2}x_2 + \dots + A_{nn}x_n = b_n \end{array}$$

$n$  equations in  
 $n$  unknowns

$$A = \begin{bmatrix} A_{1,1} & \dots & A_{1,n} \\ \vdots & \vdots & \vdots \\ A_{n,1} & \dots & A_{n,n} \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad b = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

$n \times n$                    $n \times 1$                    $n \times 1$

Can view  $A$  as a *linear transformation* of vector  $x$  to vector  $b$

- Sometimes we wish to solve for the unknowns  $x = \{x_1, \dots, x_n\}$  when  $A$  and  $b$  provide constraints

# Use of a Vector in Regression

- A design matrix
  - N samples, D features

	# hours studied	# hours playing games	# classes missed	
Student #1	10	3	0	Grade
Student #2	8	20	2	
Student #3	5	1	5	



87
75
63

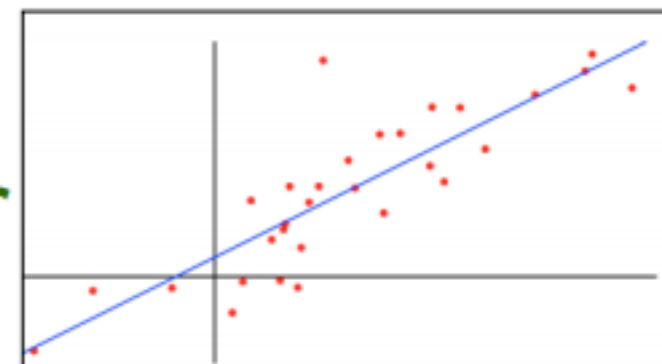
- Feature vector has three dimensions
- This is a regression problem

# Use of norm in Regression

- Linear Regression

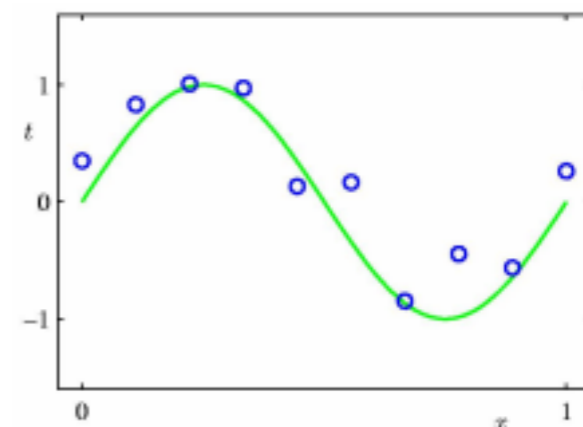
$\mathbf{x}$ : a vector,  $\mathbf{w}$ : weight vector

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_d x_d = \mathbf{w}^T \mathbf{x}$$



With nonlinear basis functions  $\phi_j$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x})$$



- Loss Function

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Second term is a weighted norm  
called a regularizer (to prevent overfitting)



# Matrix Inverse

- Inverse of square matrix  $A$  defined as  $A^{-1}A = I_n$
- We can now solve  $Ax = b$  as follows:

$$Ax = b$$

$$A^{-1}Ax = A^{-1}b$$

$$I_n x = A^{-1}b$$

$$x = A^{-1}b$$

- This depends on being able to find  $A^{-1}$
- If  $A^{-1}$  exists there are several methods for finding it

Numerically unstable, but useful for abstract analysis



# Closed-form solutions

- Two closed-form solutions

1. Matrix inversion  $\mathbf{x} = \mathbf{A}^{-1} \mathbf{b}$

2. Gaussian elimination

- If  $\mathbf{A}^{-1}$  exists, the same  $\mathbf{A}^{-1}$  can be used for any given  $\mathbf{b}$ 
  - But  $\mathbf{A}^{-1}$  cannot be represented with sufficient precision
  - It is not used in practice
- Gaussian elimination also has disadvantages
  - numerical instability (division by small no.)
  - $O(n^3)$  for  $n \times n$  matrix

# Invertibility

- Matrix can't be inverted if...
  - More rows than columns
  - More columns than rows
  - Redundant rows/columns (“linearly dependent”, “low rank”)

# Matrix Rank

- Matrix **rank** is defined as (a) maximum number of linearly independent *column* vectors in matrix or (b) maximum number of linearly independent *row* vectors in matrix (equivalent definitions)
- Linear independence/dependence

$$\mathbf{a} = \begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$$

$$\mathbf{b} = \begin{bmatrix} 4 & 5 & 6 \end{bmatrix}$$

$$\mathbf{c} = \begin{bmatrix} 5 & 7 & 9 \end{bmatrix}$$

$$\mathbf{d} = \begin{bmatrix} 2 & 4 & 6 \end{bmatrix}$$

$$\mathbf{e} = \begin{bmatrix} 0 & 1 & 0 \end{bmatrix}$$

$$\mathbf{f} = \begin{bmatrix} 0 & 0 & 1 \end{bmatrix}$$