

# OUTLINE

---

## I. Introduction

- Artificial Intelligence
- Machine Learning
- Neural Network and Deep Learning
- Applications

## II. Background

## III. Fitting a Model

## IV. Supervised Learning

## V. Unsupervised Learning

## VI. Fantastic DNNs: How to choose them, how to train them

## VII. Machine Learning in Robot Audition

# OUTLINE

---

I. Introduction

**II. Background**

III. Fitting a Model

IV. Supervised Learning

V. Unsupervised Learning

VI. Fantastic DNNs: How to choose them, how to train them

VII. Machine Learning in Robot Audition

# OUTLINE

---

I. Introduction

## II. Background

- Multivalued Multivariate Functions
- Tensors
- Differential Calculus
- Exercises

III. Fitting a Model

IV. Supervised Learning

V. Unsupervised Learning

VI. Fantastic DNNs: How to choose them, how to train them

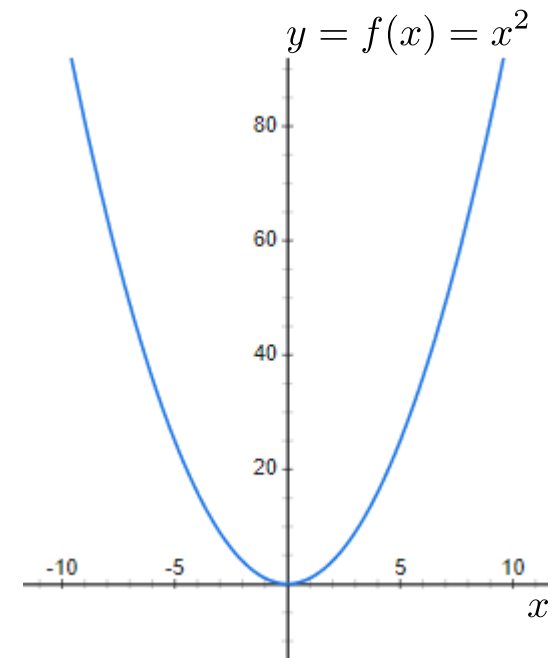
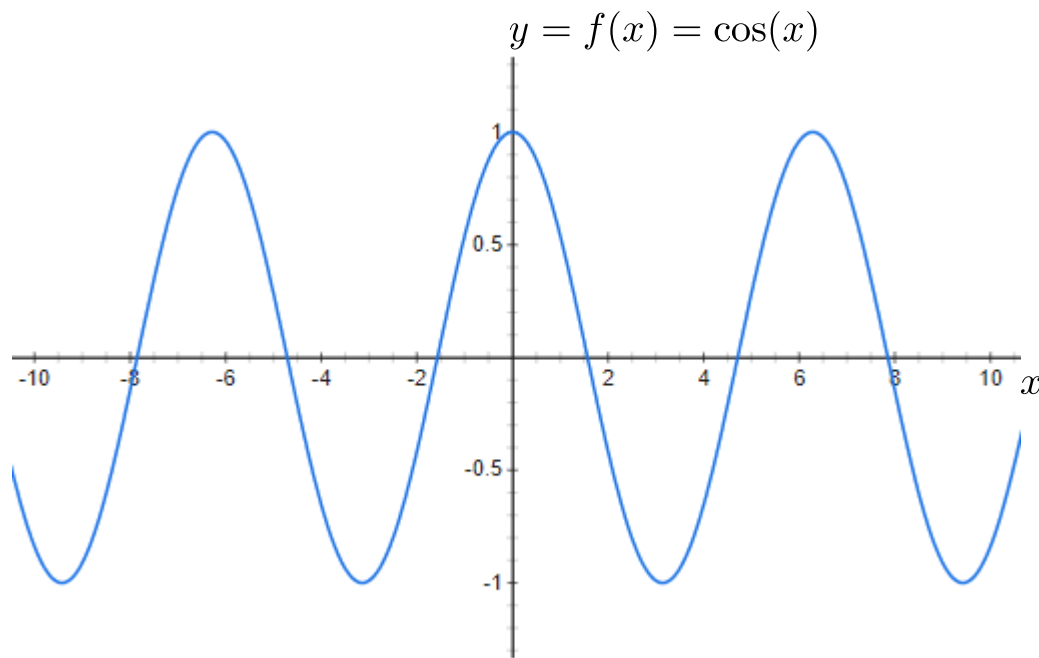
VII. Machine Learning in Robot Audition

## Simple functions

- In high-school calculus, we study functions from reals to reals, denoted  $f : \mathbb{R} \rightarrow \mathbb{R}$ .

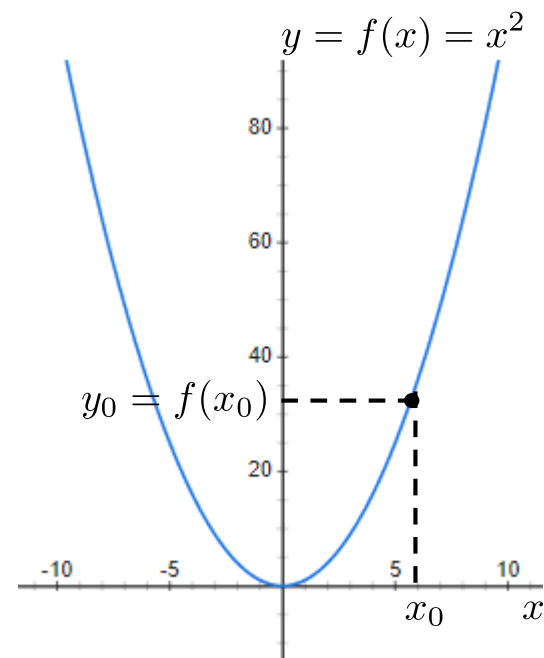
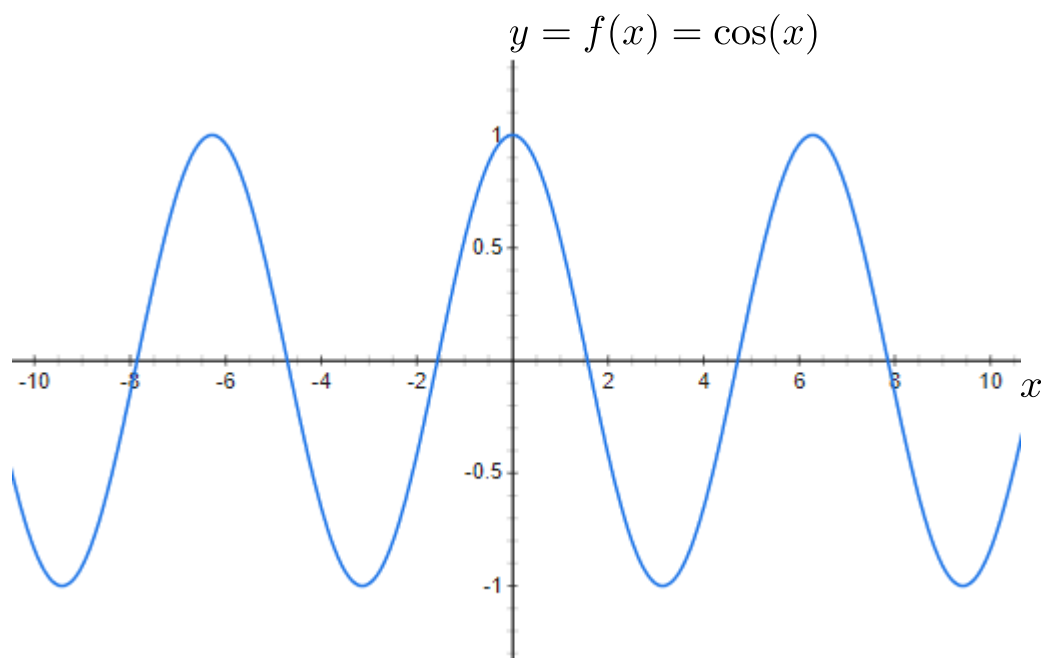
## Simple functions

- In high-school calculus, we study functions from reals to reals, denoted  $f : \mathbb{R} \rightarrow \mathbb{R}$ .
- Here are some **graphs** of functions:



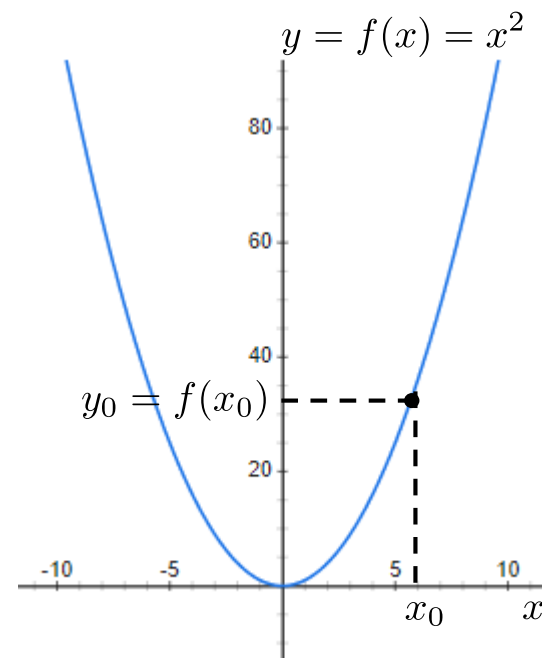
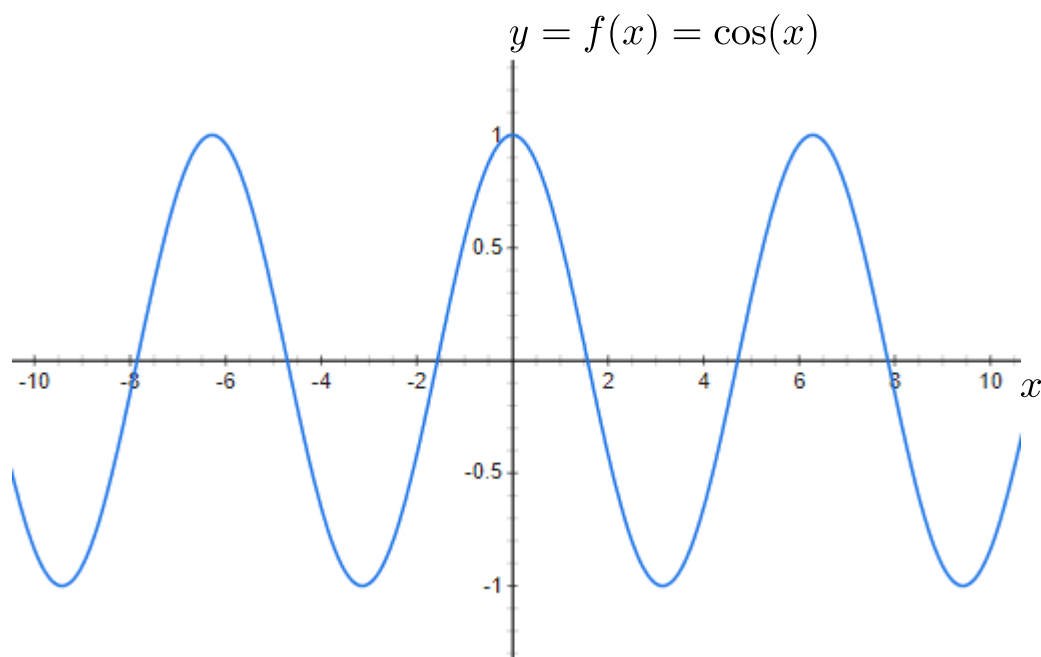
## Simple functions

- In high-school calculus, we study functions from reals to reals, denoted  $f : \mathbb{R} \rightarrow \mathbb{R}$ .
- Here are some **graphs** of functions:



## Simple functions

- In high-school calculus, we study functions from reals to reals, denoted  $f : \mathbb{R} \rightarrow \mathbb{R}$ .
- Here are some **graphs** of functions:



- **Important subtlety:** Here,  $x$  and  $y$  are **variables** that **depend** on each other,  $x_0$  and  $y_0$  are **constants**, and  $f$  is a **function** that can be **applied** to a variable or a constant.

## Multivalued functions

- This straightforwardly generalizes to **multivalued** functions

$$f : \mathbb{R} \rightarrow \mathbb{R}^N :$$

$$\mathbf{y} = f(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_N(x) \end{bmatrix} \in \mathbb{R}^N$$



## Multivalued functions

- This straightforwardly generalizes to **multivalued** functions

$$f : \mathbb{R} \rightarrow \mathbb{R}^N :$$

$$\mathbf{y} = f(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_N(x) \end{bmatrix} \in \mathbb{R}^N$$

**Example:**

$$f : \begin{cases} [0, 2\pi] & \rightarrow \mathbb{R}^2 \\ \theta & \mapsto f(\theta) = \begin{bmatrix} \cos(\theta) \\ \sin(\theta) \end{bmatrix} \end{cases}$$

## Multivalued functions

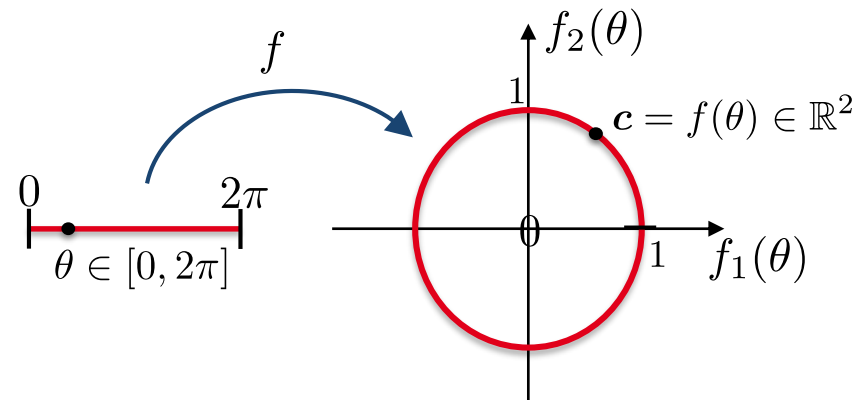
- This straightforwardly generalizes to **multivalued** functions

$$f : \mathbb{R} \rightarrow \mathbb{R}^N :$$

$$\mathbf{y} = f(x) = \begin{bmatrix} f_1(x) \\ f_2(x) \\ \vdots \\ f_N(x) \end{bmatrix} \in \mathbb{R}^N$$

**Example:**

$$f : \begin{cases} [0, 2\pi] & \rightarrow \mathbb{R}^2 \\ \theta & \mapsto f(\theta) = \begin{bmatrix} \cos(\theta) \\ \sin(\theta) \end{bmatrix} \end{cases}$$

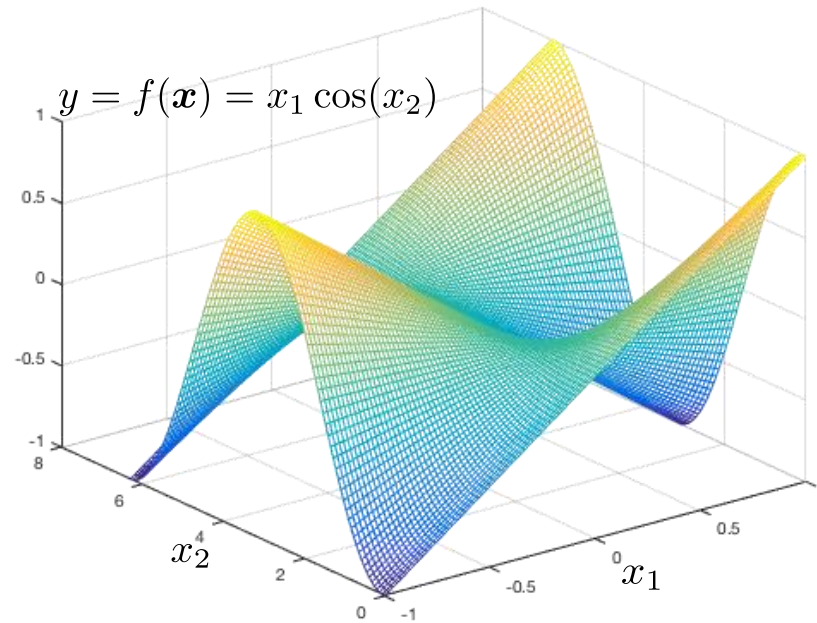


## Multivariate functions

- Another generalization are **real-valued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}$

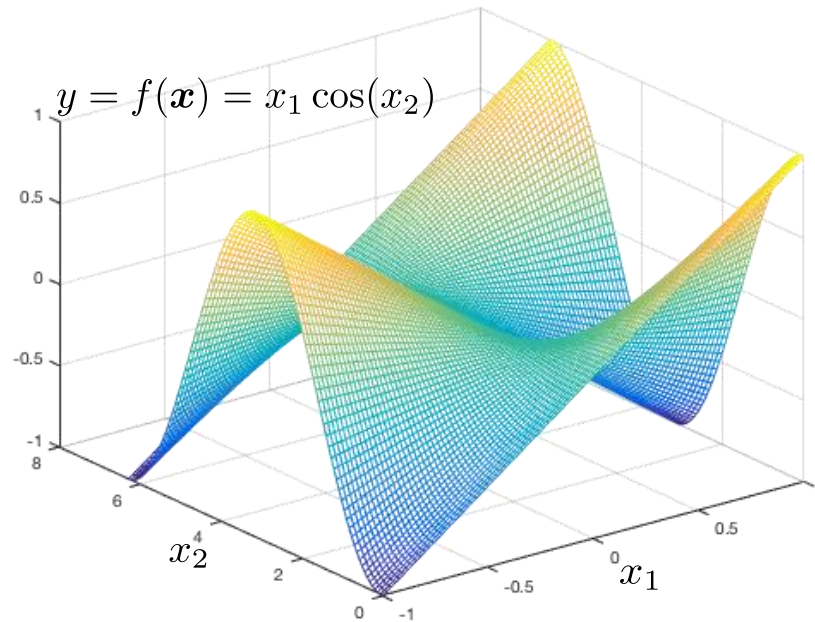
## Multivariate functions

- Another generalization are **real-valued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}$



## Multivariate functions

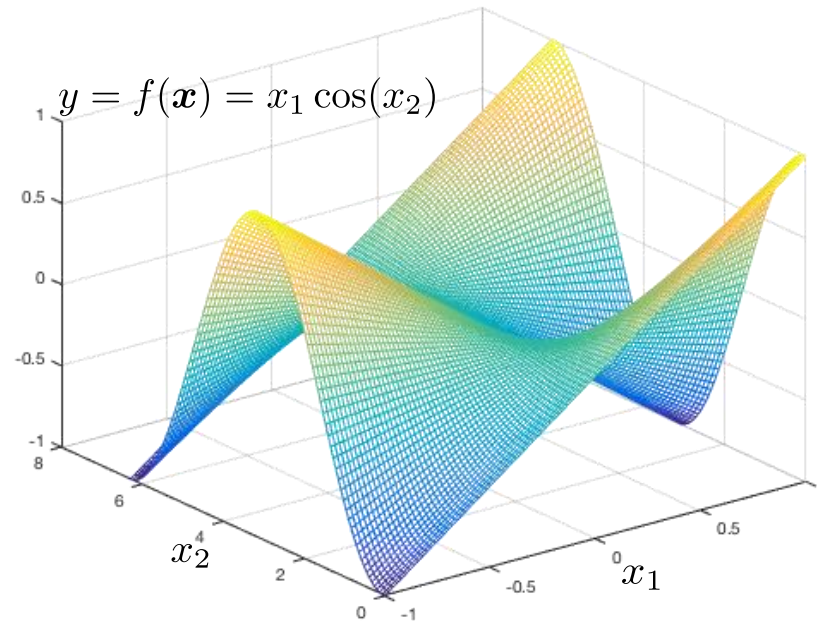
- Another generalization are **real-valued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}$



**Important examples:**

## Multivariate functions

- Another generalization are **real-valued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}$

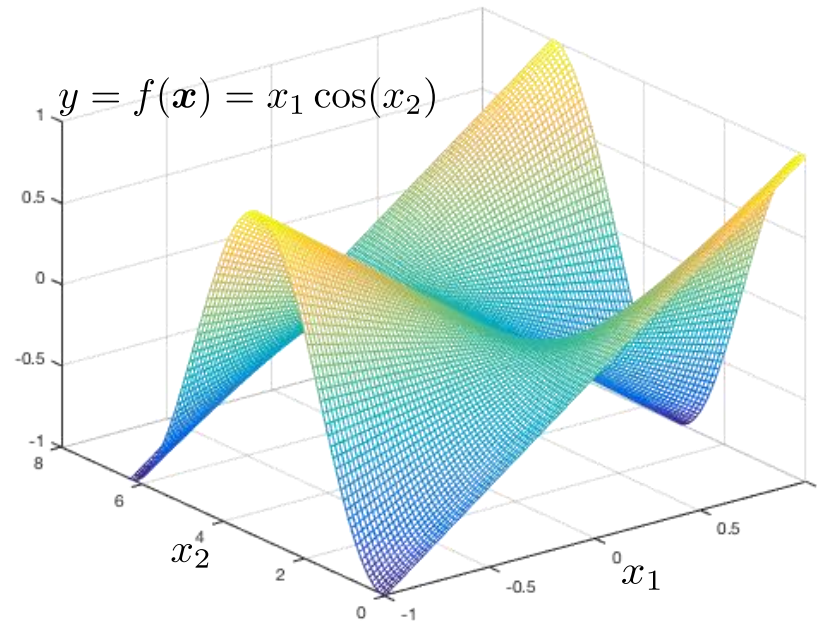


### Important examples:

- **Linear forms:**  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} = \langle \mathbf{w}, \mathbf{x} \rangle = \sum_{d=1}^D w_d x_d$

## Multivariate functions

- Another generalization are **real-valued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}$

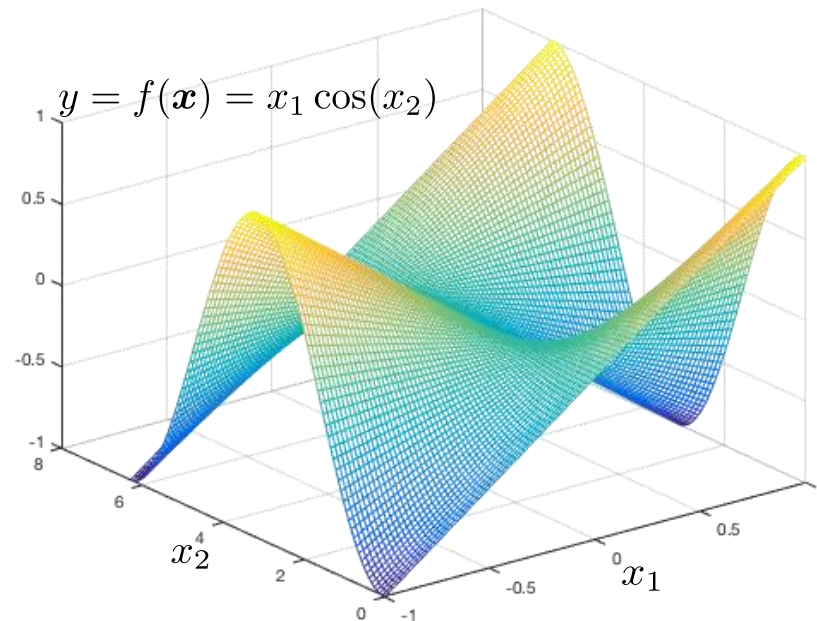


### Important examples:

- **Linear forms:**  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} = \langle \mathbf{w}, \mathbf{x} \rangle = \sum_{d=1}^D w_d x_d \rightarrow$  represented by a **row vector** :  $\mathbf{w}^\top \in \mathbb{R}^{1 \times D}$

## Multivariate functions

- Another generalization are **real-valued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}$



### Important examples:

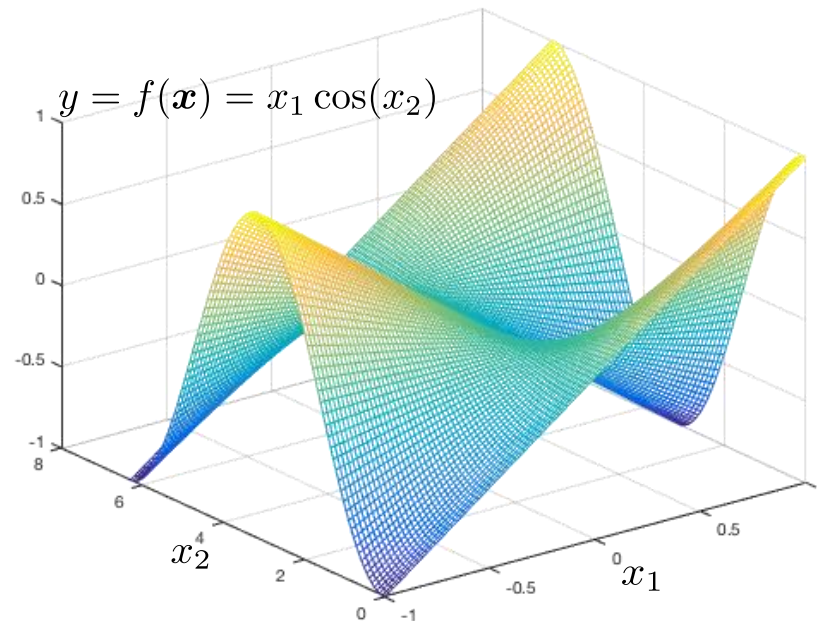
- **Linear forms:**  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} = \langle \mathbf{w}, \mathbf{x} \rangle = \sum_{d=1}^D w_d x_d \rightarrow$  represented by a **row vector** :  $\mathbf{w}^\top \in \mathbb{R}^{1 \times D}$

- **Affine forms:**  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$



## Multivariate functions

- Another generalization are **real-valued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}$



### Important examples:

- **Linear forms:**  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} = \langle \mathbf{w}, \mathbf{x} \rangle = \sum_{d=1}^D w_d x_d \rightarrow$  represented by a **row vector** :  $\mathbf{w}^\top \in \mathbb{R}^{1 \times D}$
- **Affine forms:**  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$
- The **Euclidean norm:**  $f(\mathbf{x}) = \|\mathbf{x}\|_2^2 = \mathbf{x}^\top \mathbf{x} = \sum_{d=1}^D (x_d)^2$

## Multivalued Multivariate functions

- Finally, combining the two, we get **multivalued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}^N$ .

## Multivalued Multivariate functions

- Finally, combining the two, we get **multivalued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}^N$ .

### Important examples:

- **Linear maps:**  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$

## Multivalued Multivariate functions

- Finally, combining the two, we get **multivalued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}^N$ .

### Important examples:

- **Linear maps**:  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$

They can be *represented* by **matrices** :  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{A} \in \mathbb{R}^{D \times N}$

## Multivalued Multivariate functions

- Finally, combining the two, we get **multivalued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}^N$ .

### Important examples:

- **Linear maps**:  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$

They can be *represented* by **matrices** :  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{A} \in \mathbb{R}^{D \times N}$

$$f(\mathbf{x}) = \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,D} \\ a_{2,1} & a_{2,2} & \dots & a_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N,1} & a_{L,2} & \dots & a_{N,D} \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix}$$

## Multivalued Multivariate functions

- Finally, combining the two, we get **multivalued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}^N$ .

### Important examples:

- Linear maps:**  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$

They can be *represented* by **matrices** :  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{A} \in \mathbb{R}^{D \times N}$

$$f(\mathbf{x}) = \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,D} \\ a_{2,1} & a_{2,2} & \dots & a_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N,1} & a_{L,2} & \dots & a_{N,D} \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{a}}_1^\top \\ \hat{\mathbf{a}}_2^\top \\ \vdots \\ \hat{\mathbf{a}}_N^\top \end{bmatrix} \mathbf{x}$$

## Multivalued Multivariate functions

- Finally, combining the two, we get **multivalued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}^N$ .

### Important examples:

- Linear maps:**  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$

They can be *represented* by **matrices** :  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{A} \in \mathbb{R}^{D \times N}$

$$f(\mathbf{x}) = \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,D} \\ a_{2,1} & a_{2,2} & \dots & a_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N,1} & a_{L,2} & \dots & a_{N,D} \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{a}}_1^\top \\ \hat{\mathbf{a}}_2^\top \\ \vdots \\ \hat{\mathbf{a}}_N^\top \end{bmatrix} \mathbf{x} = \begin{bmatrix} \langle \hat{\mathbf{a}}_1, \mathbf{x} \rangle \\ \langle \hat{\mathbf{a}}_2, \mathbf{x} \rangle \\ \vdots \\ \langle \hat{\mathbf{a}}_N, \mathbf{x} \rangle \end{bmatrix}$$

## Multivalued Multivariate functions

- Finally, combining the two, we get **multivalued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}^N$ .

### Important examples:

- Linear maps:**  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$

They can be *represented* by **matrices** :  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{A} \in \mathbb{R}^{D \times N}$

$$f(\mathbf{x}) = \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,D} \\ a_{2,1} & a_{2,2} & \dots & a_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N,1} & a_{L,2} & \dots & a_{N,D} \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{a}}_1^\top \\ \hat{\mathbf{a}}_2^\top \\ \vdots \\ \hat{\mathbf{a}}_N^\top \end{bmatrix} \mathbf{x} = \begin{bmatrix} \langle \hat{\mathbf{a}}_1, \mathbf{x} \rangle \\ \langle \hat{\mathbf{a}}_2, \mathbf{x} \rangle \\ \vdots \\ \langle \hat{\mathbf{a}}_N, \mathbf{x} \rangle \end{bmatrix} = \begin{bmatrix} \sum_{d=1}^D a_{1,d}x_d \\ \sum_{d=1}^D a_{2,d}x_d \\ \vdots \\ \sum_{d=1}^D a_{N,d}x_d \end{bmatrix}$$



## Multivalued Multivariate functions

- Finally, combining the two, we get **multivalued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}^N$ .

### Important examples:

- Linear maps:**  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$

They can be *represented* by **matrices** :  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{A} \in \mathbb{R}^{D \times N}$

$$f(\mathbf{x}) = \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,D} \\ a_{2,1} & a_{2,2} & \dots & a_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N,1} & a_{N,2} & \dots & a_{N,D} \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{a}}_1^\top \\ \hat{\mathbf{a}}_2^\top \\ \vdots \\ \hat{\mathbf{a}}_N^\top \end{bmatrix} \mathbf{x} = \begin{bmatrix} \langle \hat{\mathbf{a}}_1, \mathbf{x} \rangle \\ \langle \hat{\mathbf{a}}_2, \mathbf{x} \rangle \\ \vdots \\ \langle \hat{\mathbf{a}}_N, \mathbf{x} \rangle \end{bmatrix} = \begin{bmatrix} \sum_{d=1}^D a_{1,d}x_d \\ \sum_{d=1}^D a_{2,d}x_d \\ \vdots \\ \sum_{d=1}^D a_{N,d}x_d \end{bmatrix}$$

## Multivalued Multivariate functions

- Finally, combining the two, we get **multivalued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}^N$ .

### Important examples:

- Linear maps:**  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$

They can be *represented* by **matrices** :  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{A} \in \mathbb{R}^{D \times N}$

$$f(\mathbf{x}) = \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,D} \\ a_{2,1} & a_{2,2} & \dots & a_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N,1} & a_{L,2} & \dots & a_{N,D} \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{a}}_1^\top \\ \hat{\mathbf{a}}_2^\top \\ \vdots \\ \hat{\mathbf{a}}_N^\top \end{bmatrix} \mathbf{x} = \begin{bmatrix} \langle \hat{\mathbf{a}}_1, \mathbf{x} \rangle \\ \langle \hat{\mathbf{a}}, \mathbf{x} \rangle \\ \vdots \\ \langle \hat{\mathbf{a}}_N, \mathbf{x} \rangle \end{bmatrix} = \begin{bmatrix} \sum_{d=1}^D a_{1,d}x_d \\ \sum_{d=1}^D a_{2,d}x_d \\ \vdots \\ \sum_{d=1}^D a_{N,d}x_d \end{bmatrix}$$

## Multivalued Multivariate functions

- Finally, combining the two, we get **multivalued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}^N$ .

### Important examples:

- Linear maps:**  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$

They can be *represented* by **matrices** :  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{A} \in \mathbb{R}^{D \times N}$

$$f(\mathbf{x}) = \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,D} \\ a_{2,1} & a_{2,2} & \dots & a_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N,1} & a_{N,2} & \dots & a_{N,D} \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{a}}_1^\top \\ \hat{\mathbf{a}}_2^\top \\ \vdots \\ \hat{\mathbf{a}}_N^\top \end{bmatrix} \mathbf{x} = \begin{bmatrix} \langle \hat{\mathbf{a}}_1, \mathbf{x} \rangle \\ \langle \hat{\mathbf{a}}_2, \mathbf{x} \rangle \\ \vdots \\ \langle \hat{\mathbf{a}}_N, \mathbf{x} \rangle \end{bmatrix} = \begin{bmatrix} \sum_{d=1}^D a_{1,d}x_d \\ \sum_{d=1}^D a_{2,d}x_d \\ \vdots \\ \sum_{d=1}^D a_{N,d}x_d \end{bmatrix}$$

## Multivalued Multivariate functions

- Finally, combining the two, we get **multivalued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}^N$ .

### Important examples:

- Linear maps:**  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$

They can be *represented* by **matrices** :  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{A} \in \mathbb{R}^{D \times N}$

$$\begin{aligned}
 f(\mathbf{x}) &= \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,D} \\ a_{2,1} & a_{2,2} & \dots & a_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N,1} & a_{L,2} & \dots & a_{N,D} \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{a}}_1^\top \\ \hat{\mathbf{a}}_2^\top \\ \vdots \\ \hat{\mathbf{a}}_N^\top \end{bmatrix} \mathbf{x} = \begin{bmatrix} \langle \hat{\mathbf{a}}_1, \mathbf{x} \rangle \\ \langle \hat{\mathbf{a}}_2, \mathbf{x} \rangle \\ \vdots \\ \langle \hat{\mathbf{a}}_N, \mathbf{x} \rangle \end{bmatrix} = \begin{bmatrix} \sum_{d=1}^D a_{1,d}x_d \\ \sum_{d=1}^D a_{2,d}x_d \\ \vdots \\ \sum_{d=1}^D a_{N,d}x_d \end{bmatrix} \\
 &= [\mathbf{a}_1, \dots, \mathbf{a}_D] \times \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \sum_{d=1}^D \mathbf{a}_d x_d
 \end{aligned}$$

## Multivalued Multivariate functions

- Finally, combining the two, we get **multivalued multivariate functions**, i.e.,  $f : \mathbb{R}^D \rightarrow \mathbb{R}^N$ .

### Important examples:

- Linear maps:**  $f(\mathbf{x} + \mathbf{y}) = f(\mathbf{x}) + f(\mathbf{y})$

They can be *represented* by **matrices** :  $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ ,  $\mathbf{A} \in \mathbb{R}^{D \times N}$

$$\begin{aligned}
 f(\mathbf{x}) &= \begin{bmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,D} \\ a_{2,1} & a_{2,2} & \dots & a_{2,D} \\ \vdots & \vdots & \ddots & \vdots \\ a_{N,1} & a_{L,2} & \dots & a_{N,D} \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \begin{bmatrix} \hat{\mathbf{a}}_1^\top \\ \hat{\mathbf{a}}_2^\top \\ \vdots \\ \hat{\mathbf{a}}_N^\top \end{bmatrix} \mathbf{x} = \begin{bmatrix} \langle \hat{\mathbf{a}}_1, \mathbf{x} \rangle \\ \langle \hat{\mathbf{a}}_2, \mathbf{x} \rangle \\ \vdots \\ \langle \hat{\mathbf{a}}_N, \mathbf{x} \rangle \end{bmatrix} = \begin{bmatrix} \sum_{d=1}^D a_{1,d}x_d \\ \sum_{d=1}^D a_{2,d}x_d \\ \vdots \\ \sum_{d=1}^D a_{N,d}x_d \end{bmatrix} \\
 &= [\mathbf{a}_1, \dots, \mathbf{a}_D] \times \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} = \sum_{d=1}^D \mathbf{a}_d x_d
 \end{aligned}$$

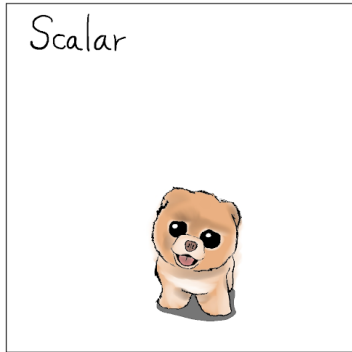
- Affine maps:**  $f(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$ ,  $\mathbf{A} \in \mathbb{R}^{N \times D}$ ,  $\mathbf{b} = \mathbb{R}^N$

# Tensors

- Matrices and vectors can be generalized to **Tensors**

# Tensors

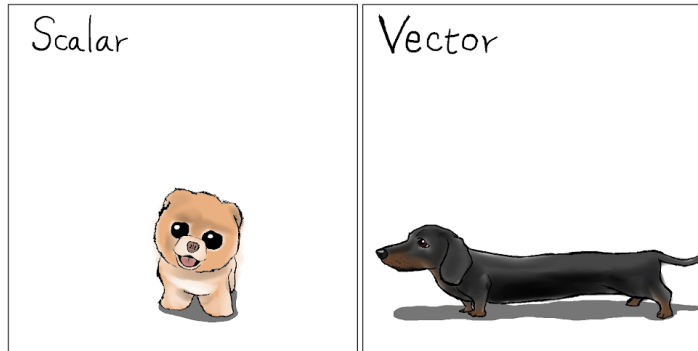
- Matrices and vectors can be generalized to **Tensors**



$a \in \mathbb{R}$   
(0-way tensor)

# Tensors

- Matrices and vectors can be generalized to **Tensors**



$a \in \mathbb{R}$   
(0-way tensor)

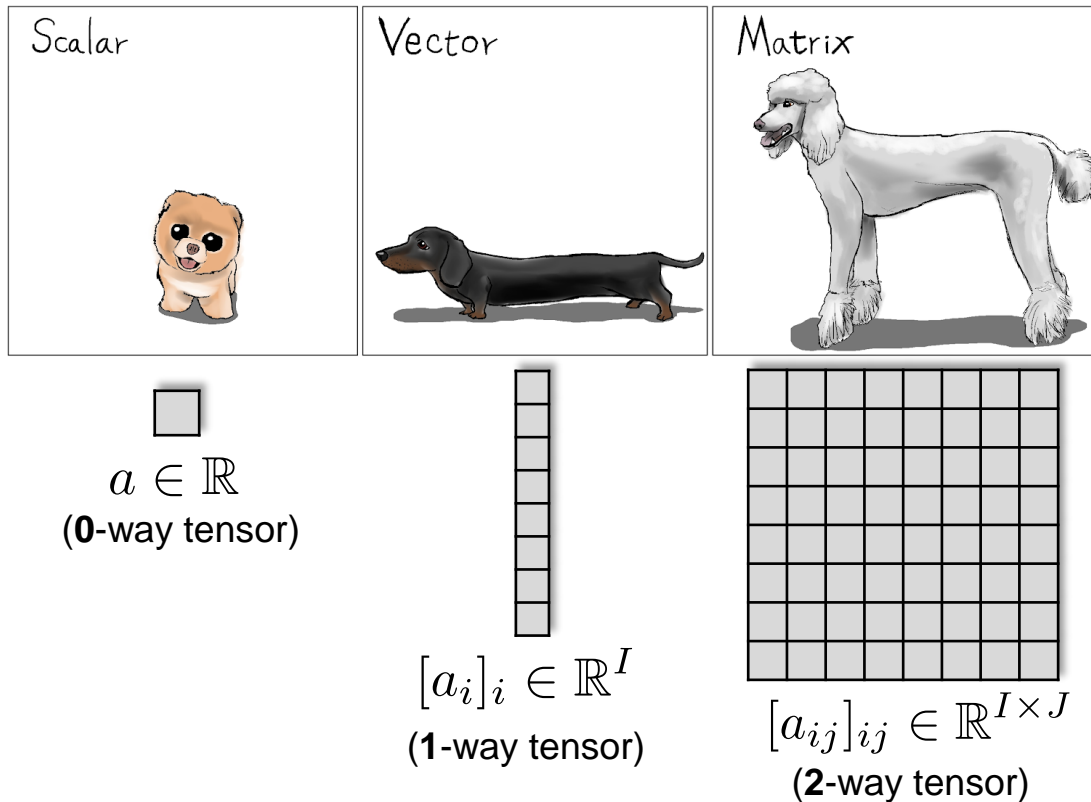


$[a_i]_i \in \mathbb{R}^I$   
(1-way tensor)





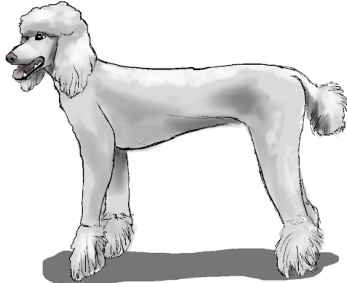



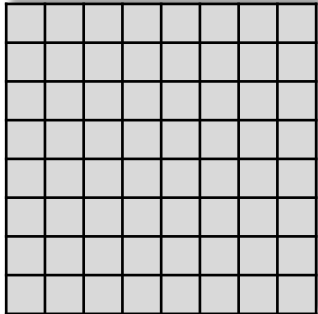
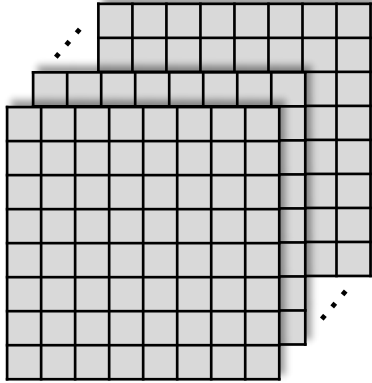
# Tensors

- Matrices and vectors can be generalized to **Tensors**



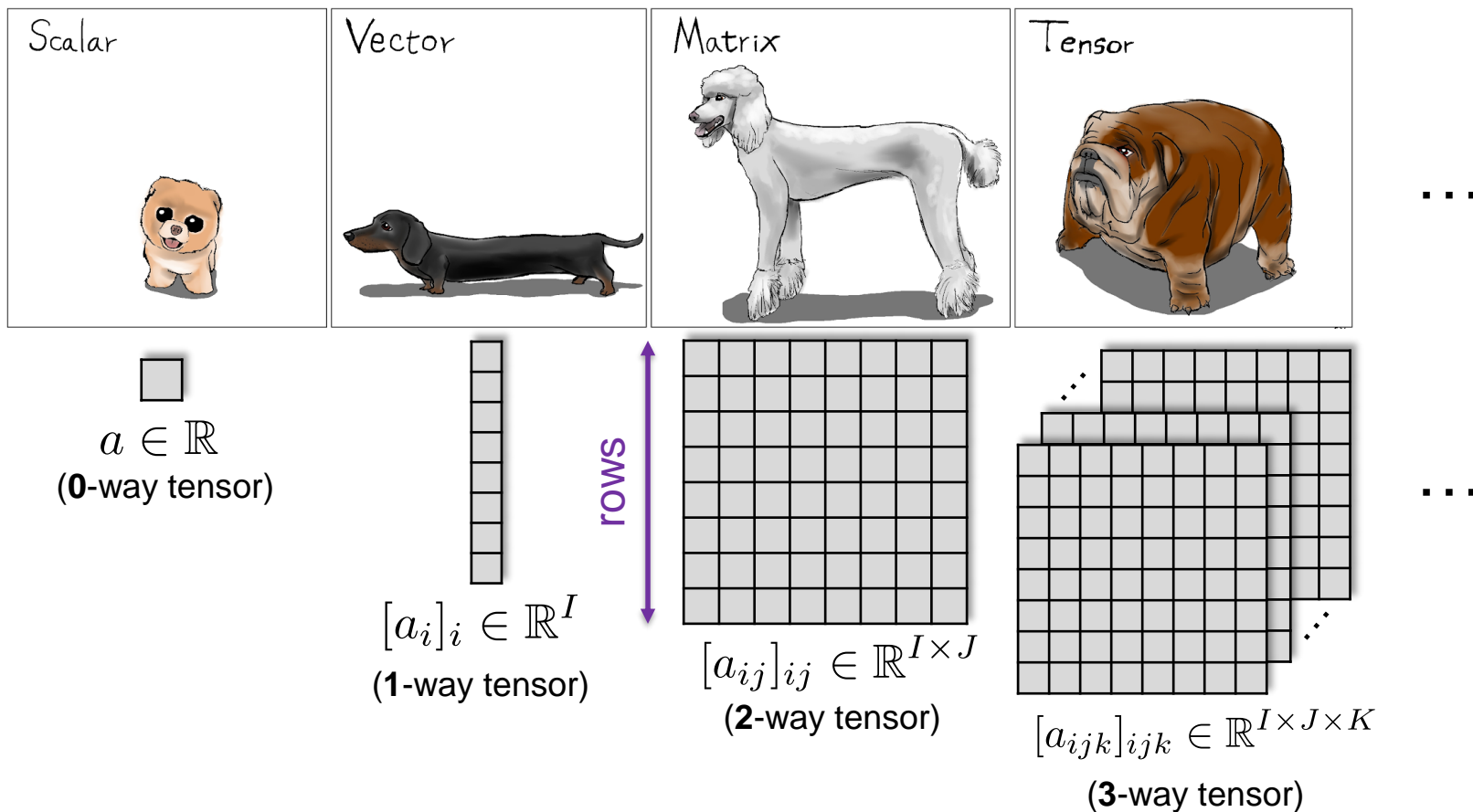
# Tensors

- Matrices and vectors can be generalized to **Tensors**

| Scalar  | Vector  | Matrix   | Tensor   | ... |
|---|---|--|--|-----|
|  |  |  |   | ... |
|  |  |  |  | ... |
| $a \in \mathbb{R}$<br>(0-way tensor)  | $[a_i]_i \in \mathbb{R}^I$<br>(1-way tensor)                                      | $[a_{ij}]_{ij} \in \mathbb{R}^{I \times J}$<br>(2-way tensor)                      | $[a_{ijk}]_{ijk} \in \mathbb{R}^{I \times J \times K}$<br>(3-way tensor)             |     |

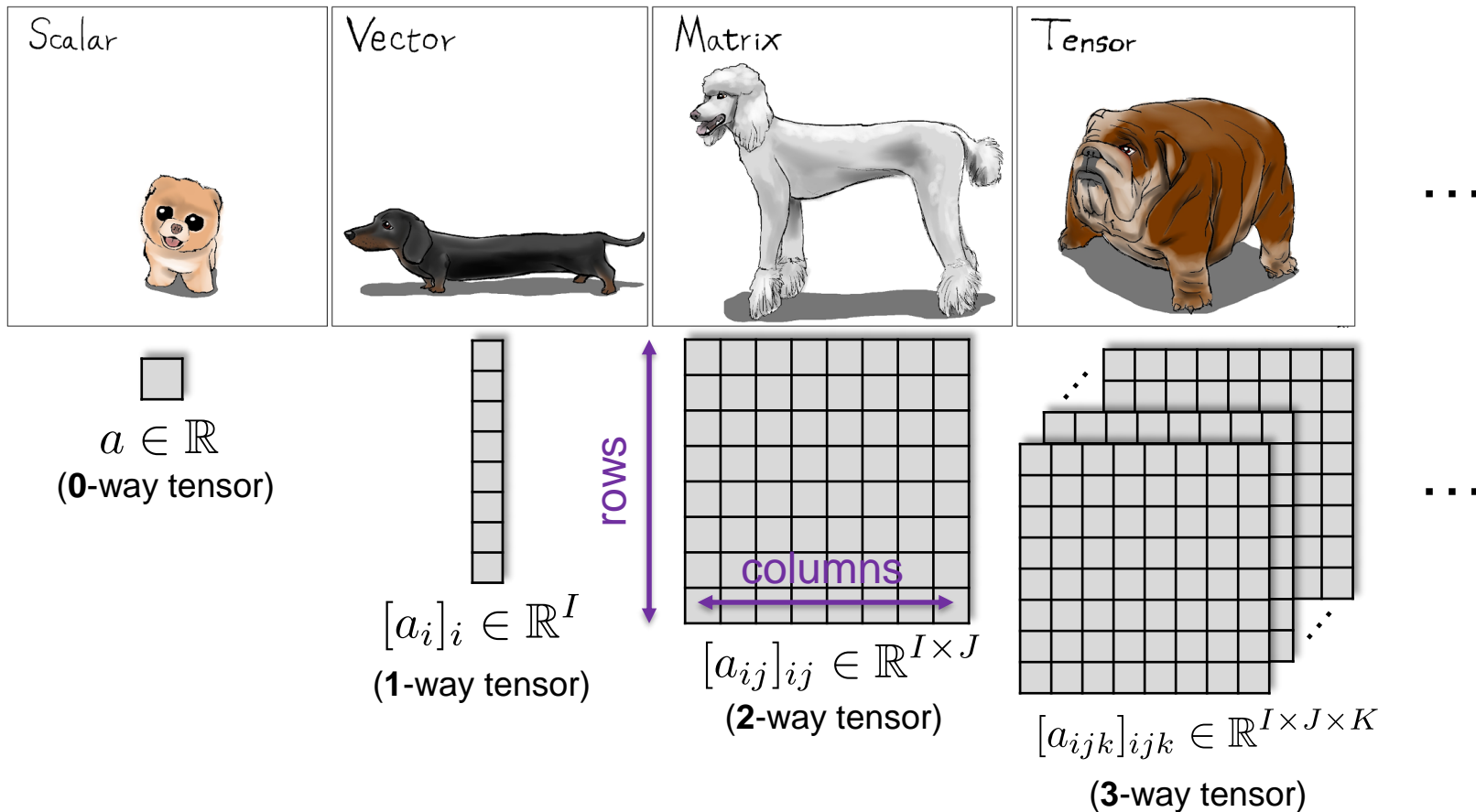
### Tensors

- Matrices and vectors can be generalized to **Tensors**



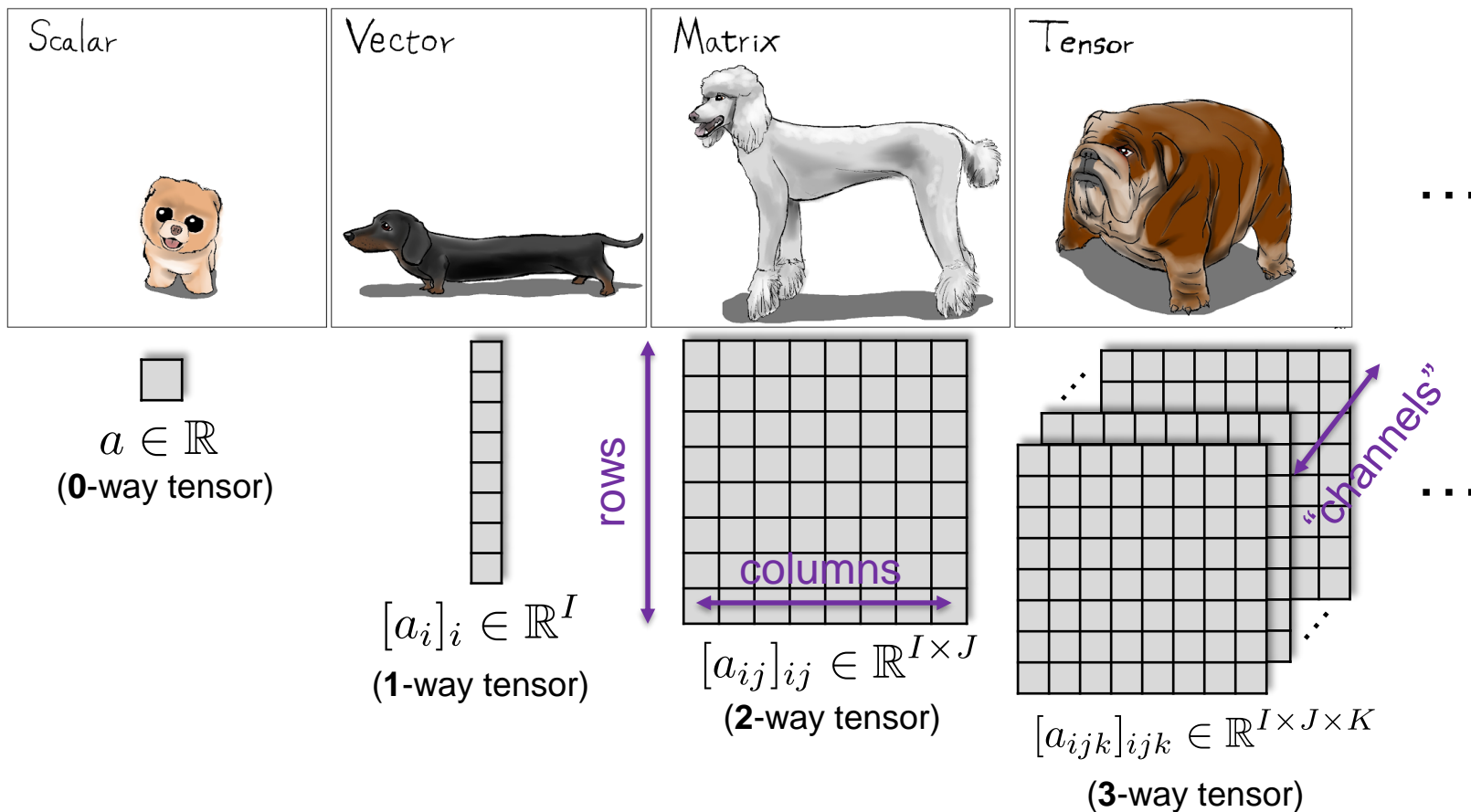
# Tensors

- Matrices and vectors can be generalized to **Tensors**



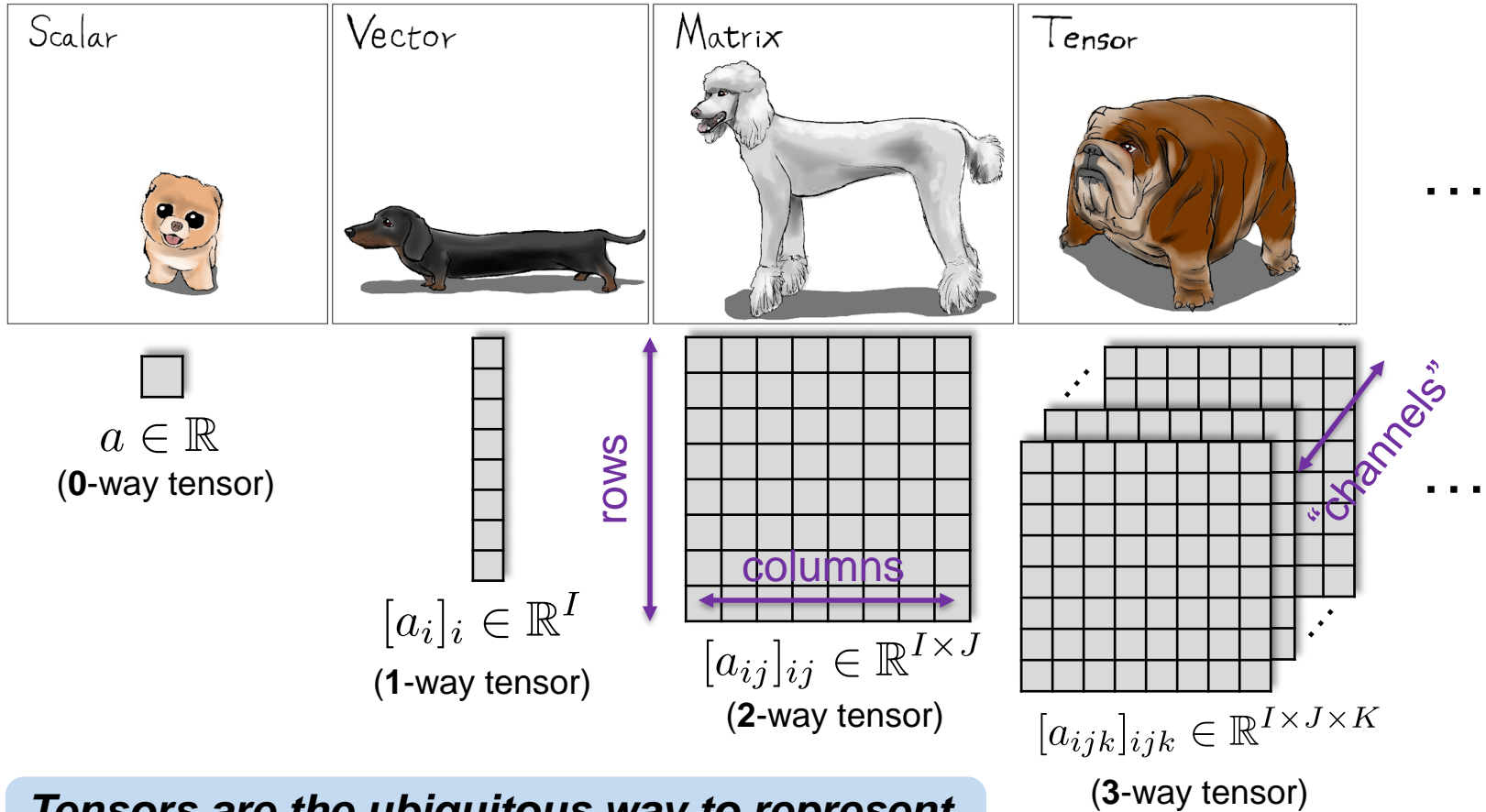
# Tensors

- Matrices and vectors can be generalized to **Tensors**



# Tensors

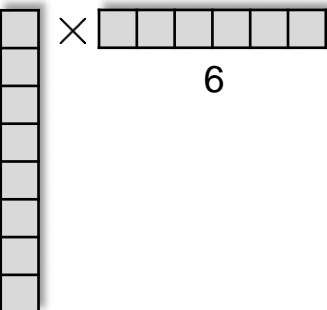
- Matrices and vectors can be generalized to **Tensors**

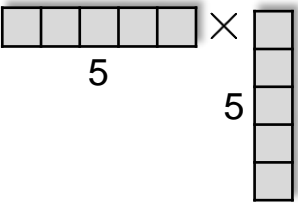


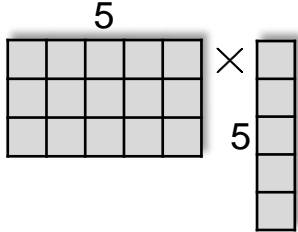
***Tensors are the ubiquitous way to represent data in modern deep learning frameworks***

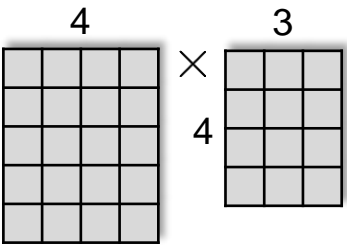
### Tensors

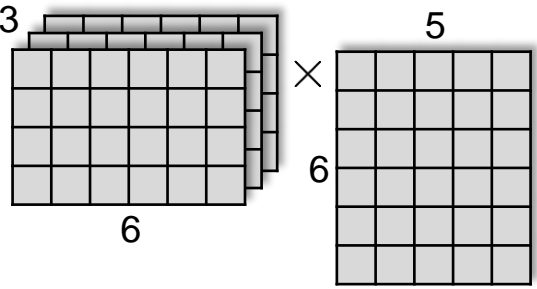
- Some vector/matrix/tensor operations:

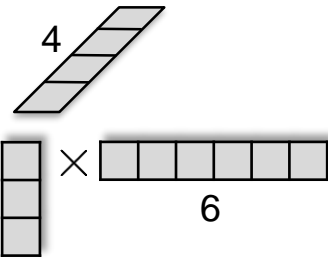
a)  = ?

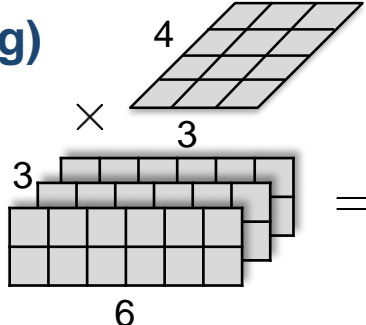
b)  = ?

c)  = ?

d)  = ?

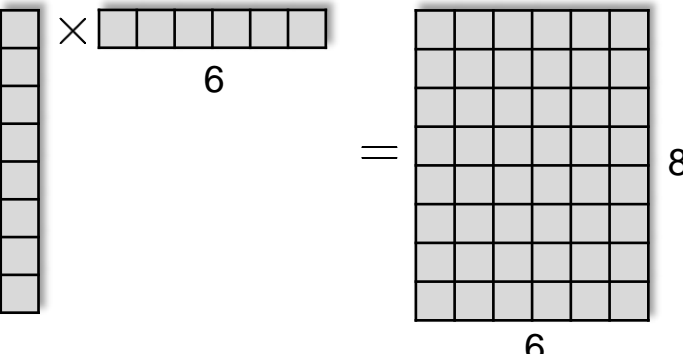
f)  = ?

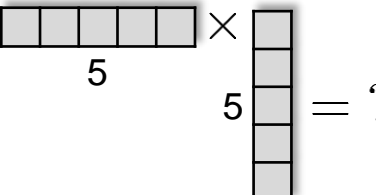
e)  = ?

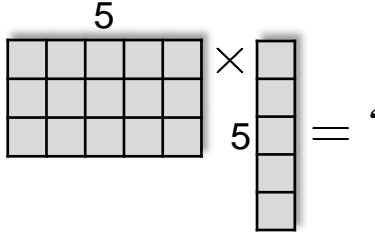
g)  = ?

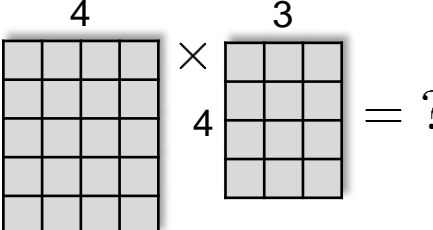
### Tensors

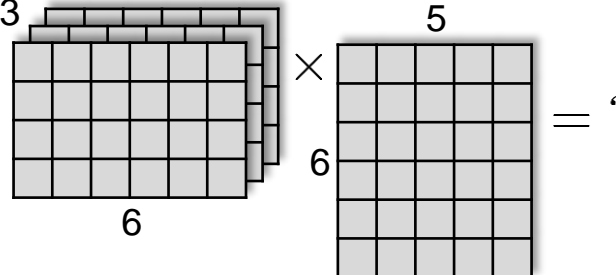
- Some vector/matrix/tensor operations:

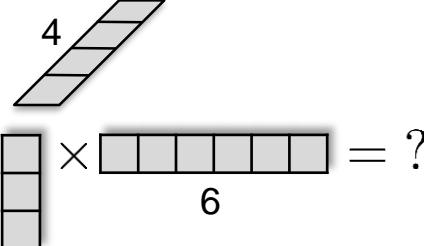
**a)**   $8 \times 6 = 8 \times 6$

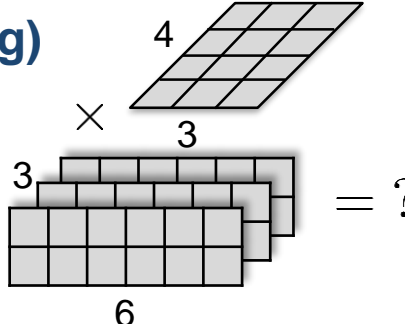
**b)**   $5 \times 5 = ?$

**c)**   $3 \times 5 \times 5 = ?$

**d)**   $5 \times 4 \times 4 \times 3 = ?$

**f)**   $3 \times 4 \times 6 \times 6 \times 5 = ?$

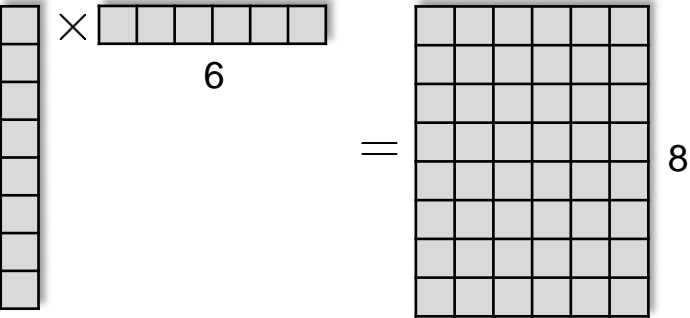
**e)**   $3 \times 4 \times 3 \times 6 = ?$

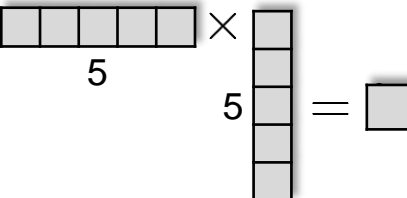
**g)**   $2 \times 3 \times 6 \times 4 \times 3 = ?$

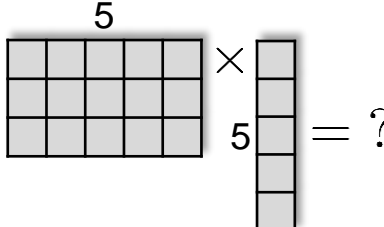


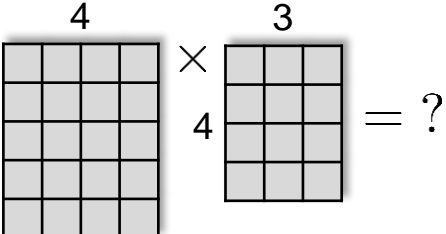
### Tensors

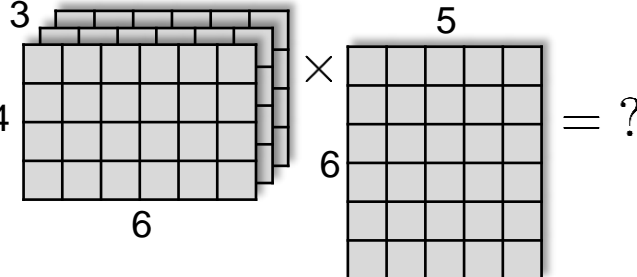
- Some vector/matrix/tensor operations:

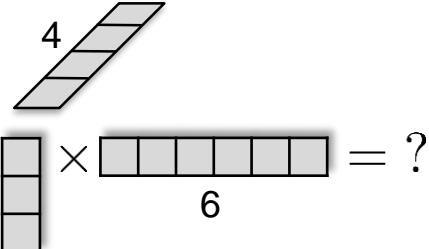
**a)**   $8 \times 6 = 8 \times 6$

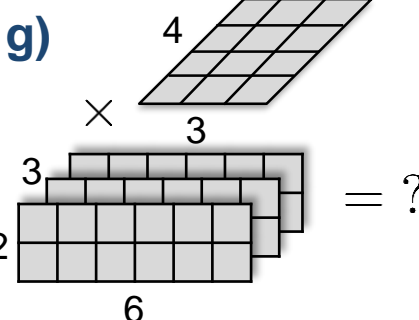
**b)**   $5 \times 5 = 1$

**c)**   $3 \times 5 \times 5 = ?$

**d)**   $4 \times 5 \times 3 \times 4 = ?$

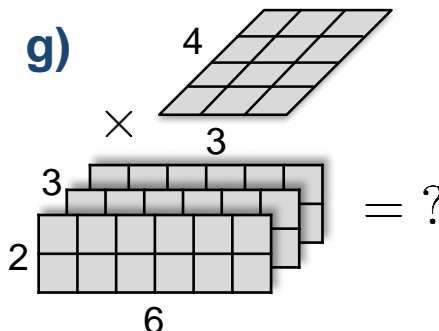
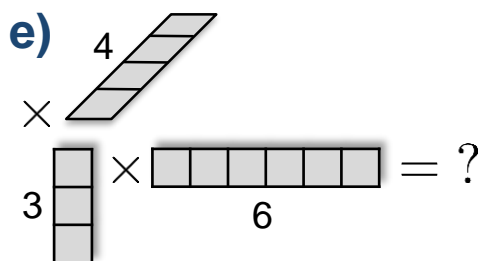
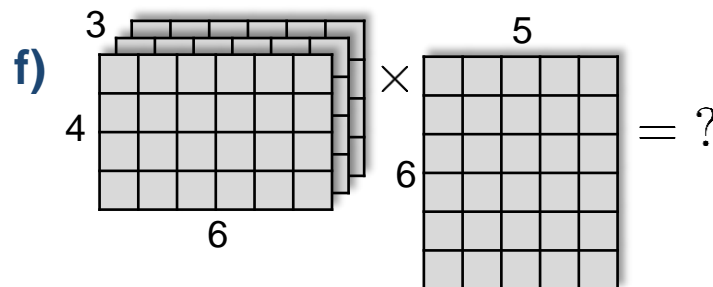
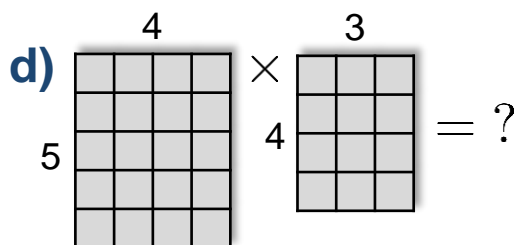
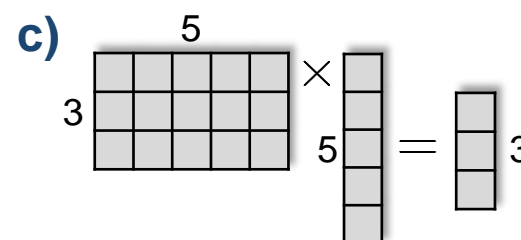
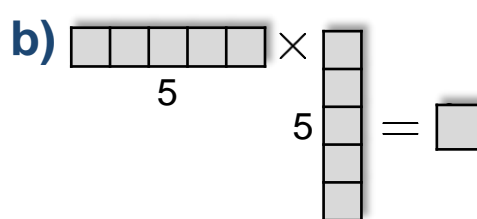
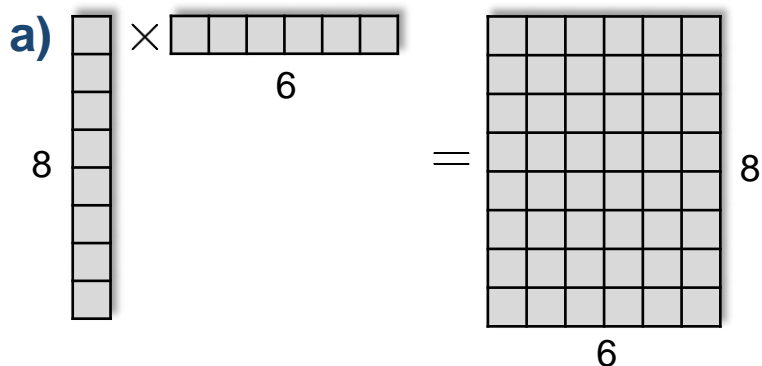
**f)**   $3 \times 4 \times 6 \times 5 \times 6 = ?$

**e)**   $3 \times 4 \times 3 \times 6 = ?$

**g)**   $2 \times 3 \times 6 \times 4 \times 3 = ?$

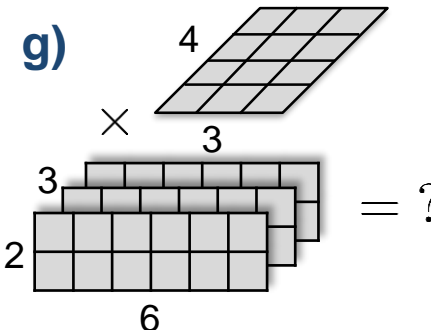
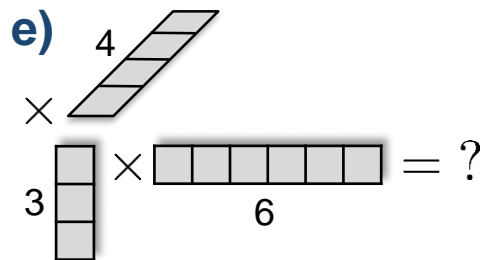
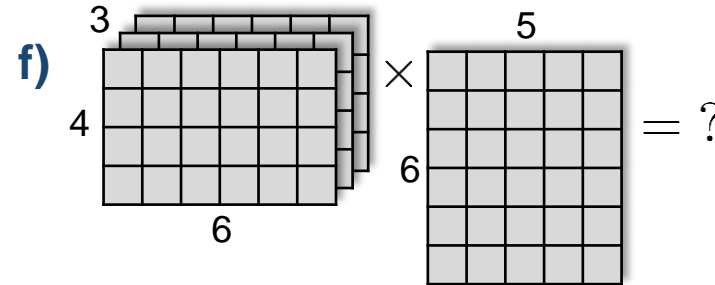
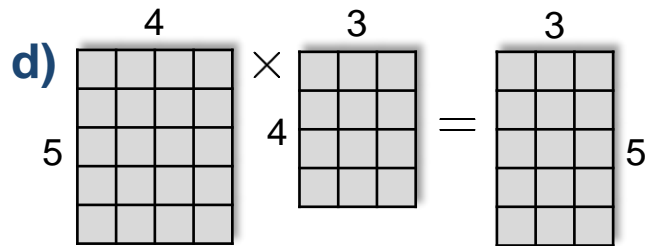
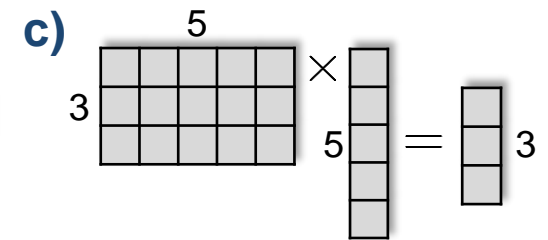
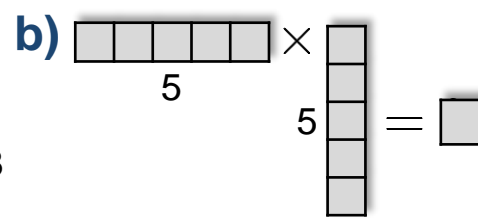
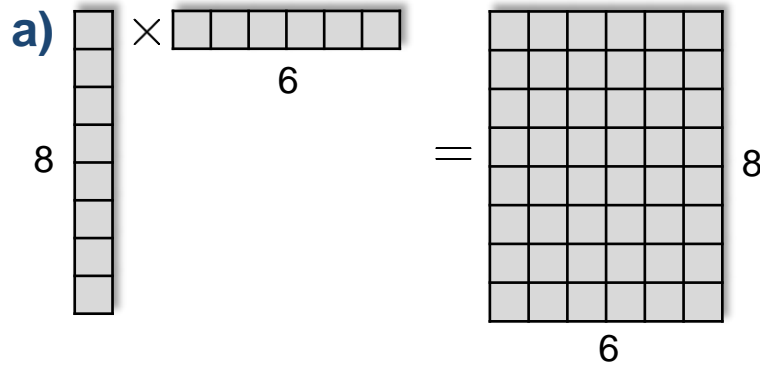
### Tensors

- Some vector/matrix/tensor operations:



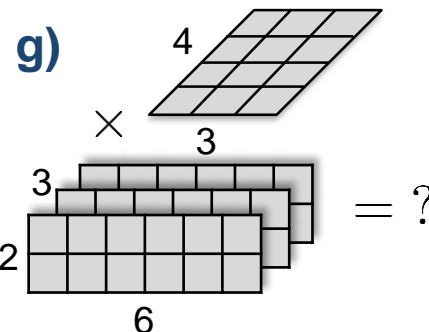
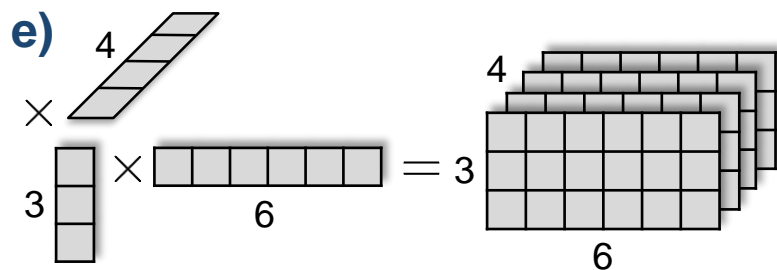
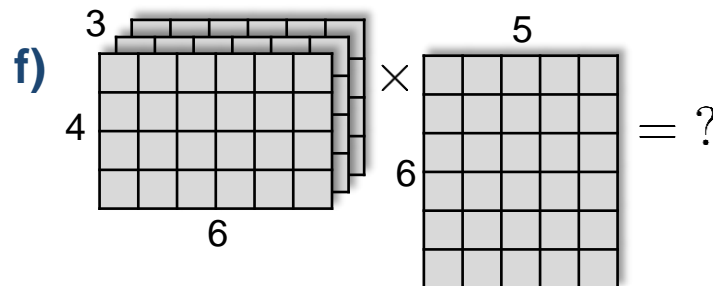
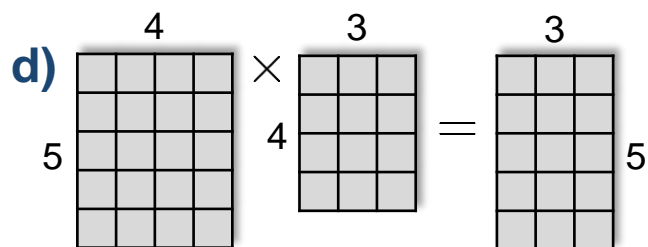
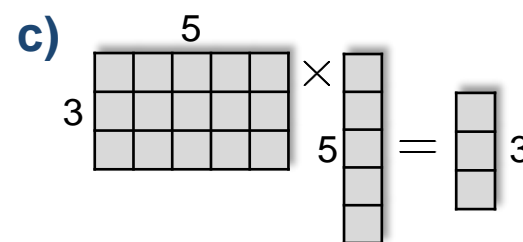
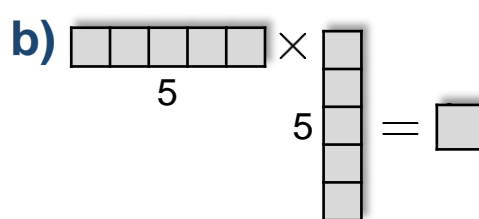
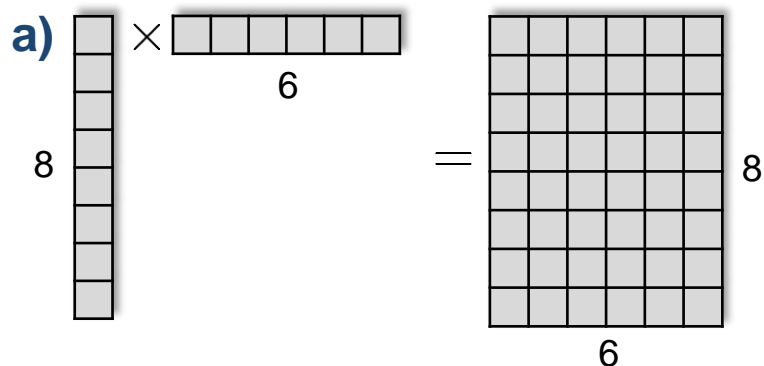
### Tensors

- Some vector/matrix/tensor operations:



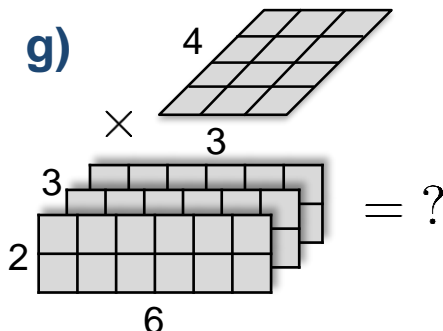
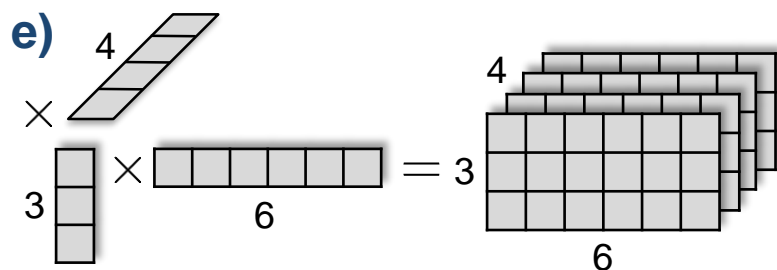
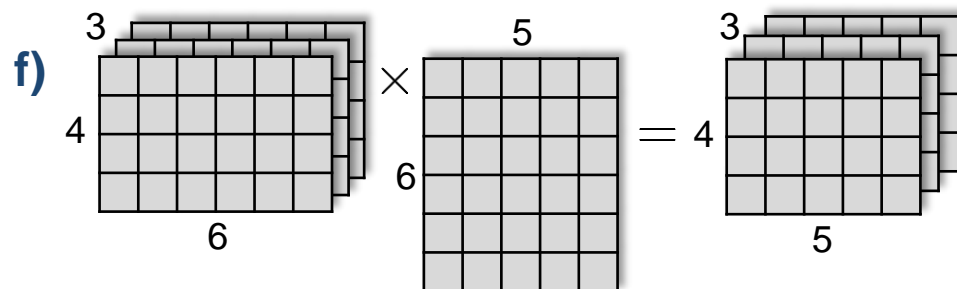
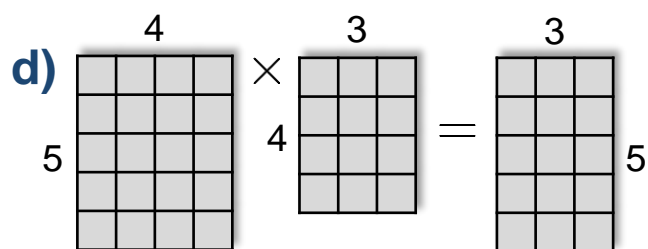
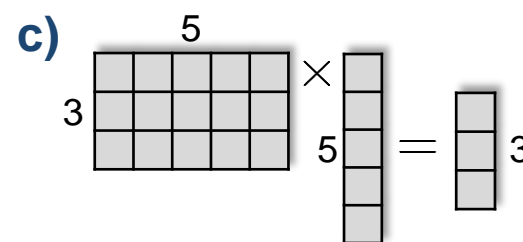
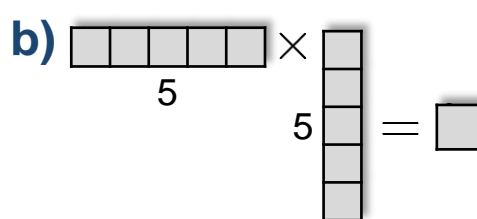
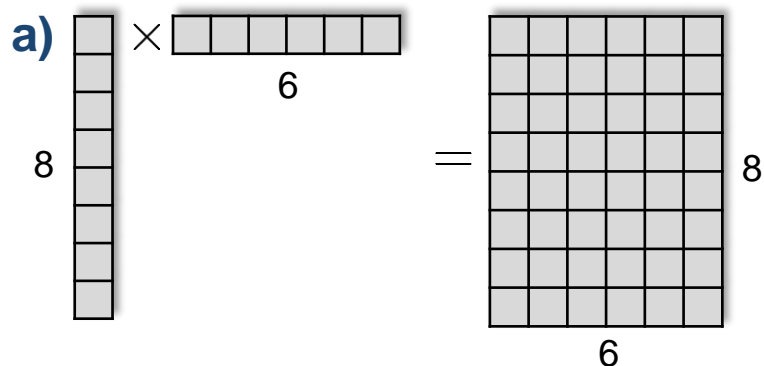
### Tensors

- Some vector/matrix/tensor operations:



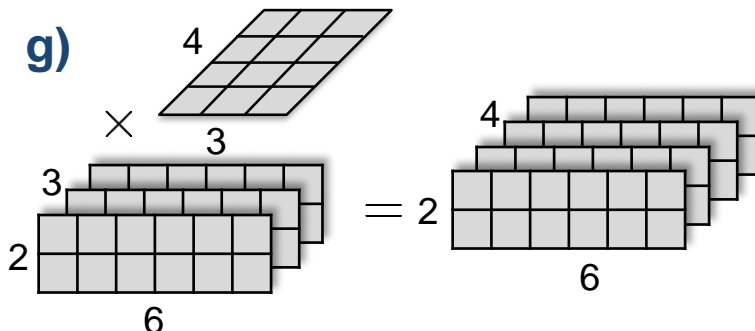
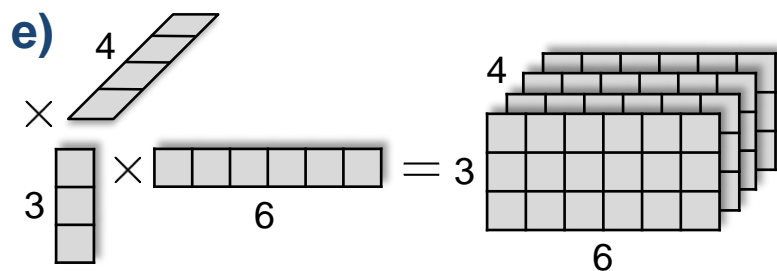
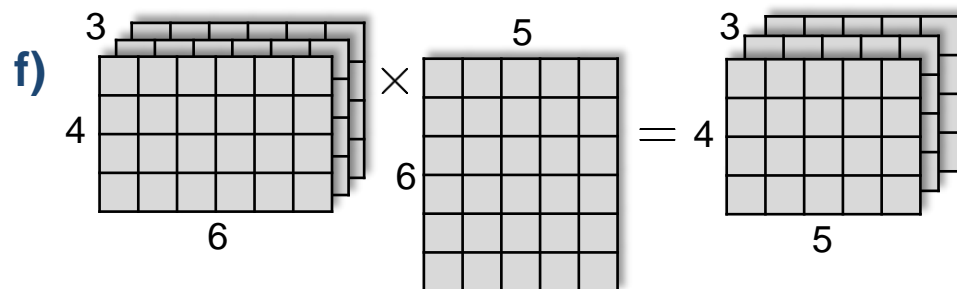
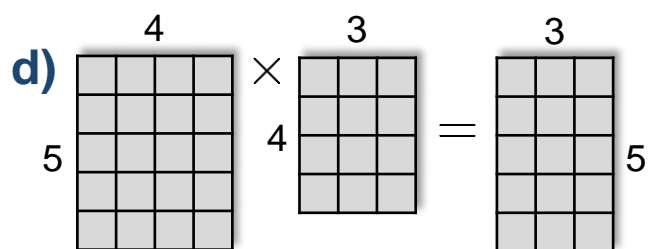
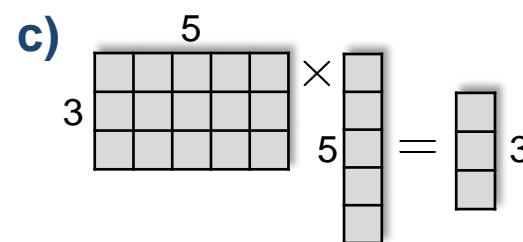
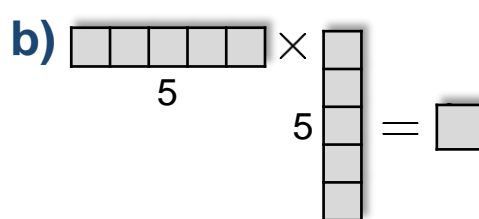
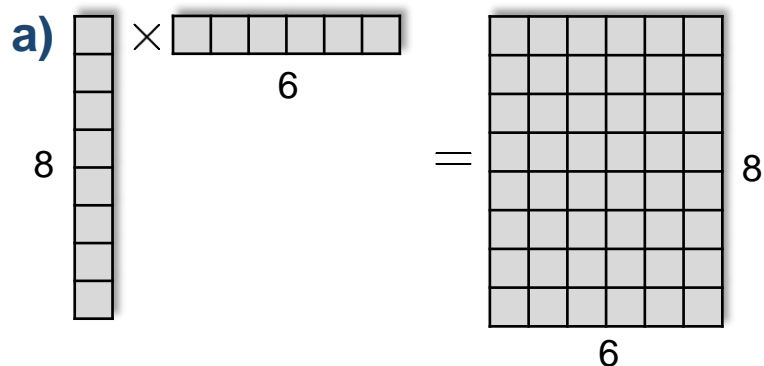
### Tensors

- Some vector/matrix/tensor operations:



### Tensors

- Some vector/matrix/tensor operations:



## Generalizing the derivative

- Let  $y = f(x)$ ,  $x \in \mathbb{R}^D$ ,  $y \in \mathbb{R}^N$ . How can we **define**  $\frac{dy}{dx}$  ?

## Generalizing the derivative

- Let  $y = f(x)$ ,  $x \in \mathbb{R}^D$ ,  $y \in \mathbb{R}^N$ . How can we **define**  $\frac{dy}{dx}$  ?
- Let  $\delta_{x_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^N \\ \mathbf{h} & \mapsto f(x_0 + \mathbf{h}) - f(x_0) \end{cases}$  “the variation of  $f$  when moving by a step  $\mathbf{h}$  from  $x_0$ .”



## Generalizing the derivative

- Let  $y = f(x)$ ,  $x \in \mathbb{R}^D$ ,  $y \in \mathbb{R}^N$ . How can we **define**  $\frac{dy}{dx}$  ?
- Let  $\delta_{x_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^N \\ \mathbf{h} & \mapsto f(x_0 + \mathbf{h}) - f(x_0) \end{cases}$  “the variation of  $f$  when moving by a step  $\mathbf{h}$  from  $x_0$ .”
- $\left. \frac{dy}{dx} \right|_{x=x_0}$  is the **linear approx.** of  $\delta_{x_0}$  for  $\mathbf{h}$  *infinitesimally* small.

## Generalizing the derivative

- Let  $y = f(x)$ ,  $x \in \mathbb{R}^D$ ,  $y \in \mathbb{R}^N$ . How can we **define**  $\frac{dy}{dx}$  ?
- Let  $\delta_{x_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^N \\ \mathbf{h} & \mapsto f(x_0 + \mathbf{h}) - f(x_0) \end{cases}$  “the variation of  $f$  when moving by a step  $\mathbf{h}$  from  $x_0$ .”
- $\left. \frac{dy}{dx} \right|_{x=x_0}$  is the **linear approx.** of  $\delta_{x_0}$  for  $\mathbf{h}$  **infinitesimally small**.
- We call this the **total derivative** of  $f$ , or of  $y$ , with resp. to  $x$ , at  $x_0$ .

## Generalizing the derivative

- Let  $y = f(x)$ ,  $x \in \mathbb{R}^D$ ,  $y \in \mathbb{R}^N$ . How can we **define**  $\frac{dy}{dx}$  ?
- Let  $\delta_{x_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^N \\ \mathbf{h} & \mapsto f(x_0 + \mathbf{h}) - f(x_0) \end{cases}$  “the variation of  $f$  when moving by a step  $\mathbf{h}$  from  $x_0$ .”
- $\left. \frac{dy}{dx} \right|_{x=x_0}$  is the **linear approx.** of  $\delta_{x_0}$  for  $\mathbf{h}$  *infinitesimally* small.
- We call this the **total derivative** of  $f$ , or of  $y$ , with resp. to  $x$ , at  $x_0$ .
- Since  $\left. \frac{dy}{dx} \right|_{x_0}$  is a **linear map** from  $\mathbb{R}^D$  to  $\mathbb{R}^N$ , it can be *represented* by a **matrix**,  $\left. \frac{dy}{dx} \right|_{x_0} = \mathbf{J}_x[f](x_0) \in \mathbb{R}^{N \times D}$ , called the **Jacobian**.

## Generalizing the derivative

- Let  $y = f(x)$ ,  $x \in \mathbb{R}^D$ ,  $y \in \mathbb{R}^N$ . How can we **define**  $\frac{dy}{dx}$  ?
- Let  $\delta_{x_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^N \\ \mathbf{h} & \mapsto f(x_0 + \mathbf{h}) - f(x_0) \end{cases}$  “the variation of  $f$  when moving by a step  $\mathbf{h}$  from  $x_0$ .”
- $\frac{dy}{dx} \Big|_{x=x_0}$  is the **linear approx.** of  $\delta_{x_0}$  for  $\mathbf{h}$  *infinitesimally small*.
- We call this the **total derivative** of  $f$ , or of  $y$ , with resp. to  $x$ , at  $x_0$ .
- Since  $\frac{dy}{dx} \Big|_{x_0}$  is a **linear map** from  $\mathbb{R}^D$  to  $\mathbb{R}^N$ , it can be *represented* by a **matrix**,  $\frac{dy}{dx} \Big|_{x_0} = \mathbf{J}_x[f](x_0) \in \mathbb{R}^{N \times D}$ , called the **Jacobian**.

Formulas:

$$\frac{dy}{dx} \Big|_{x_0} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_1}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_1}{\partial x_D} \Big|_{x_{D,0}} \\ \frac{\partial y_2}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_2}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_2}{\partial x_D} \Big|_{x_{D,0}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_N}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_N}{\partial x_D} \Big|_{x_{D,0}} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1(x_{1,0})}{\partial x_1} & \frac{\partial f_1(x_{2,0})}{\partial x_2} & \cdots & \frac{\partial f_1(x_{D,0})}{\partial D_1} \\ \frac{\partial f_2(x_{1,0})}{\partial x_1} & \frac{\partial f_2(x_{2,0})}{\partial x_2} & \cdots & \frac{\partial f_2(x_{D,0})}{\partial D_1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_N(x_{1,0})}{\partial x_1} & \frac{\partial f_N(x_{2,0})}{\partial x_2} & \cdots & \frac{\partial f_N(x_{D,0})}{\partial D_1} \end{bmatrix} = \mathbf{J}_x[f](x_0)$$

## Special Case 1

What does it give for a simple function  $f : \mathbb{R}^1 \rightarrow \mathbb{R}^1$  ?

## Special Case 1

What does it give for a simple function  $f : \mathbb{R}^1 \rightarrow \mathbb{R}^1$  ?

- In high-school we learn:

$$f'(x_0) = \lim_{h \rightarrow \infty} \frac{f(x_0 + h) - f(x_0)}{h}$$

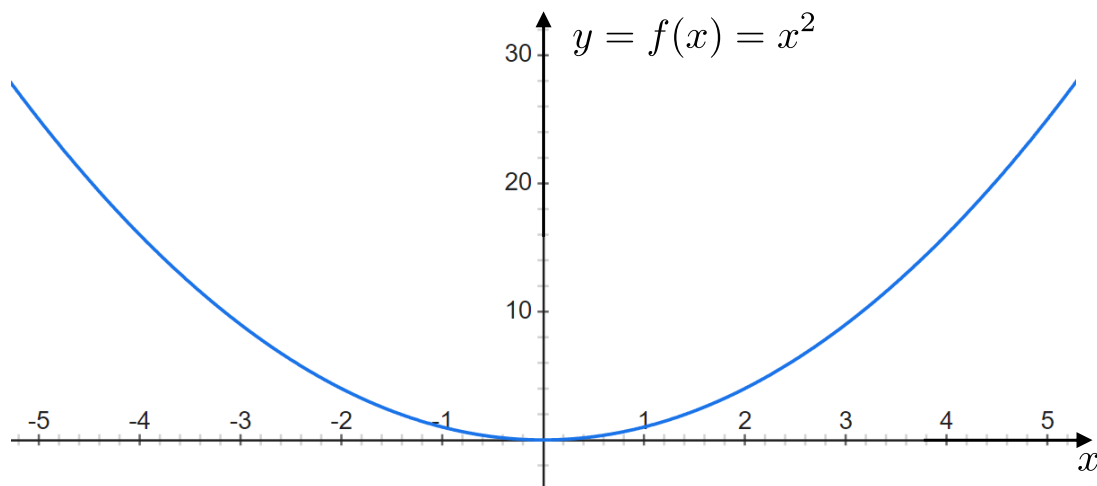
## Special Case 1

What does it give for a simple function  $f : \mathbb{R}^1 \rightarrow \mathbb{R}^1$  ?

- In high-school we learn:

$$f'(x_0) = \lim_{h \rightarrow \infty} \frac{f(x_0 + h) - f(x_0)}{h}$$

- **Ex.**  $y = f(x) = x^2$



## Special Case 1

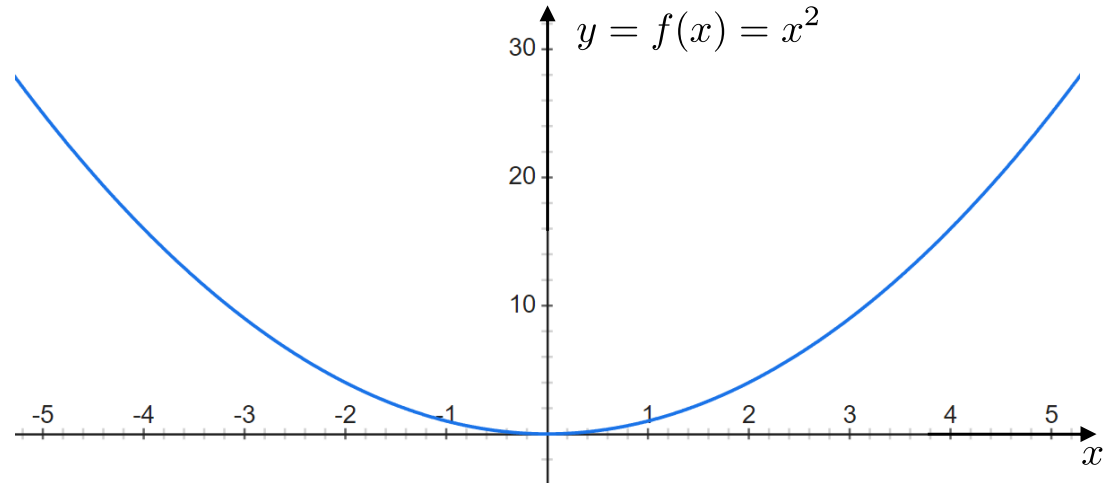
What does it give for a simple function  $f : \mathbb{R}^1 \rightarrow \mathbb{R}^1$  ?

- In high-school we learn:

$$f'(x_0) = \lim_{h \rightarrow \infty} \frac{f(x_0 + h) - f(x_0)}{h}$$

- **Ex.**  $y = f(x) = x^2$

$$f'(x_0) = 2x_0 = \left. \frac{dy}{dx} \right|_{x_0}$$





## Special Case 1

What does it give for a simple function  $f : \mathbb{R}^1 \rightarrow \mathbb{R}^1$  ?

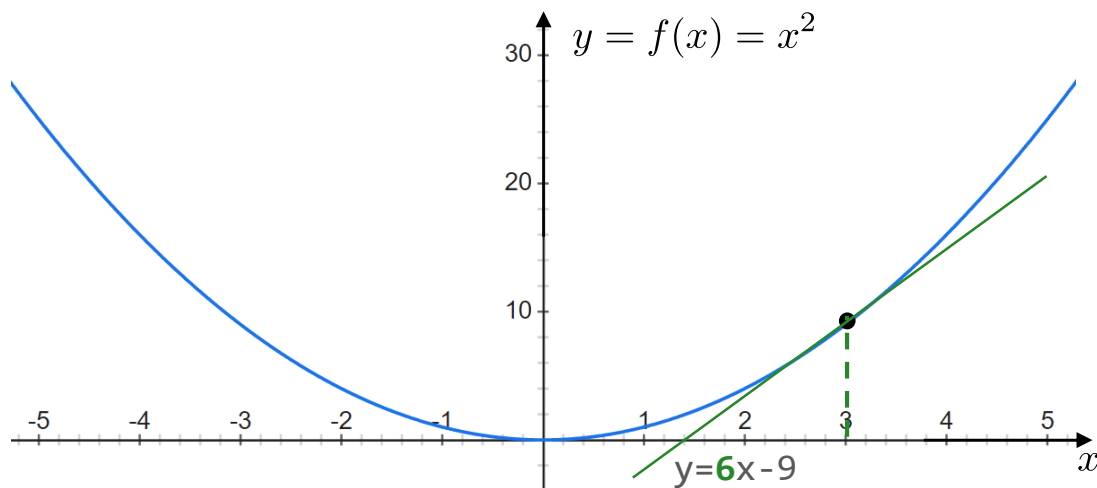
- In high-school we learn:

$$f'(x_0) = \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}$$

- Ex.**  $y = f(x) = x^2$

$$f'(x_0) = 2x_0 = \left. \frac{dy}{dx} \right|_{x_0}$$

$$f'(3) = 6$$



## Special Case 1

What does it give for a simple function  $f : \mathbb{R}^1 \rightarrow \mathbb{R}^1$  ?

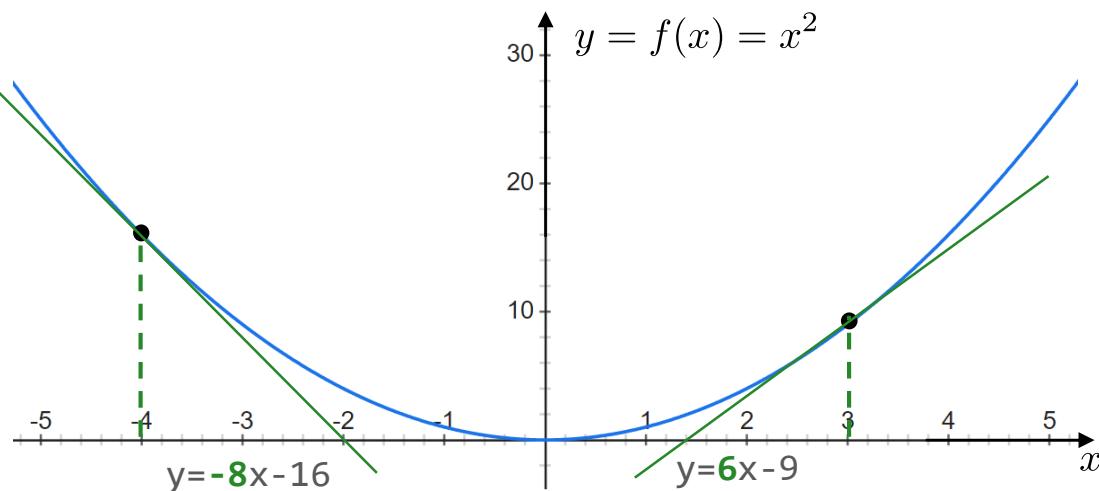
- In high-school we learn:

$$f'(x_0) = \lim_{h \rightarrow \infty} \frac{f(x_0 + h) - f(x_0)}{h}$$

- Ex.**  $y = f(x) = x^2$

$$f'(x_0) = 2x_0 = \left. \frac{dy}{dx} \right|_{x_0}$$

$$f'(3) = 6, \quad f'(-4) = -8$$



# Special Case 1

What does it give for a simple function  $f : \mathbb{R}^1 \rightarrow \mathbb{R}^1$  ?

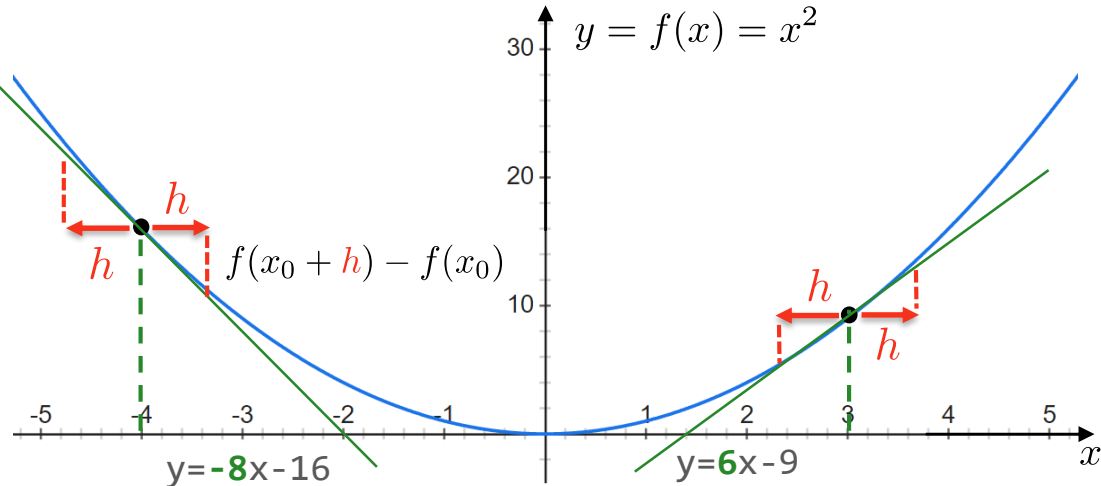
- In high-school we learn:

$$f'(x_0) = \lim_{h \rightarrow \infty} \frac{f(x_0 + h) - f(x_0)}{h}$$

- Ex.**  $y = f(x) = x^2$

$$f'(x_0) = 2x_0 = \left. \frac{dy}{dx} \right|_{x_0}$$

$$f'(3) = 6, \quad f'(-4) = -8$$



- The derivative of  $f$  at  $x_0$  may be viewed as the **linear map** that approximates

$$\delta_{x_0} : \begin{cases} \mathbb{R}^1 & \rightarrow & \mathbb{R}^1 \\ h & \mapsto & f(x_0 + h) - f(x_0) \end{cases} \text{ for } h \text{ infinitesimally small.}$$

*“the variation of  $f$  when moving to the left or to the right from  $x_0$ .”*

# Special Case 1

What does it give for a simple function  $f : \mathbb{R}^1 \rightarrow \mathbb{R}^1$  ?

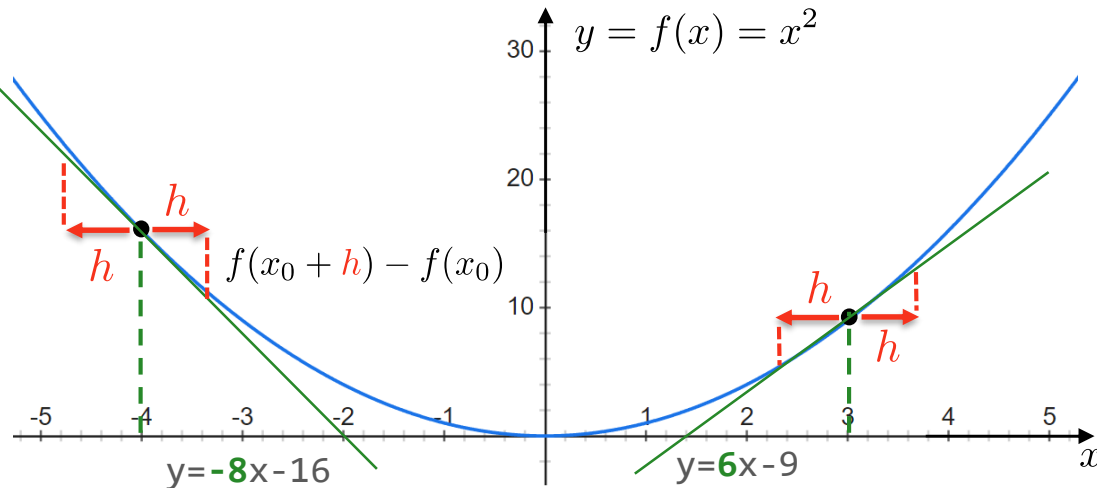
- In high-school we learn:

$$f'(x_0) = \lim_{h \rightarrow \infty} \frac{f(x_0 + h) - f(x_0)}{h}$$

- Ex.**  $y = f(x) = x^2$

$$f'(x_0) = 2x_0 = \left. \frac{dy}{dx} \right|_{x_0}$$

$$f'(3) = 6, f'(-4) = -8$$



- The derivative of  $f$  at  $x_0$  may be viewed as the **linear map** that approximates

$$\delta_{x_0} : \begin{cases} \mathbb{R}^1 & \rightarrow \mathbb{R}^1 \\ h & \mapsto f(x_0 + h) - f(x_0) \end{cases} \text{ for } h \text{ infinitesimally small.}$$

“the variation of  $f$  when moving to the **left** or to the **right** from  $x_0$ .”

Indeed:  $\lim_{h \rightarrow \infty} f(x_0 + h) - f(x_0) = f'(x_0)h !$

# Special Case 1

What does it give for a simple function  $f : \mathbb{R}^1 \rightarrow \mathbb{R}^1$  ?

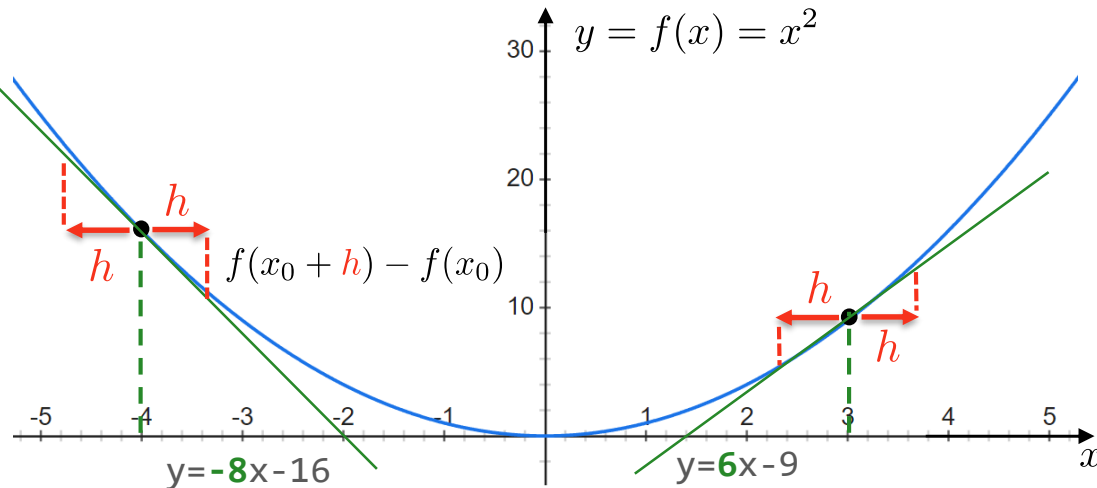
- In high-school we learn:

$$f'(x_0) = \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}$$

- Ex.**  $y = f(x) = x^2$

$$f'(x_0) = 2x_0 = \left. \frac{dy}{dx} \right|_{x_0}$$

$$f'(3) = 6, f'(-4) = -8$$



- The derivative of  $f$  at  $x_0$  may be viewed as the **linear map** that approximates

$$\delta_{x_0} : \begin{cases} \mathbb{R}^1 & \rightarrow \mathbb{R}^1 \\ h & \mapsto f(x_0 + h) - f(x_0) \end{cases} \text{ for } h \text{ infinitesimally small.}$$

“the variation of  $f$  when moving to the **left** or to the **right** from  $x_0$ .”

Indeed:  $\lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h} = f'(x_0)$   $f'(x_0)h$  !

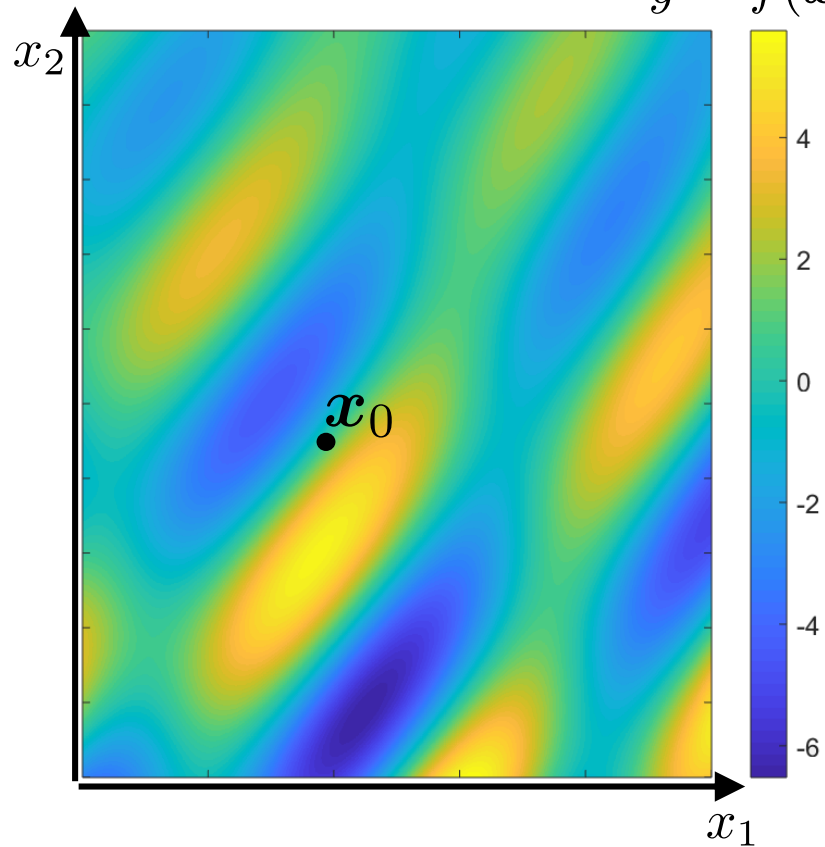
- As a **linear map**,  $f'(x_0)$  can be seen as a **0-way tensor**, in  $\mathbb{R}^{1 \times 1}$  !



## Special Case 2

What about a **real-valued** multivariate function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^1$  ?

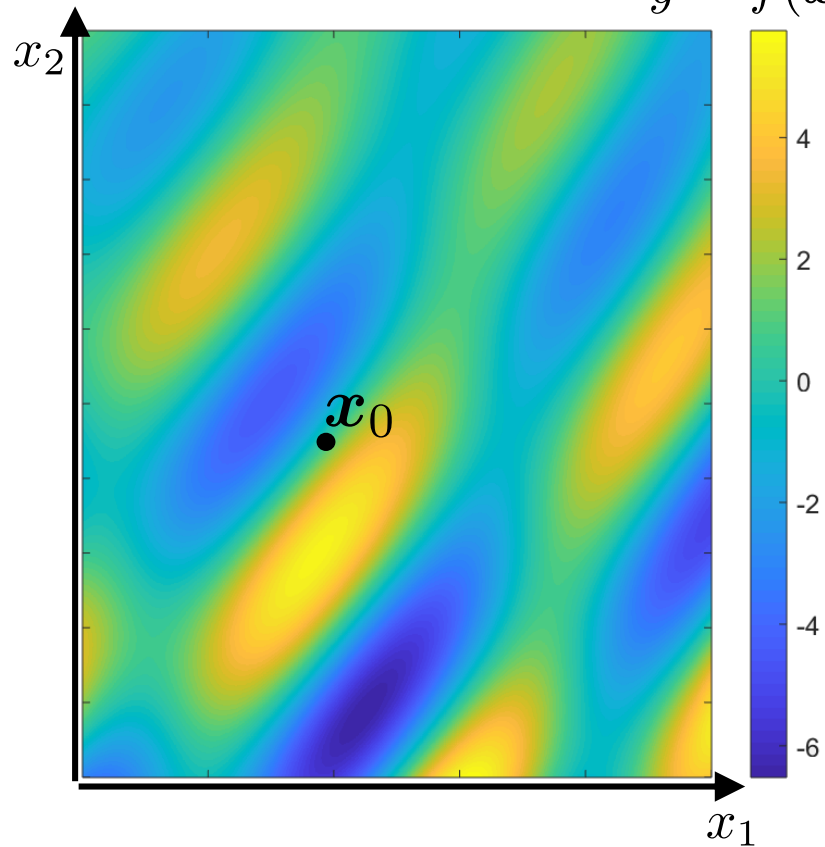
$$y = f(\mathbf{x})$$



## Special Case 2

What about a **real-valued** multivariate function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^1$  ?

$$y = f(\mathbf{x})$$



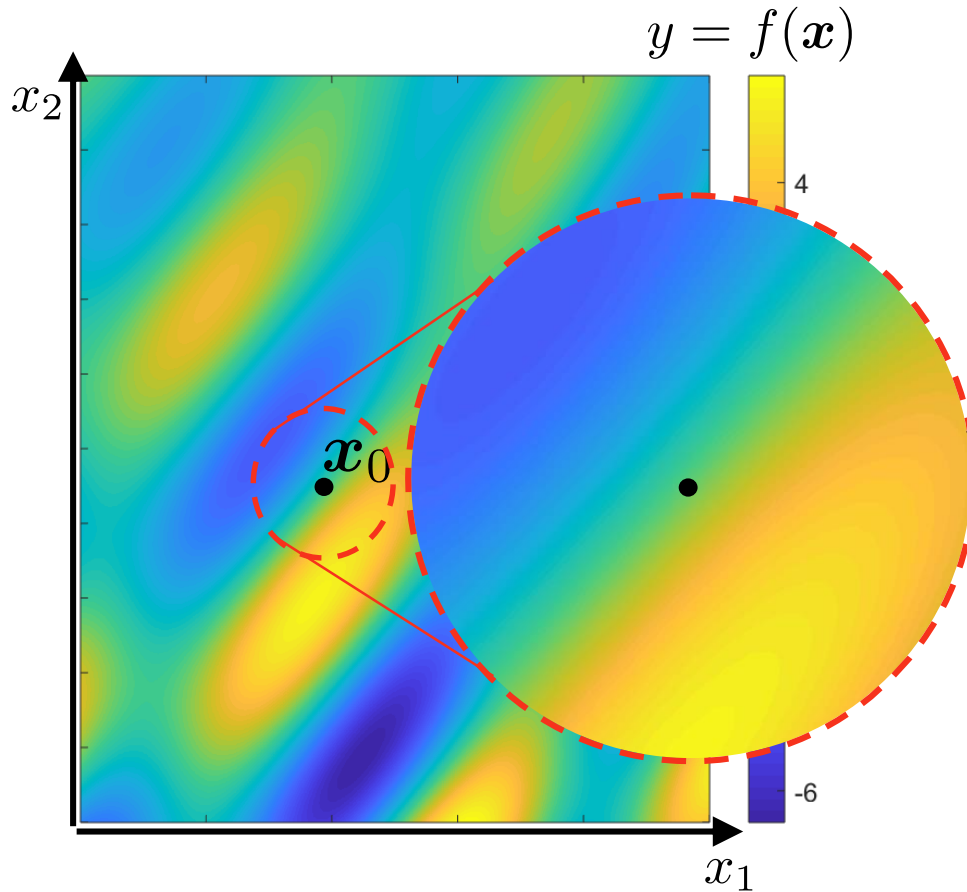
- The **total derivative** of  $f$  at  $x_0$  is the **linear form** that approximates

$$\delta_{x_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^1 \\ \mathbf{h} & \mapsto f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) \end{cases}$$

for  $\mathbf{h}$  infinitesimally small.

## Special Case 2

What about a **real-valued** multivariate function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^1$  ?



- The **total derivative** of  $f$  at  $\mathbf{x}_0$  is the **linear form** that approximates

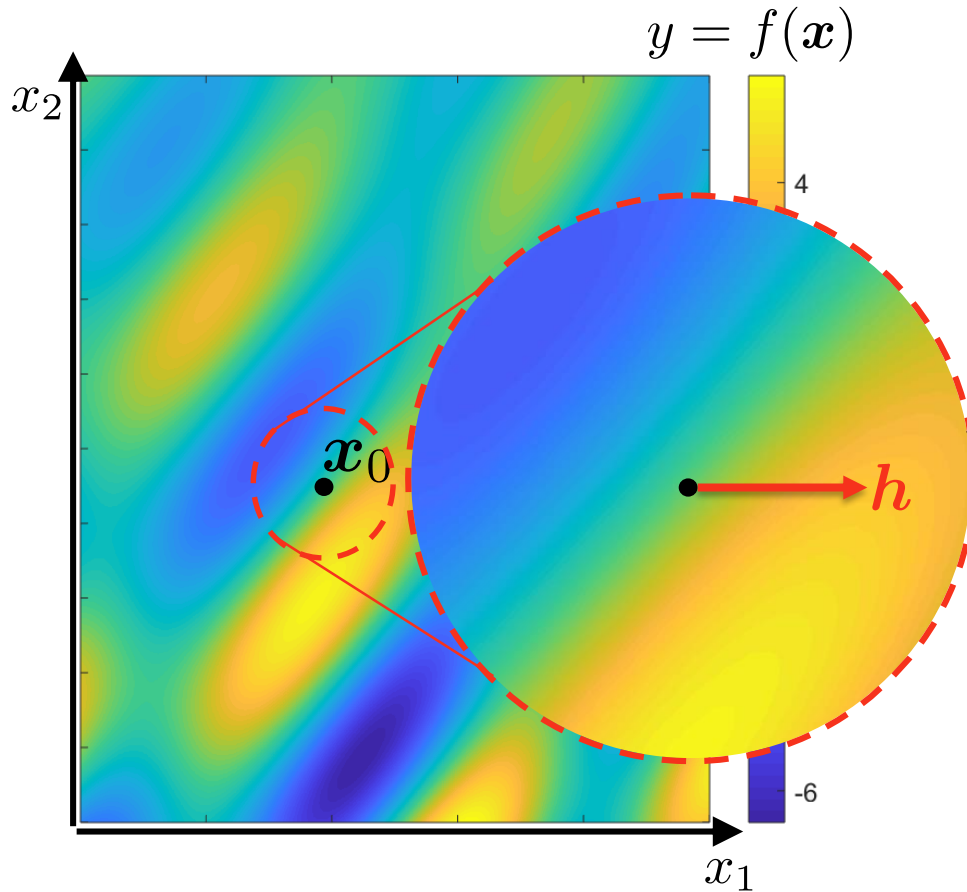
$$\delta_{\mathbf{x}_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^1 \\ \mathbf{h} & \mapsto f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) \end{cases}$$

for  $\mathbf{h}$  infinitesimally small.



## Special Case 2

What about a **real-valued** multivariate function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^1$  ?



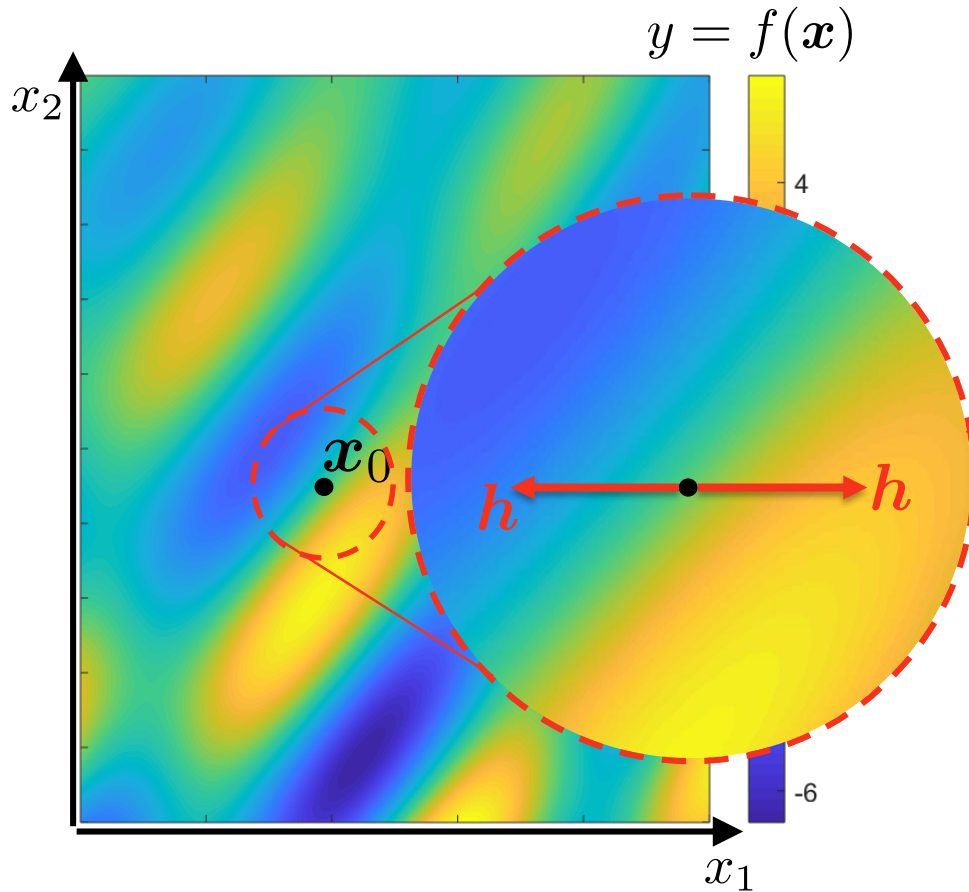
- The **total derivative** of  $f$  at  $\mathbf{x}_0$  is the **linear form** that approximates

$$\delta_{\mathbf{x}_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^1 \\ \mathbf{h} & \mapsto f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) \end{cases}$$

for  $\mathbf{h}$  infinitesimally small.

## Special Case 2

What about a **real-valued** multivariate function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^1$  ?



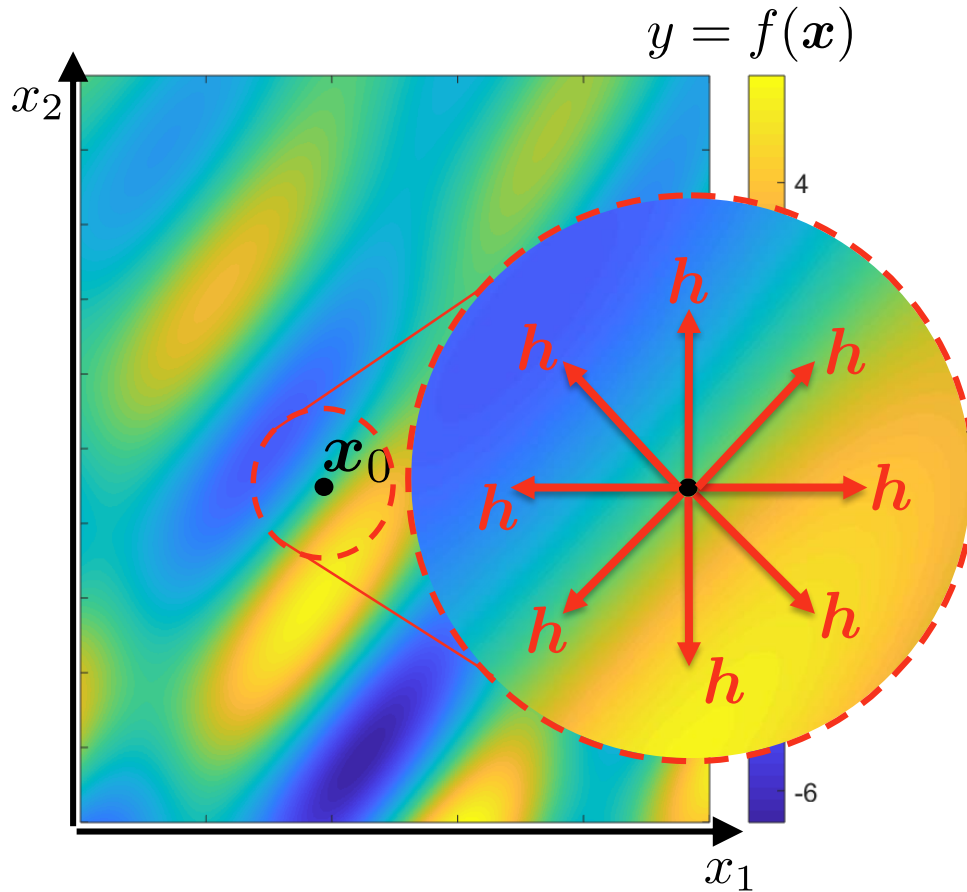
- The **total derivative** of  $f$  at  $\mathbf{x}_0$  is the **linear form** that approximates

$$\delta_{\mathbf{x}_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^1 \\ \mathbf{h} & \mapsto f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) \end{cases}$$

for  $\mathbf{h}$  infinitesimally small.

## Special Case 2

What about a **real-valued** multivariate function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^1$  ?



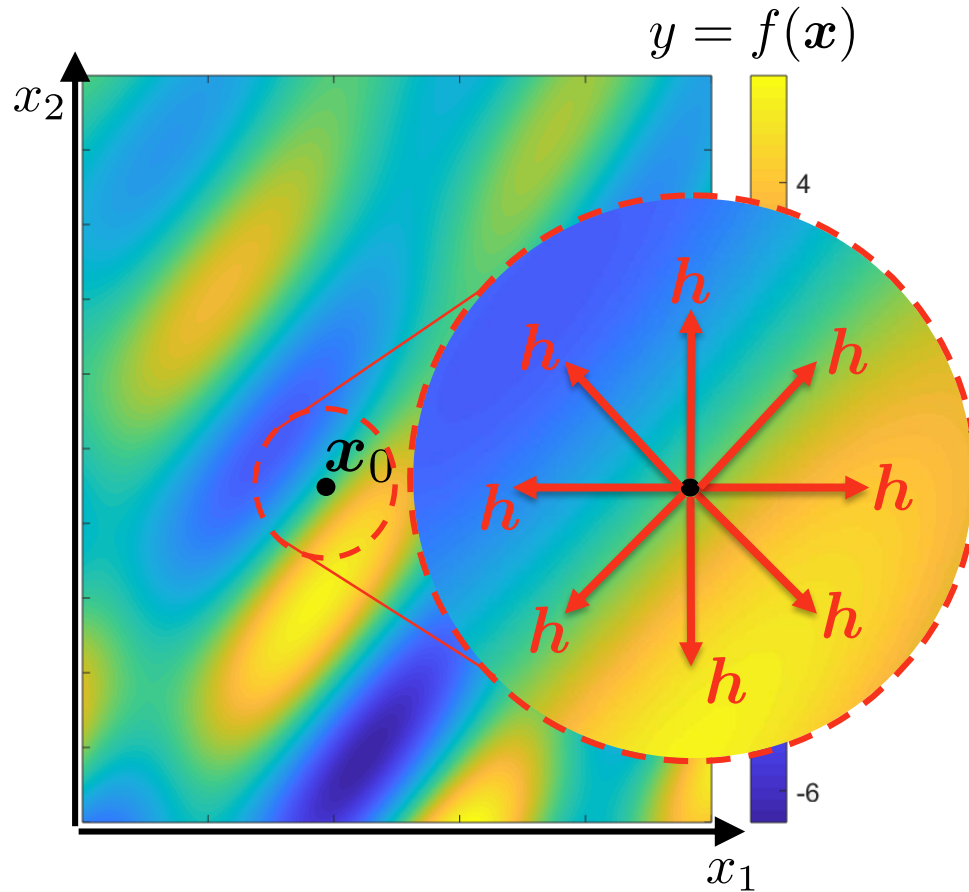
- The **total derivative** of  $f$  at  $\mathbf{x}_0$  is the **linear form** that approximates

$$\delta_{\mathbf{x}_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^1 \\ \mathbf{h} & \mapsto f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) \end{cases}$$

for  $\mathbf{h}$  infinitesimally small.

## Special Case 2

What about a **real-valued** multivariate function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^1$  ?



- The **total derivative** of  $f$  at  $\mathbf{x}_0$  is the **linear form** that approximates

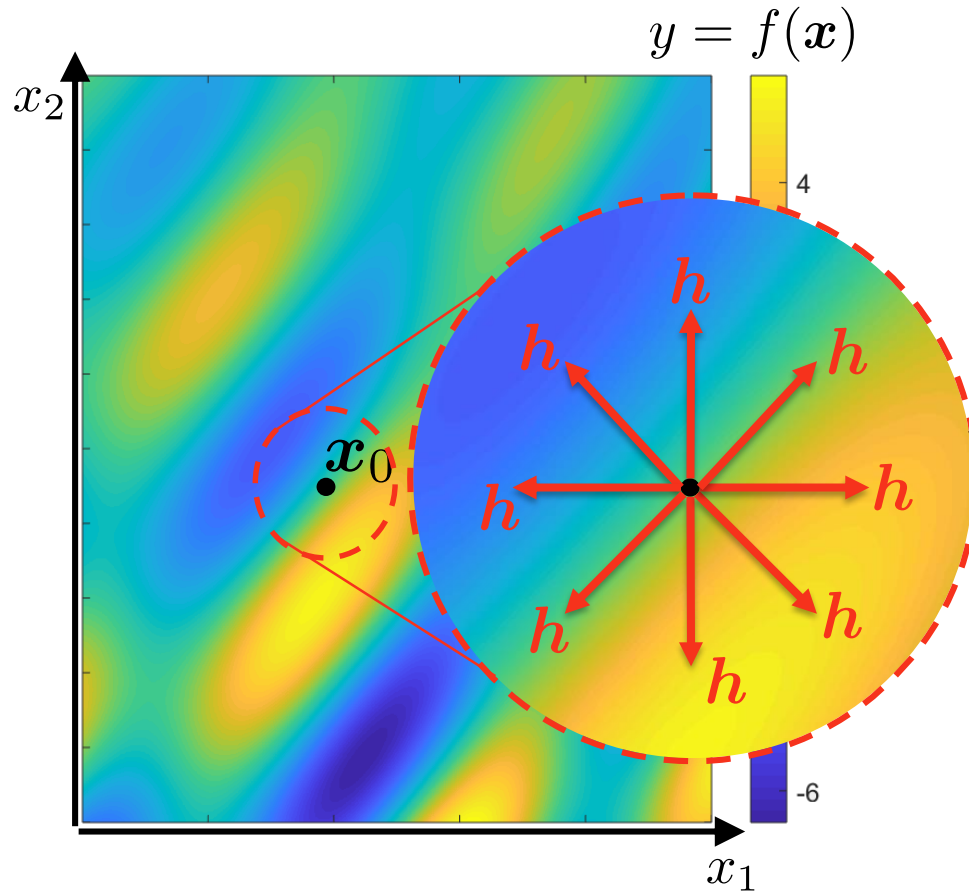
$$\delta_{\mathbf{x}_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^1 \\ \mathbf{h} & \mapsto f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) \end{cases}$$

for  $\mathbf{h}$  infinitesimally small.

- As a linear form, it can be represented by a **row vector**:  $\delta_{\mathbf{x}_0}(\mathbf{h}) \approx \mathbf{w}^\top \mathbf{h}$

## Special Case 2

What about a **real-valued** multivariate function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^1$  ?



- The **total derivative** of  $f$  at  $\mathbf{x}_0$  is the **linear form** that approximates

$$\delta_{\mathbf{x}_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^1 \\ \mathbf{h} & \mapsto f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) \end{cases}$$

for  $\mathbf{h}$  infinitesimally small.

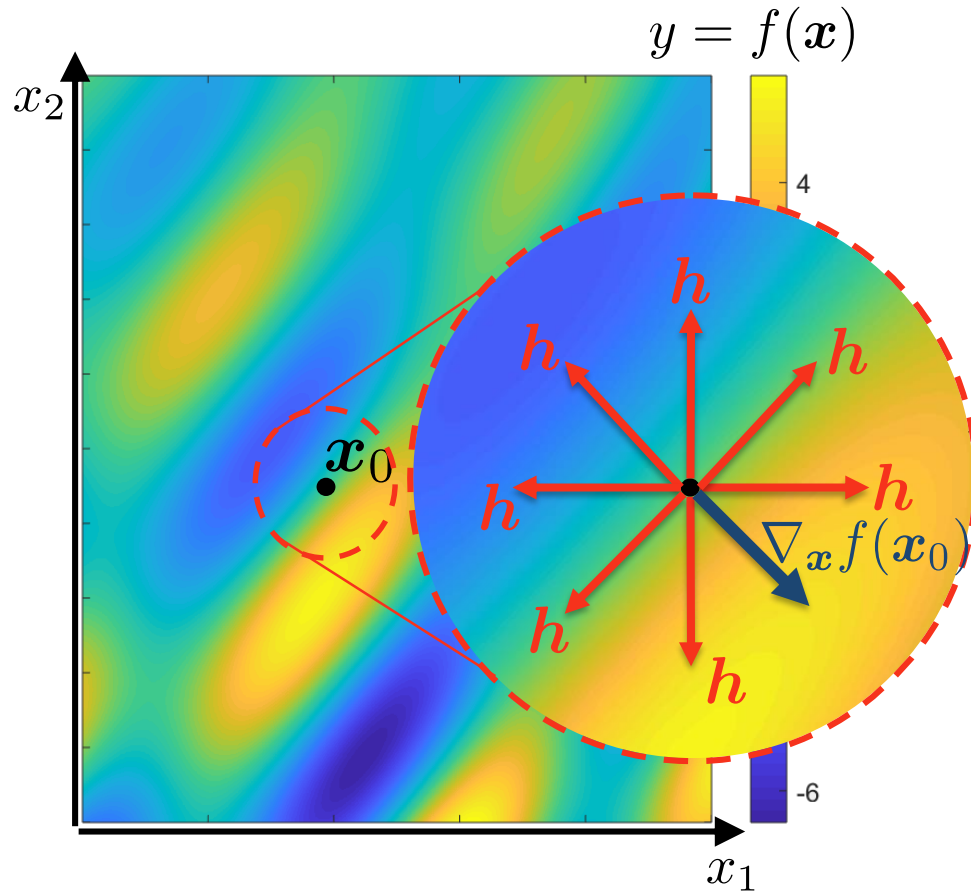
- As a linear form, it can be represented by a **row vector**:  $\delta_{\mathbf{x}_0}(\mathbf{h}) \approx \mathbf{w}^\top \mathbf{h}$
- The **transpose** of this vector is called the **gradient** of  $f$  at  $\mathbf{x}_0$  denoted:

$$\mathbf{w} = \nabla_{\mathbf{x}} f(\mathbf{x}_0) \in \mathbb{R}^D$$



## Special Case 2

What about a **real-valued** multivariate function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^1$  ?



- The **total derivative** of  $f$  at  $x_0$  is the **linear form** that approximates

$$\delta_{x_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^1 \\ \mathbf{h} & \mapsto f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) \end{cases}$$

for  $\mathbf{h}$  infinitesimally small.

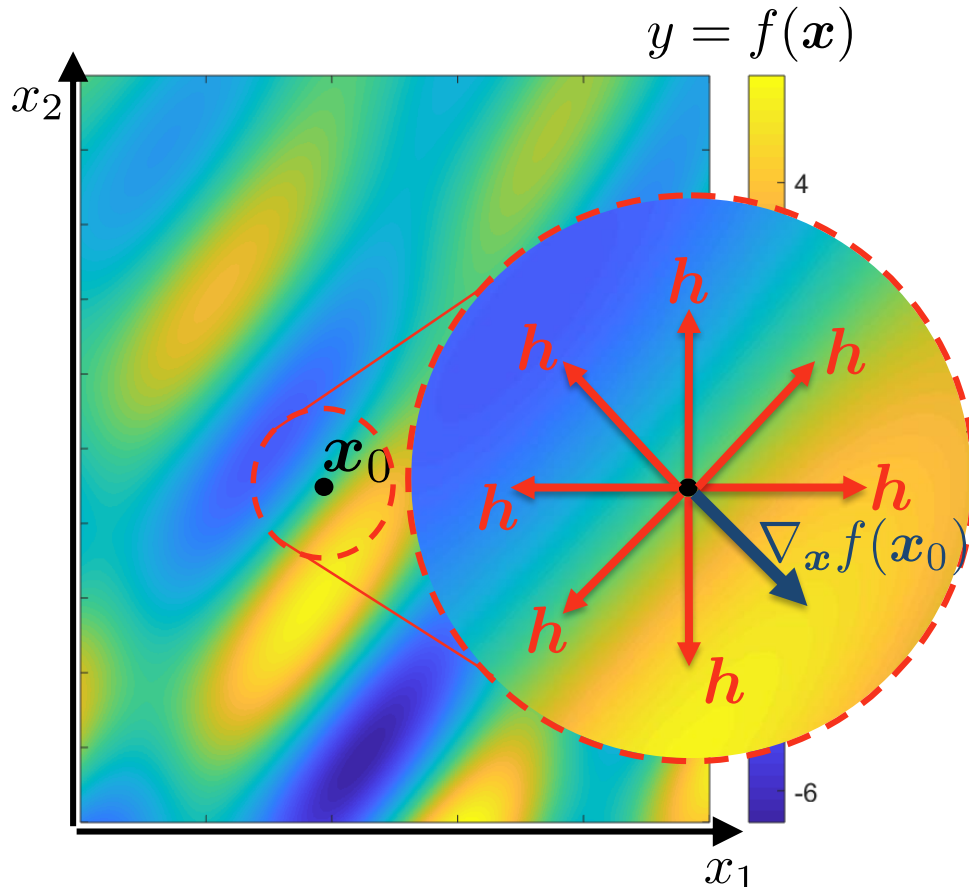
- As a linear form, it can be represented by a **row vector**:  $\delta_{x_0}(\mathbf{h}) \approx \mathbf{w}^\top \mathbf{h}$
- The **transpose** of this vector is called the **gradient** of  $f$  at  $x_0$  denoted:

$$\mathbf{w} = \nabla_x f(\mathbf{x}_0) \in \mathbb{R}^D$$

- It can be interpreted as the **direction and rate of fastest increase** of  $f$  at  $x_0$

## Special Case 2

What about a **real-valued** multivariate function  $f : \mathbb{R}^D \rightarrow \mathbb{R}^1$  ?



- The **total derivative** of  $f$  at  $x_0$  is the **linear form** that approximates

$$\delta_{x_0} : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^1 \\ \mathbf{h} & \mapsto f(\mathbf{x}_0 + \mathbf{h}) - f(\mathbf{x}_0) \end{cases}$$

for  $\mathbf{h}$  infinitesimally small.

- As a linear form, it can be represented by a **row vector**:  $\delta_{x_0}(\mathbf{h}) \approx \mathbf{w}^\top \mathbf{h}$
- The **transpose** of this vector is called the **gradient** of  $f$  at  $x_0$  denoted:

$$\mathbf{w} = \nabla_x f(\mathbf{x}_0) \in \mathbb{R}^D$$

- It can be interpreted as the **direction and rate of fastest increase** of  $f$  at  $x_0$

- The formula is  $\nabla_x f(\mathbf{x}_0) = \left[ \frac{\partial y_1}{\partial x_1} \Big|_{x_{1,0}}, \dots, \frac{\partial y_1}{\partial x_D} \Big|_{x_{D,0}} \right]^\top = \frac{dy}{dx} \Big|_{x_0}^\top$

## Summary ( for $y = f(x)$ )

| Domain of $x$ | Domain of $y$ | Total Derivative at $x_0$ |
|---------------|---------------|---------------------------|
|               |               |                           |
|               |               |                           |
|               |               |                           |
|               |               |                           |



## Summary ( for $y = f(x)$ )

| Domain of $x$ | Domain of $y$            | Total Derivative at $x_0$ |
|---------------|--------------------------|---------------------------|
| $\mathbb{R}$  | $\rightarrow \mathbb{R}$ |                           |
|               |                          |                           |
|               |                          |                           |
|               |                          |                           |

## Summary ( for $y = f(x)$ )

| Domain of $x$ | Domain of $y$ | Total Derivative at $x_0$  |
|---------------|---------------|--|
| $\mathbb{R}$  | $\mathbb{R}$  | $\left. \frac{dy}{dx} \right _{x_0} = f'(x_0) \in \mathbb{R}^{1 \times 1}$ |
|               |               |  |
|               |               |  |
|               |               |  |

## Summary ( for $y = f(x)$ )

| Domain of $x$ | Domain of $y$  | Total Derivative at $x_0$  |
|---------------|----------------|--|
| $\mathbb{R}$  | $\mathbb{R}$   | $\left. \frac{dy}{dx} \right _{x_0} = f'(x_0) \in \mathbb{R}^{1 \times 1}$ |
| $\mathbb{R}$  | $\mathbb{R}^N$ |  |
|               |                |  |
|               |                |  |

## Summary ( for $y = f(x)$ )

| Domain of $x$ | Domain of $y$  | Total Derivative at $x_0$   |
|---------------|----------------|---|
| $\mathbb{R}$  | $\mathbb{R}$   | $\left. \frac{dy}{dx} \right _{x_0} = f'(x_0) \in \mathbb{R}^{1 \times 1}$                            |
| $\mathbb{R}$  | $\mathbb{R}^N$ | $\left. \frac{dy}{dx} \right _{x_0} = [f'_1(x_0), \dots, f'_N(x_0)]^\top \in \mathbb{R}^{N \times 1}$ |
|               |                |   |
|               |                |   |

## Summary ( for $y = f(x)$ )

| Domain of $x$  | Domain of $y$  | Total Derivative at $x_0$   |
|----------------|----------------|---|
| $\mathbb{R}$   | $\mathbb{R}$   | $\left. \frac{dy}{dx} \right _{x_0} = f'(x_0) \in \mathbb{R}^{1 \times 1}$                            |
| $\mathbb{R}$   | $\mathbb{R}^N$ | $\left. \frac{dy}{dx} \right _{x_0} = [f'_1(x_0), \dots, f'_N(x_0)]^\top \in \mathbb{R}^{N \times 1}$ |
| $\mathbb{R}^D$ | $\mathbb{R}$   |   |
|                |                |   |

## Summary ( for $y = f(x)$ )

| Domain of $x$  | Domain of $y$  | Total Derivative at $x_0$   |
|----------------|----------------|---|
| $\mathbb{R}$   | $\mathbb{R}$   | $\left. \frac{dy}{dx} \right _{x_0} = f'(x_0) \in \mathbb{R}^{1 \times 1}$                            |
| $\mathbb{R}$   | $\mathbb{R}^N$ | $\left. \frac{dy}{dx} \right _{x_0} = [f'_1(x_0), \dots, f'_N(x_0)]^\top \in \mathbb{R}^{N \times 1}$ |
| $\mathbb{R}^D$ | $\mathbb{R}$   | $\left. \frac{dy}{dx} \right _{x_0} = \nabla_x f(x_0)^\top \in \mathbb{R}^{1 \times D}$               |
|                |                |   |

## Summary ( for $y = f(x)$ )

| Domain of $x$  | Domain of $y$  | Total Derivative at $x_0$   |
|----------------|----------------|---|
| $\mathbb{R}$   | $\mathbb{R}$   | $\left. \frac{dy}{dx} \right _{x_0} = f'(x_0) \in \mathbb{R}^{1 \times 1}$                            |
| $\mathbb{R}$   | $\mathbb{R}^N$ | $\left. \frac{dy}{dx} \right _{x_0} = [f'_1(x_0), \dots, f'_N(x_0)]^\top \in \mathbb{R}^{N \times 1}$ |
| $\mathbb{R}^D$ | $\mathbb{R}$   | $\left. \frac{dy}{dx} \right _{x_0} = \nabla_x f(x_0)^\top \in \mathbb{R}^{1 \times D}$               |
| $\mathbb{R}^D$ | $\mathbb{R}^N$ |   |

## Summary ( for $y = f(x)$ )

| Domain of $x$  | Domain of $y$  | Total Derivative at $x_0$   |
|----------------|----------------|---|
| $\mathbb{R}$   | $\mathbb{R}$   | $\left. \frac{dy}{dx} \right _{x_0} = f'(x_0) \in \mathbb{R}^{1 \times 1}$                                  |
| $\mathbb{R}$   | $\mathbb{R}^N$ | $\left. \frac{dy}{dx} \right _{x_0} = [f'_1(x_0), \dots, f'_N(x_0)]^\top \in \mathbb{R}^{N \times 1}$       |
| $\mathbb{R}^D$ | $\mathbb{R}$   | $\left. \frac{dy}{dx} \right _{x_0} = \nabla_{\mathbf{x}} f(\mathbf{x}_0)^\top \in \mathbb{R}^{1 \times D}$ |
| $\mathbb{R}^D$ | $\mathbb{R}^N$ | $\left. \frac{dy}{dx} \right _{x_0} = \mathbf{J}_{\mathbf{x}}[f](\mathbf{x}_0) \in \mathbb{R}^{N \times D}$ |



## Summary ( for $y = f(x)$ )

| Domain of $x$  | Domain of $y$  | Total Derivative at $x_0$   |
|----------------|----------------|---|
| $\mathbb{R}$   | $\mathbb{R}$   | $\left. \frac{dy}{dx} \right _{x_0} = f'(x_0) \in \mathbb{R}^{1 \times 1}$                                  |
| $\mathbb{R}$   | $\mathbb{R}^N$ | $\left. \frac{dy}{dx} \right _{x_0} = [f'_1(x_0), \dots, f'_N(x_0)]^\top \in \mathbb{R}^{N \times 1}$       |
| $\mathbb{R}^D$ | $\mathbb{R}$   | $\left. \frac{dy}{dx} \right _{x_0} = \nabla_{\mathbf{x}} f(\mathbf{x}_0)^\top \in \mathbb{R}^{1 \times D}$ |
| $\mathbb{R}^D$ | $\mathbb{R}^N$ | $\left. \frac{dy}{dx} \right _{x_0} = \mathbf{J}_{\mathbf{x}}[f](\mathbf{x}_0) \in \mathbb{R}^{N \times D}$ |

**Note:** All of this generalizes naturally to **matrix** or **tensor variables**, yielding **tensor total derivatives**.

## The Chain Rule

•High school review:  $(g \circ f)'(x_0) = (g' \circ f)(x_0) \cdot f'(x_0)$

## The Chain Rule

- High school review:  $(g \circ f)'(x_0) = (g' \circ f)(x_0) \cdot f'(x_0)$
- By letting  $y = f(x)$  and  $z = g(y) = g(f(x))$  this writes:

$$\left. \frac{dz}{dx} \right|_{x_0} = \left. \frac{dz}{dy} \right|_{f(x_0)} \cdot \left. \frac{dy}{dx} \right|_{x_0}$$

## The Chain Rule

- High school review:  $(g \circ f)'(x_0) = (g' \circ f)(x_0) \cdot f'(x_0)$
- By letting  $y = f(x)$  and  $z = g(y) = g(f(x))$  this writes:

$$\frac{dz}{dx} \Big|_{x_0} = \frac{dz}{dy} \Big|_{f(x_0)} \cdot \frac{dy}{dx} \Big|_{x_0}$$

## The Chain Rule

• High school review:  $(g \circ f)'(x_0) = (g' \circ f)(x_0) \cdot f'(x_0)$

• By letting  $y = f(x)$  and  $z = g(y) = g(f(x))$  this writes:

$$\frac{dz}{dx} \Big|_{x_0} = \frac{dz}{dy} \Big|_{f(x_0)} \cdot \frac{dy}{dx} \Big|_{x_0}$$

## The Chain Rule

- High school review:  $(g \circ f)'(x_0) = (g' \circ f)(x_0) \cdot f'(x_0)$
- By letting  $y = f(x)$  and  $z = g(y) = g(f(x))$  this writes:

$$\left. \frac{dz}{dx} \right|_{x_0} = \left. \frac{dz}{dy} \right|_{f(x_0)} \cdot \left. \frac{dy}{dx} \right|_{x_0}$$

## The Chain Rule

- High school review:  $(g \circ f)'(x_0) = (g' \circ f)(x_0) \cdot f'(x_0)$
- By letting  $y = f(x)$  and  $z = g(y) = g(f(x))$  this writes:

$$\left. \frac{dz}{dx} \right|_{x_0} = \left. \frac{dz}{dy} \right|_{f(x_0)} \cdot \left. \frac{dy}{dx} \right|_{x_0} \quad \text{“the constant } x_0 \text{ ‘flows’ through the variables”}$$

## The Chain Rule

•High school review:  $(g \circ f)'(x_0) = (g' \circ f)(x_0) \cdot f'(x_0)$

•By letting  $y = f(x)$  and  $z = g(y) = g(f(x))$  this writes:

$$\left. \frac{dz}{dx} \right|_{x_0} = \left. \frac{dz}{dy} \right|_{f(x_0)} \cdot \left. \frac{dy}{dx} \right|_{x_0} \quad \text{“the constant } x_0 \text{ ‘flows’ through the variables”}$$

•This generalizes immediately to total derivatives!



## The Chain Rule

- High school review:  $(g \circ f)'(x_0) = (g' \circ f)(x_0) \cdot f'(x_0)$
- By letting  $y = f(x)$  and  $z = g(y) = g(f(x))$  this writes:

$$\left. \frac{dz}{dx} \right|_{x_0} = \left. \frac{dz}{dy} \right|_{f(x_0)} \cdot \left. \frac{dy}{dx} \right|_{x_0} \quad \text{“the constant } x_0 \text{ ‘flows’ through the variables”}$$

- This generalizes immediately to total derivatives!

- Let:
- $x \in \mathbb{R}^D$
  - $y = f(x) \in \mathbb{R}$
  - $z = g(y) = g(f(x)) \in \mathbb{R}^N$
  - $w = h(z) = h(g(f(x))) \in \mathbb{R}^M$

# The Chain Rule

- High school review:  $(g \circ f)'(x_0) = (g' \circ f)(x_0) \cdot f'(x_0)$
- By letting  $y = f(x)$  and  $z = g(y) = g(f(x))$  this writes:

$$\frac{dz}{dx} \Big|_{x_0} = \frac{dz}{dy} \Big|_{f(x_0)} \cdot \frac{dy}{dx} \Big|_{x_0} \quad \text{“the constant } x_0 \text{ flows’ through the variables”}$$

- This generalizes immediately to total derivatives!

Let:

- $x \in \mathbb{R}^D$
- $y = f(x) \in \mathbb{R}$
- $z = g(y) = g(f(x)) \in \mathbb{R}^N$
- $w = h(z) = h(g(f(x))) \in \mathbb{R}^M$

$$\left. \begin{array}{l} \mathbb{R}^{M \times D} \\ \mathbb{R}^{M \times N} \\ \mathbb{R}^{N \times 1} \\ \mathbb{R}^{1 \times D} \end{array} \right\} \frac{dw}{dx} \Big|_{x_0} = \frac{dw}{dz} \Big|_{z_0} \times \frac{dz}{dy} \Big|_{y_0} \times \frac{dy}{dx} \Big|_{x_0}$$

# The Chain Rule

- High school review:  $(g \circ f)'(x_0) = (g' \circ f)(x_0) \cdot f'(x_0)$
- By letting  $y = f(x)$  and  $z = g(y) = g(f(x))$  this writes:

$$\frac{dz}{dx} \Big|_{x_0} = \frac{dz}{dy} \Big|_{f(x_0)} \cdot \frac{dy}{dx} \Big|_{x_0} \quad \text{“the constant } x_0 \text{ flows’ through the variables”}$$

- This generalizes immediately to total derivatives!

Let:

- $x \in \mathbb{R}^D$
- $y = f(x) \in \mathbb{R}$
- $z = g(y) = g(f(x)) \in \mathbb{R}^N$
- $w = h(z) = h(g(f(x))) \in \mathbb{R}^M$

$$\left. \begin{array}{l} \mathbb{R}^{M \times D} \\ \mathbb{R}^{M \times N} \\ \mathbb{R}^{N \times 1} \\ \mathbb{R}^{1 \times D} \end{array} \right\} \frac{dw}{dx} \Big|_{x_0} = \frac{dw}{dz} \Big|_{z_0} \times \frac{dz}{dy} \Big|_{y_0} \times \frac{dy}{dx} \Big|_{x_0}$$

- In functional notations:

$$\mathbf{J}_x [h \circ g \circ f](x_0) = \mathbf{J}_z [h](g \circ f(x_0)) \times \begin{bmatrix} g'_1 \circ f(x_0) \\ \vdots \\ g'_N \circ f(x_0) \end{bmatrix} \times \nabla_x f(x_0)^\top \quad \text{😬}$$

# The Chain Rule

- High school review:  $(g \circ f)'(x_0) = (g' \circ f)(x_0) \cdot f'(x_0)$
- By letting  $y = f(x)$  and  $z = g(y) = g(f(x))$  this writes:

$$\frac{dz}{dx} \Big|_{x_0} = \frac{dz}{dy} \Big|_{f(x_0)} \cdot \frac{dy}{dx} \Big|_{x_0} \quad \text{“the constant } x_0 \text{ flows’ through the variables”}$$

- This generalizes immediately to total derivatives!

Let:

- $x \in \mathbb{R}^D$
- $y = f(x) \in \mathbb{R}$
- $z = g(y) = g(f(x)) \in \mathbb{R}^N$
- $w = h(z) = h(g(f(x))) \in \mathbb{R}^M$

$$\left. \begin{array}{l} \mathbb{R}^{M \times D} \\ \mathbb{R}^{M \times N} \\ \mathbb{R}^{N \times 1} \\ \mathbb{R}^{1 \times D} \end{array} \right\} \frac{dw}{dx} \Big|_{x_0} = \frac{dw}{dz} \Big|_{z_0} \times \frac{dz}{dy} \Big|_{y_0} \times \frac{dy}{dx} \Big|_{x_0}$$

- In functional notations:

$$\mathbf{J}_x [h \circ g \circ f](x_0) = \mathbf{J}_z [h](g \circ f(x_0)) \times \begin{bmatrix} g'_1 \circ f(x_0) \\ \vdots \\ g'_N \circ f(x_0) \end{bmatrix} \times \nabla_x f(x_0)^\top \quad \text{😬}$$

- All of this has natural generalizations to **tensors**

# The Chain Rule

- High school review:  $(g \circ f)'(x_0) = (g' \circ f)(x_0) \cdot f'(x_0)$
- By letting  $y = f(x)$  and  $z = g(y) = g(f(x))$  this writes:

$$\frac{dz}{dx} \Big|_{x_0} = \frac{dz}{dy} \Big|_{f(x_0)} \cdot \frac{dy}{dx} \Big|_{x_0} \quad \text{“the constant } x_0 \text{ flows’ through the variables”}$$

- This generalizes immediately to total derivatives!

Let:

- $x \in \mathbb{R}^D$
- $y = f(x) \in \mathbb{R}$
- $z = g(y) = g(f(x)) \in \mathbb{R}^N$
- $w = h(z) = h(g(f(x))) \in \mathbb{R}^M$

$$\left. \begin{array}{l} \mathbb{R}^{M \times D} \\ \uparrow \\ \frac{dw}{dx} \Big|_{x_0} \end{array} \right\} = \left. \begin{array}{l} \mathbb{R}^{M \times N} \\ \uparrow \\ \frac{dw}{dz} \Big|_{z_0} \end{array} \right\} \times \left. \begin{array}{l} \mathbb{R}^{N \times 1} \\ \uparrow \\ \frac{dz}{dy} \Big|_{y_0} \end{array} \right\} \times \left. \begin{array}{l} \mathbb{R}^{1 \times D} \\ \uparrow \\ \frac{dy}{dx} \Big|_{x_0} \end{array} \right\}$$

- In functional notations:

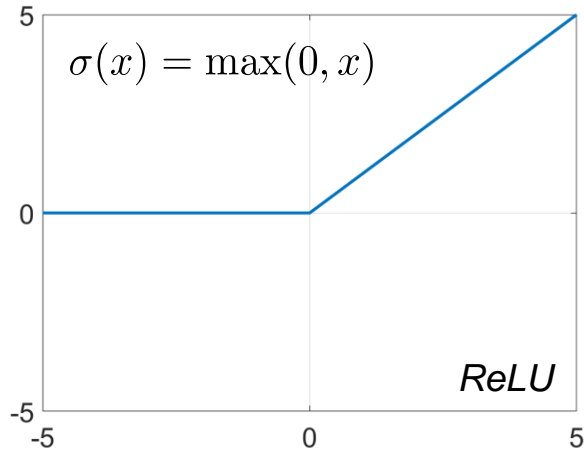
$$\mathbf{J}_x [h \circ g \circ f](x_0) = \mathbf{J}_z [h](g \circ f(x_0)) \times \begin{bmatrix} g'_1 \circ f(x_0) \\ \vdots \\ g'_N \circ f(x_0) \end{bmatrix} \times \nabla_x f(x_0)^\top \quad \text{😬}$$

- All of this has natural generalizations to **tensors**

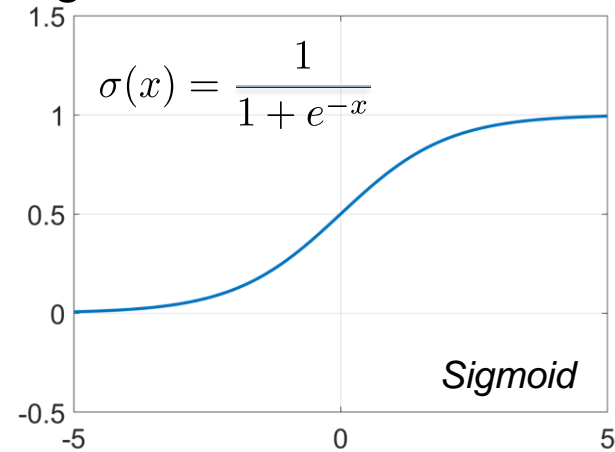
- All of this generalizes to **partial derivatives**, e.g.,  $\frac{\partial z}{\partial x} \Big|_{x_0} = \frac{\partial z}{\partial y} \Big|_{y_0} \times \frac{\partial y}{\partial x} \Big|_{x_0}$

## Exercises: *Simple Derivatives*

- Calculate the derivatives of the following functions:



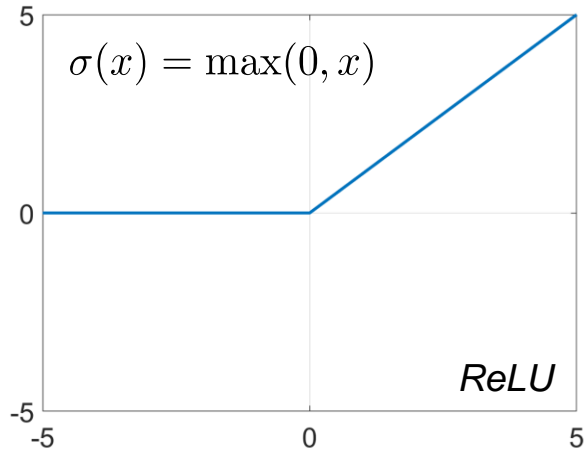
$$\sigma'(x) = ?$$



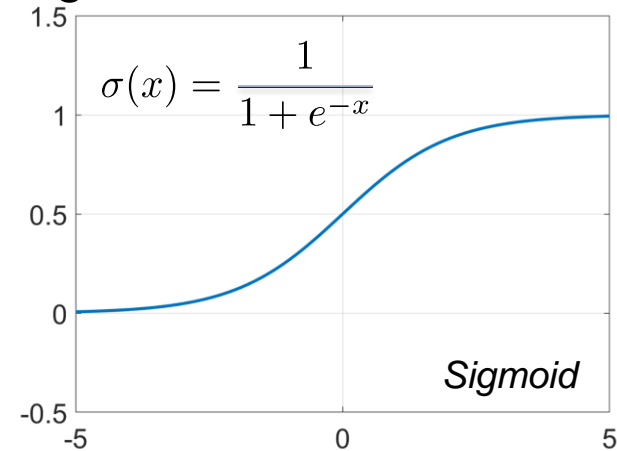
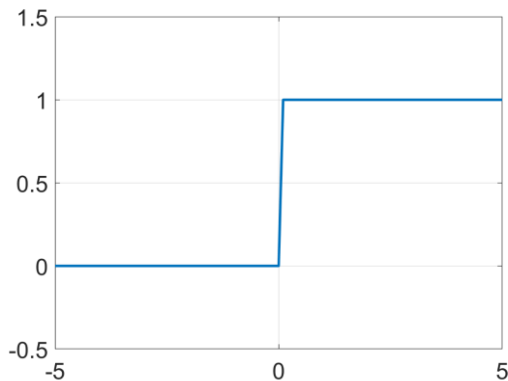
$$\sigma'(x) = ?$$

## Exercises: *Simple Derivatives*

- Calculate the derivatives of the following functions:



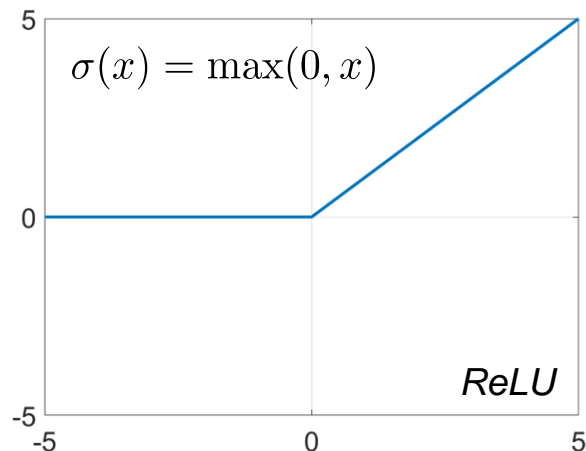
$$\sigma'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x > 0 \end{cases} = H(x)$$



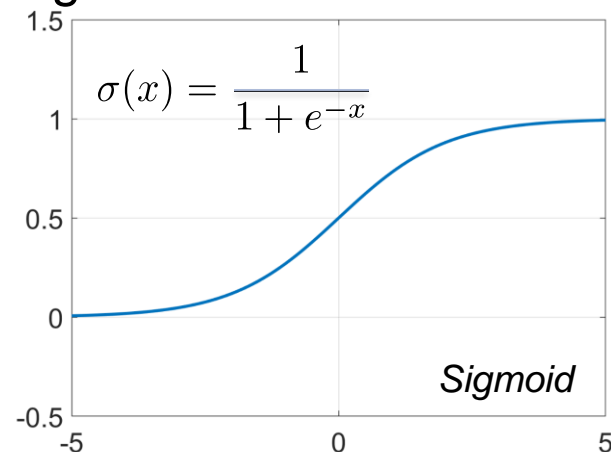
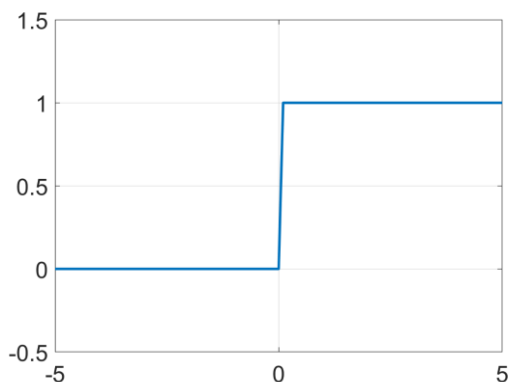
$$\sigma'(x) = ?$$

## Exercises: *Simple Derivatives*

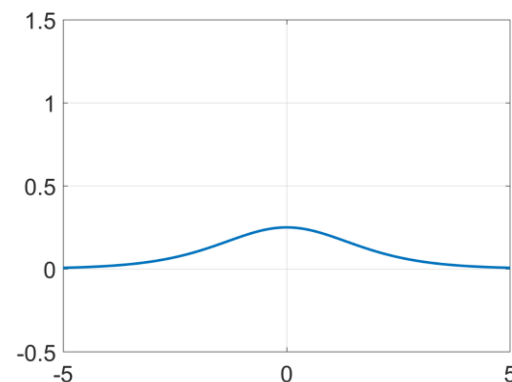
- Calculate the derivatives of the following functions:



$$\sigma'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x > 0 \end{cases} = H(x)$$



$$\sigma'(x) = \frac{e^{-x}}{(1 + e^{-x})^2} = \sigma(x)(1 - \sigma(x))$$





## Exercises: *Gradients*

$$\nabla_{\mathbf{x}} f(\mathbf{x}_0) = \left[ \frac{\partial f}{\partial x_1}(x_{0,1}), \dots, \frac{\partial f}{\partial x_D}(x_{0,D}) \right]^\top$$

- Calculate the gradient of the following functions at  $\mathbf{x}_0$ :

$$f(\mathbf{x}) = b + \mathbf{w}^\top \mathbf{x} \quad \Rightarrow \quad \nabla_{\mathbf{x}} f(\mathbf{x}_0) = ?$$

$$g(\mathbf{x}) = \|\mathbf{x}\|_2^2 \quad \Rightarrow \quad \nabla_{\mathbf{x}} g(\mathbf{x}_0) = ?$$

## Exercises: *Gradients*

$$\nabla_{\mathbf{x}} f(\mathbf{x}_0) = \left[ \frac{\partial f}{\partial x_1}(x_{0,1}), \dots, \frac{\partial f}{\partial x_D}(x_{0,D}) \right]^\top$$

- Calculate the gradient of the following functions at  $\mathbf{x}_0$  :

$$f(\mathbf{x}) = b + \mathbf{w}^\top \mathbf{x} = b + \sum_{d=1}^D w_d x_d \quad \Rightarrow \quad \nabla_{\mathbf{x}} f(\mathbf{x}_0) = ?$$

$$g(\mathbf{x}) = \|\mathbf{x}\|_2^2 = \sum_{d=1}^D (x_d)^2 \quad \Rightarrow \quad \nabla_{\mathbf{x}} g(\mathbf{x}_0) = ?$$

## Exercises: *Gradients*

$$\nabla_{\mathbf{x}} f(\mathbf{x}_0) = \left[ \frac{\partial f}{\partial x_1}(x_{0,1}), \dots, \frac{\partial f}{\partial x_D}(x_{0,D}) \right]^\top$$

- Calculate the gradient of the following functions at  $\mathbf{x}_0$ :

$$f(\mathbf{x}) = b + \mathbf{w}^\top \mathbf{x} = b + \sum_{d=1}^D w_d x_d \quad \Rightarrow \quad \nabla_{\mathbf{x}} f(\mathbf{x}_0) = [w_1, \dots, w_D]^\top = \mathbf{w}$$

$$g(\mathbf{x}) = \|\mathbf{x}\|_2^2 = \sum_{d=1}^D (x_d)^2 \quad \Rightarrow \quad \nabla_{\mathbf{x}} g(\mathbf{x}_0) = ?$$

## Exercises: *Gradients*

$$\nabla_{\mathbf{x}} f(\mathbf{x}_0) = \left[ \frac{\partial f}{\partial x_1}(x_{0,1}), \dots, \frac{\partial f}{\partial x_D}(x_{0,D}) \right]^\top$$

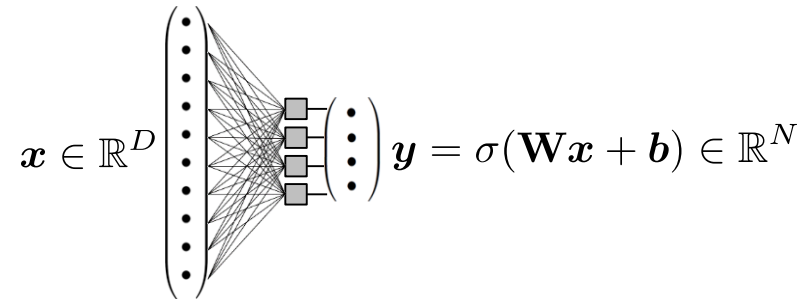
- Calculate the gradient of the following functions at  $\mathbf{x}_0$ :

$$f(\mathbf{x}) = b + \mathbf{w}^\top \mathbf{x} = b + \sum_{d=1}^D w_d x_d \quad \Rightarrow \quad \nabla_{\mathbf{x}} f(\mathbf{x}_0) = [w_1, \dots, w_D]^\top = \mathbf{w}$$

$$g(\mathbf{x}) = \|\mathbf{x}\|_2^2 = \sum_{d=1}^D (x_d)^2 \quad \Rightarrow \quad \nabla_{\mathbf{x}} g(\mathbf{x}_0) = [2x_{0,1}, \dots, 2x_{0,D}]^\top \\ = 2\mathbf{x}_0$$

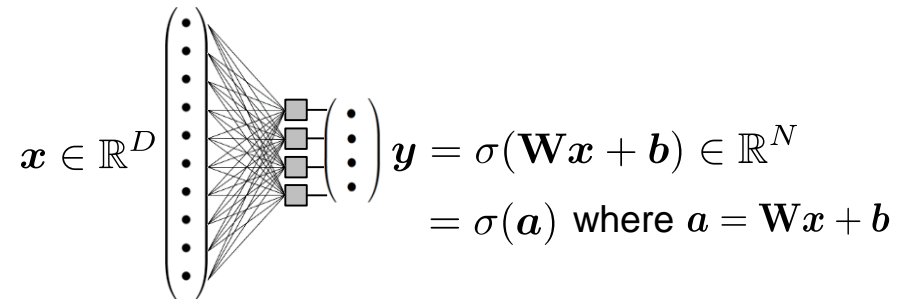
## Exercises: *Jacobians*

*Remember: the simple perceptron*



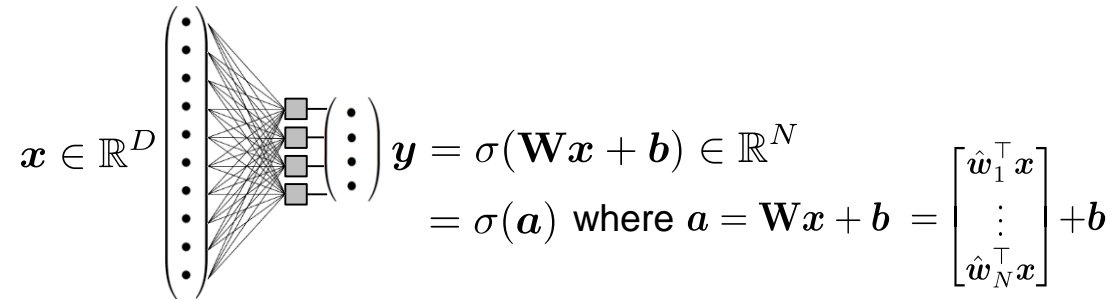
## Exercises: *Jacobians*

*Remember: the simple perceptron*



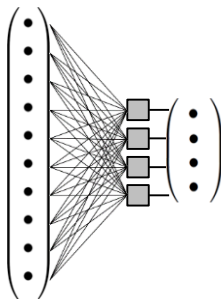
## Exercises: *Jacobians*

**Remember:** the *simple perceptron*



## Exercises: *Jacobians*

**Remember:** the *simple perceptron*

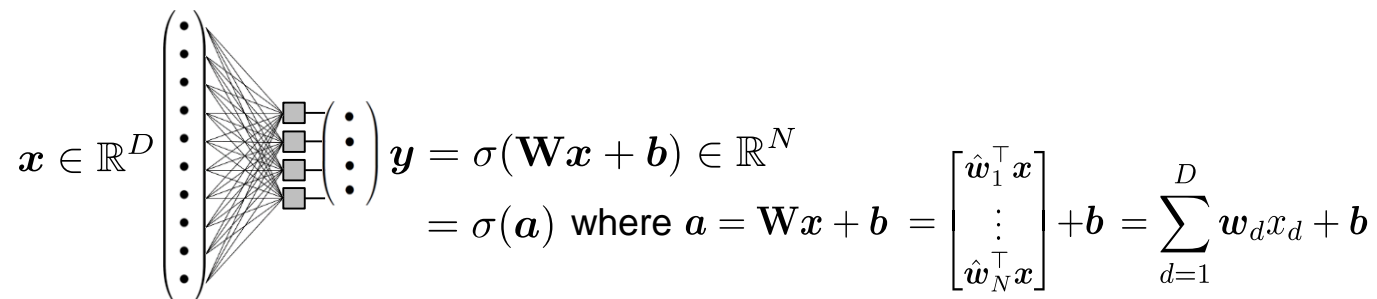


$$\begin{aligned}
 \mathbf{x} \in \mathbb{R}^D & \quad \mathbf{y} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b}) \in \mathbb{R}^N \\
 & = \sigma(\mathbf{a}) \text{ where } \mathbf{a} = \mathbf{W}\mathbf{x} + \mathbf{b} = \begin{bmatrix} \hat{\mathbf{w}}_1^\top \mathbf{x} \\ \vdots \\ \hat{\mathbf{w}}_N^\top \mathbf{x} \end{bmatrix} + \mathbf{b} = \sum_{d=1}^D \mathbf{w}_d x_d + \mathbf{b}
 \end{aligned}$$



## Exercises: *Jacobians*

**Remember:** the *simple perceptron*



- Calculate the following Jacobians:

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{a}}{\partial w_d} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} =$$

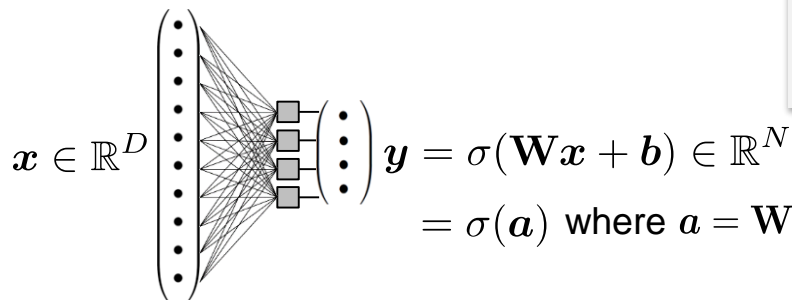
$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} \Big|_{\mathbf{b}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial w_d} \Big|_{w_d} =$$

### Exercises: *Jacobians*

**Remember:** the *simple perceptron*



$$= \sigma(\mathbf{a}) \text{ where } \mathbf{a} = \mathbf{W}\mathbf{x} + \mathbf{b} = \begin{bmatrix} \hat{w}_1^\top \mathbf{x} \\ \vdots \\ \hat{w}_N^\top \mathbf{x} \end{bmatrix} + \mathbf{b} = \sum_{d=1}^D \mathbf{w}_d x_d + \mathbf{b}$$

$$\frac{d\mathbf{y}}{d\mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_1}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_1}{\partial x_D} \Big|_{x_{D,0}} \\ \frac{\partial y_2}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_2}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_2}{\partial x_D} \Big|_{x_{D,0}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_N}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_N}{\partial x_D} \Big|_{x_{D,0}} \end{bmatrix}$$

- Calculate the following Jacobians:

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{a}}{\partial \mathbf{w}_d} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} =$$

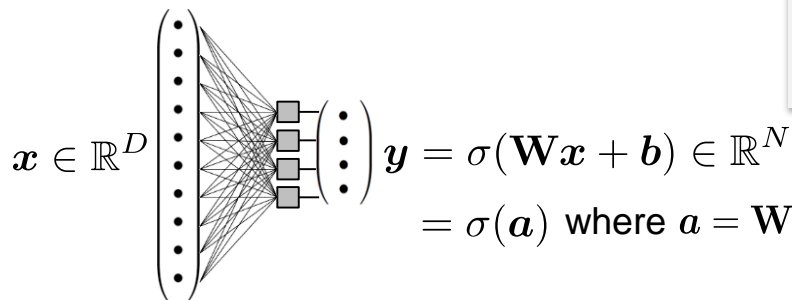
$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} \Big|_{\mathbf{b}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{w}_d} \Big|_{\mathbf{w}_d} =$$

### Exercises: *Jacobians*

**Remember:** the *simple perceptron*



$$= \sigma(\mathbf{a}) \text{ where } \mathbf{a} = \mathbf{W}\mathbf{x} + \mathbf{b} = \begin{bmatrix} \hat{w}_1^\top \mathbf{x} \\ \vdots \\ \hat{w}_N^\top \mathbf{x} \end{bmatrix} + \mathbf{b} = \sum_{d=1}^D \mathbf{w}_d x_d + \mathbf{b}$$

$$\frac{d\mathbf{y}}{d\mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_1}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_1}{\partial x_D} \Big|_{x_{D,0}} \\ \frac{\partial y_2}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_2}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_2}{\partial x_D} \Big|_{x_{D,0}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_N}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_N}{\partial x_D} \Big|_{x_{D,0}} \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{x}} f_1(\mathbf{x}_0)^\top \\ \nabla_{\mathbf{x}} f_2(\mathbf{x}_0)^\top \\ \vdots \\ \nabla_{\mathbf{x}} f_N(\mathbf{x}_0)^\top \end{bmatrix}$$

- Calculate the following Jacobians:

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{a}}{\partial \mathbf{w}_d} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} =$$

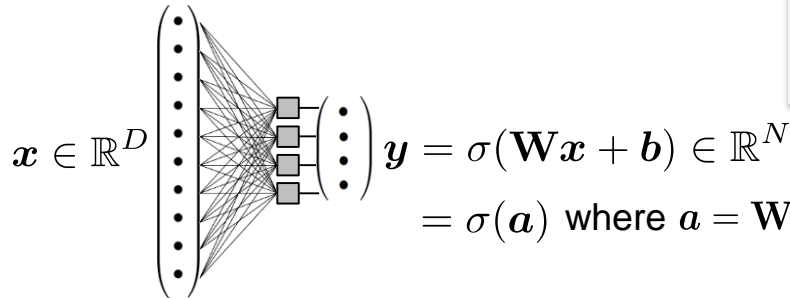
$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} \Big|_{\mathbf{b}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{w}_d} \Big|_{\mathbf{w}_d} =$$

### Exercises: *Jacobians*

**Remember:** the *simple perceptron*



$$= \sigma(\mathbf{a}) \text{ where } \mathbf{a} = \mathbf{W}\mathbf{x} + \mathbf{b} = \begin{bmatrix} \hat{w}_1^\top \mathbf{x} \\ \vdots \\ \hat{w}_N^\top \mathbf{x} \end{bmatrix} + \mathbf{b} = \sum_{d=1}^D \mathbf{w}_d x_d + \mathbf{b}$$

$$\frac{d\mathbf{y}}{d\mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_1}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_1}{\partial x_D} \Big|_{x_{D,0}} \\ \frac{\partial y_2}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_2}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_2}{\partial x_D} \Big|_{x_{D,0}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_N}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_N}{\partial x_D} \Big|_{x_{D,0}} \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{x}} f_1(\mathbf{x}_0)^\top \\ \nabla_{\mathbf{x}} f_2(\mathbf{x}_0)^\top \\ \vdots \\ \nabla_{\mathbf{x}} f_N(\mathbf{x}_0)^\top \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} \times \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0}$$

- Calculate the following Jacobians:

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \quad \frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{\mathbf{x}_0} = \quad \frac{\partial \mathbf{a}}{\partial \mathbf{w}_d} \Big|_{\mathbf{x}_0} = \quad \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} =$$

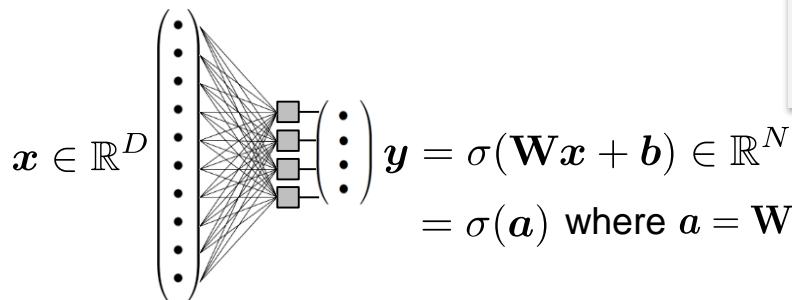
$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} \Big|_{\mathbf{b}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{w}_d} \Big|_{\mathbf{w}_d} =$$

### Exercises: *Jacobians*

**Remember: the simple perceptron**



$$= \sigma(\mathbf{a}) \text{ where } \mathbf{a} = \mathbf{W}x + \mathbf{b} = \begin{bmatrix} \hat{w}_1^\top x \\ \vdots \\ \hat{w}_N^\top x \end{bmatrix} + \mathbf{b} = \sum_{d=1}^D \mathbf{w}_d x_d + \mathbf{b}$$

$$\frac{d\mathbf{y}}{d\mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_1}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_1}{\partial x_D} \Big|_{x_{D,0}} \\ \frac{\partial y_2}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_2}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_2}{\partial x_D} \Big|_{x_{D,0}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_N}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_N}{\partial x_D} \Big|_{x_{D,0}} \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{x}} f_1(\mathbf{x}_0)^\top \\ \nabla_{\mathbf{x}} f_2(\mathbf{x}_0)^\top \\ \vdots \\ \nabla_{\mathbf{x}} f_N(\mathbf{x}_0)^\top \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} \times \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0}$$

- Calculate the following Jacobians:

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \hat{w}_1^\top \\ \vdots \\ \hat{w}_N^\top \end{bmatrix} = \mathbf{W} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{\mathbf{x}_0} = \quad \frac{\partial \mathbf{a}}{\partial \mathbf{w}_d} \Big|_{\mathbf{x}_0} = \quad \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} =$$

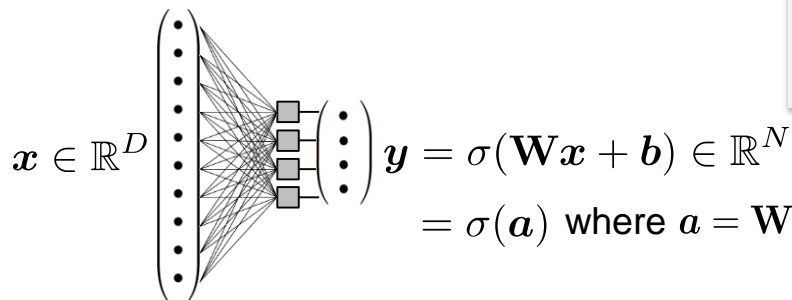
$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} \Big|_{\mathbf{b}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{w}_d} \Big|_{\mathbf{w}_d} =$$

### Exercises: *Jacobians*

**Remember: the simple perceptron**



$$= \sigma(\mathbf{a}) \text{ where } \mathbf{a} = \mathbf{W}\mathbf{x} + \mathbf{b} = \begin{bmatrix} \hat{w}_1^\top \mathbf{x} \\ \vdots \\ \hat{w}_N^\top \mathbf{x} \end{bmatrix} + \mathbf{b} = \sum_{d=1}^D \mathbf{w}_d x_d + \mathbf{b}$$

$$\frac{d\mathbf{y}}{d\mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_1}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_1}{\partial x_D} \Big|_{x_{D,0}} \\ \frac{\partial y_2}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_2}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_2}{\partial x_D} \Big|_{x_{D,0}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_N}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_N}{\partial x_D} \Big|_{x_{D,0}} \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{x}} f_1(\mathbf{x}_0)^\top \\ \nabla_{\mathbf{x}} f_2(\mathbf{x}_0)^\top \\ \vdots \\ \nabla_{\mathbf{x}} f_N(\mathbf{x}_0)^\top \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} \times \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0}$$

- Calculate the following Jacobians:

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \hat{w}_1^\top \\ \vdots \\ \hat{w}_N^\top \end{bmatrix} = \mathbf{W} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{\mathbf{x}_0} = \mathbf{I}_N \quad \frac{\partial \mathbf{a}}{\partial \mathbf{w}_d} \Big|_{\mathbf{x}_0} = \quad \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} =$$

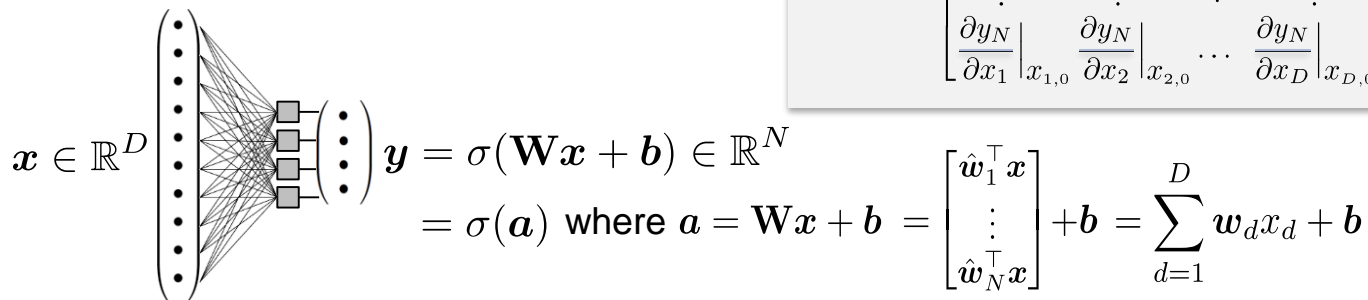
$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} \Big|_{\mathbf{b}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{w}_d} \Big|_{\mathbf{w}_d} =$$

### Exercises: *Jacobians*

**Remember: the simple perceptron**



$$\frac{d\mathbf{y}}{d\mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_1}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_1}{\partial x_D} \Big|_{x_{D,0}} \\ \frac{\partial y_2}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_2}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_2}{\partial x_D} \Big|_{x_{D,0}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_N}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_N}{\partial x_D} \Big|_{x_{D,0}} \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{x}} f_1(\mathbf{x}_0)^\top \\ \nabla_{\mathbf{x}} f_2(\mathbf{x}_0)^\top \\ \vdots \\ \nabla_{\mathbf{x}} f_N(\mathbf{x}_0)^\top \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} \times \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0}$$

- Calculate the following Jacobians:

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \hat{w}_1^\top \\ \vdots \\ \hat{w}_N^\top \end{bmatrix} = \mathbf{W} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{\mathbf{x}_0} = \mathbf{I}_N \quad \frac{\partial \mathbf{a}}{\partial w_d} \Big|_{\mathbf{x}_0} = x_d \mathbf{I}_N \quad \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} =$$

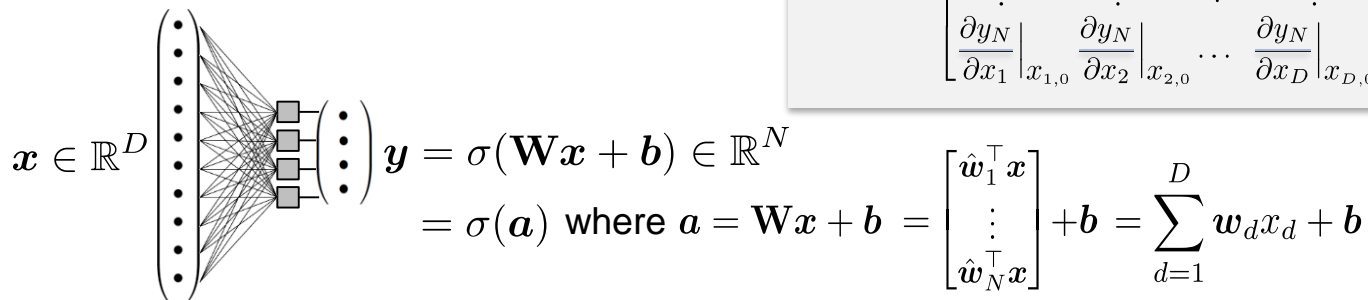
$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} \Big|_{\mathbf{b}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial w_d} \Big|_{w_d} =$$

### Exercises: *Jacobians*

**Remember: the simple perceptron**



$$\frac{d\mathbf{y}}{d\mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_1}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_1}{\partial x_D} \Big|_{x_{D,0}} \\ \frac{\partial y_2}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_2}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_2}{\partial x_D} \Big|_{x_{D,0}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_N}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_N}{\partial x_D} \Big|_{x_{D,0}} \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{x}} f_1(\mathbf{x}_0)^\top \\ \nabla_{\mathbf{x}} f_2(\mathbf{x}_0)^\top \\ \vdots \\ \nabla_{\mathbf{x}} f_N(\mathbf{x}_0)^\top \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} \times \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0}$$

- Calculate the following Jacobians:

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \hat{w}_1^\top \\ \vdots \\ \hat{w}_N^\top \end{bmatrix} = \mathbf{W} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{\mathbf{x}_0} = \mathbf{I}_N \quad \frac{\partial \mathbf{a}}{\partial w_d} \Big|_{\mathbf{x}_0} = x_d \mathbf{I}_N \quad \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} = \begin{bmatrix} \sigma'(a_{0,1}) & 0 & \cdots & 0 \\ 0 & \sigma'(a_{0,2}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma'(a_{0,N}) \end{bmatrix} = \text{diag}[\sigma'(\mathbf{a}_0)]$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} =$$

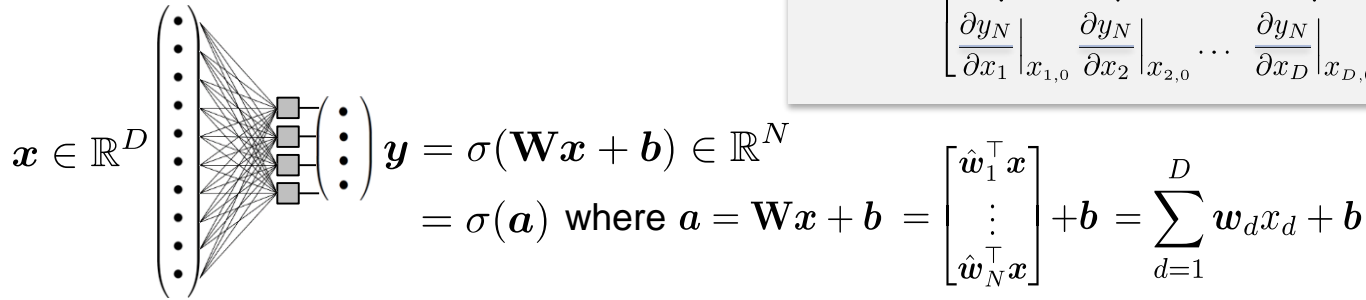
$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} \Big|_{\mathbf{b}_0} =$$

$$\frac{\partial \mathbf{y}}{\partial w_d} \Big|_{w_d} =$$



### Exercises: *Jacobians*

**Remember: the simple perceptron**



$$\frac{d\mathbf{y}}{d\mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_1}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_1}{\partial x_D} \Big|_{x_{D,0}} \\ \frac{\partial y_2}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_2}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_2}{\partial x_D} \Big|_{x_{D,0}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_N}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_N}{\partial x_D} \Big|_{x_{D,0}} \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{x}} f_1(\mathbf{x}_0)^\top \\ \nabla_{\mathbf{x}} f_2(\mathbf{x}_0)^\top \\ \vdots \\ \nabla_{\mathbf{x}} f_N(\mathbf{x}_0)^\top \end{bmatrix}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} \times \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0}$$

- Calculate the following Jacobians:

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \hat{\mathbf{w}}_1^\top \\ \vdots \\ \hat{\mathbf{w}}_N^\top \end{bmatrix} = \mathbf{W} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{\mathbf{x}_0} = \mathbf{I}_N \quad \frac{\partial \mathbf{a}}{\partial \mathbf{w}_d} \Big|_{\mathbf{x}_0} = x_d \mathbf{I}_N \quad \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} = \begin{bmatrix} \sigma'(a_{0,1}) & 0 & \cdots & 0 \\ 0 & \sigma'(a_{0,2}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma'(a_{0,N}) \end{bmatrix} = \text{diag}[\sigma'(\mathbf{a}_0)]$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} \times \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \text{diag}[\sigma'(\mathbf{W}\mathbf{x}_0 + \mathbf{b})] \mathbf{W}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} \Big|_{\mathbf{b}_0} =$$

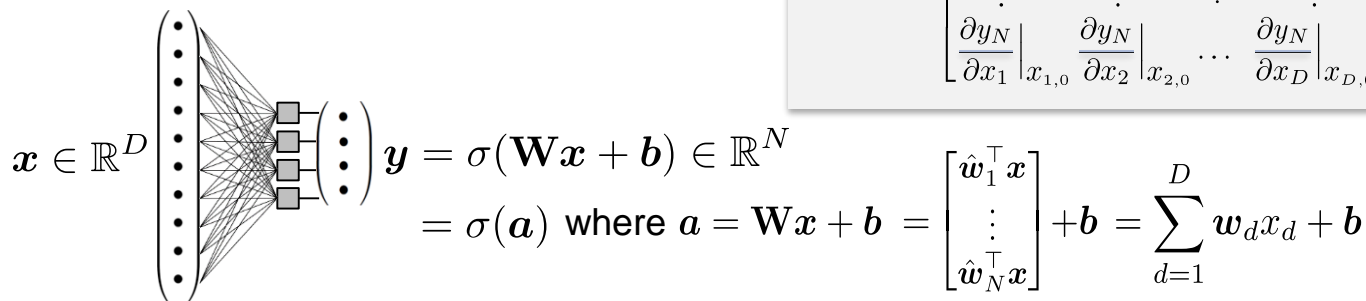
$$\frac{\partial \mathbf{y}}{\partial \mathbf{w}_d} \Big|_{\mathbf{w}_d} =$$

### Exercises: *Jacobians*

**Remember: the simple perceptron**

$$\frac{dy}{dx} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_1}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_1}{\partial x_D} \Big|_{x_{D,0}} \\ \frac{\partial y_2}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_2}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_2}{\partial x_D} \Big|_{x_{D,0}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_N}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_N}{\partial x_D} \Big|_{x_{D,0}} \end{bmatrix} = \begin{bmatrix} \nabla_{\mathbf{x}} f_1(\mathbf{x}_0)^\top \\ \nabla_{\mathbf{x}} f_2(\mathbf{x}_0)^\top \\ \vdots \\ \nabla_{\mathbf{x}} f_N(\mathbf{x}_0)^\top \end{bmatrix}$$

$$\frac{\partial y}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \frac{\partial y}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} \times \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0}$$



- Calculate the following Jacobians:

$$\frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \begin{bmatrix} \hat{\mathbf{w}}_1^\top \\ \vdots \\ \hat{\mathbf{w}}_N^\top \end{bmatrix} = \mathbf{W} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{\mathbf{x}_0} = \mathbf{I}_N \quad \frac{\partial \mathbf{a}}{\partial w_d} \Big|_{\mathbf{x}_0} = x_d \mathbf{I}_N \quad \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} = \begin{bmatrix} \sigma'(a_{0,1}) & 0 & \cdots & 0 \\ 0 & \sigma'(a_{0,2}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma'(a_{0,N}) \end{bmatrix} = \text{diag}[\sigma'(\mathbf{a}_0)]$$

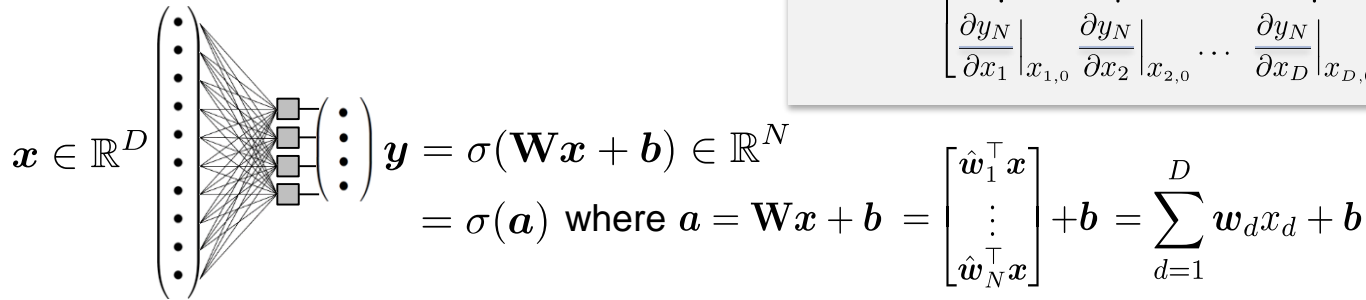
$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{W}\mathbf{x}_0 + \mathbf{b}} \times \frac{\partial \mathbf{a}}{\partial \mathbf{x}} \Big|_{\mathbf{x}_0} = \text{diag}[\sigma'(\mathbf{W}\mathbf{x}_0 + \mathbf{b})] \mathbf{W}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{b}} \Big|_{\mathbf{b}_0} = \frac{\partial \mathbf{y}}{\partial \mathbf{a}} \Big|_{\mathbf{W}\mathbf{x} + \mathbf{b}_0} \times \frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{\mathbf{b}_0} = \text{diag}[\sigma'(\mathbf{W}\mathbf{x} + \mathbf{b}_0)]$$

$$\frac{\partial \mathbf{y}}{\partial w_d} \Big|_{w_d} =$$

### Exercises: *Jacobians*

**Remember: the simple perceptron**



$$\frac{dy}{dx} \Big|_{x_0} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_1}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_1}{\partial x_D} \Big|_{x_{D,0}} \\ \frac{\partial y_2}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_2}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_2}{\partial x_D} \Big|_{x_{D,0}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_N}{\partial x_1} \Big|_{x_{1,0}} & \frac{\partial y_N}{\partial x_2} \Big|_{x_{2,0}} & \cdots & \frac{\partial y_N}{\partial x_D} \Big|_{x_{D,0}} \end{bmatrix} = \begin{bmatrix} \nabla_x f_1(\mathbf{x}_0)^\top \\ \nabla_x f_2(\mathbf{x}_0)^\top \\ \vdots \\ \nabla_x f_N(\mathbf{x}_0)^\top \end{bmatrix}$$

$$\frac{\partial y}{\partial x} \Big|_{x_0} = \frac{\partial y}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} \times \frac{\partial \mathbf{a}}{\partial x} \Big|_{x_0}$$

- Calculate the following Jacobians:

$$\frac{\partial \mathbf{a}}{\partial x} \Big|_{x_0} = \begin{bmatrix} \hat{w}_1^T \\ \vdots \\ \hat{w}_N^T \end{bmatrix} = \mathbf{W} \quad \frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{x_0} = \mathbf{I}_N \quad \frac{\partial \mathbf{a}}{\partial w_d} \Big|_{x_0} = x_d \mathbf{I}_N \quad \frac{\partial y}{\partial \mathbf{a}} \Big|_{\mathbf{a}_0} = \begin{bmatrix} \sigma'(a_{0,1}) & 0 & \cdots & 0 \\ 0 & \sigma'(a_{0,2}) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma'(a_{0,N}) \end{bmatrix} = \text{diag}[\sigma'(\mathbf{a}_0)]$$

$$\frac{\partial y}{\partial x} \Big|_{x_0} = \frac{\partial y}{\partial \mathbf{a}} \Big|_{\mathbf{W}x_0 + \mathbf{b}} \times \frac{\partial \mathbf{a}}{\partial x} \Big|_{x_0} = \text{diag}[\sigma'(\mathbf{W}x_0 + \mathbf{b})] \mathbf{W}$$

$$\frac{\partial y}{\partial \mathbf{b}} \Big|_{b_0} = \frac{\partial y}{\partial \mathbf{a}} \Big|_{\mathbf{W}x + \mathbf{b}_0} \times \frac{\partial \mathbf{a}}{\partial \mathbf{b}} \Big|_{b_0} = \text{diag}[\sigma'(\mathbf{W}x + \mathbf{b}_0)]$$

$$\frac{\partial y}{\partial w_d} \Big|_{w_d} = \frac{\partial y}{\partial \mathbf{a}} \Big|_{\mathbf{W}_d x + \mathbf{b}} \times \frac{\partial \mathbf{a}}{\partial w_d} \Big|_{w_d} = \text{diag}[\sigma'(\mathbf{W}_d x + \mathbf{b})] x_d$$

## Exercise: *Least Squares*

- Calculate the gradient of  $f(\boldsymbol{\theta}) = \|\mathbf{W}\boldsymbol{\theta} - \mathbf{y}\|_2^2$  where  $\mathbf{W} \in \mathbb{R}^{N \times D}$

## Exercise: *Least Squares*

- Calculate the gradient of  $f(\boldsymbol{\theta}) = \|\mathbf{W}\boldsymbol{\theta} - \mathbf{y}\|_2^2$  where  $\mathbf{W} \in \mathbb{R}^{N \times D}$

**Hint:** Let  $\mathbf{a} = \mathbf{W}\boldsymbol{\theta} - \mathbf{y}$ ,  $g(\mathbf{a}) = \|\mathbf{a}\|_2^2$  and  $r = g(\mathbf{a}) = f(\boldsymbol{\theta})$ .

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0)^\top = \left. \frac{\partial r}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_0} = ?$$

## Exercise: *Least Squares*

- Calculate the gradient of  $f(\boldsymbol{\theta}) = \|\mathbf{W}\boldsymbol{\theta} - \mathbf{y}\|_2^2$  where  $\mathbf{W} \in \mathbb{R}^{N \times D}$

**Hint:** Let  $\mathbf{a} = \mathbf{W}\boldsymbol{\theta} - \mathbf{y}$ ,  $g(\mathbf{a}) = \|\mathbf{a}\|_2^2$  and  $r = g(\mathbf{a}) = f(\boldsymbol{\theta})$ .

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0)^\top = \left. \frac{\partial r}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_0} = \left. \frac{\partial r}{\partial \mathbf{a}} \right|_{\mathbf{W}\boldsymbol{\theta}_0 - \mathbf{y}} \times \left. \frac{\partial \mathbf{a}}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_0}$$

## Exercise: *Least Squares*

- Calculate the gradient of  $f(\boldsymbol{\theta}) = \|\mathbf{W}\boldsymbol{\theta} - \mathbf{y}\|_2^2$  where  $\mathbf{W} \in \mathbb{R}^{N \times D}$

**Hint:** Let  $\mathbf{a} = \mathbf{W}\boldsymbol{\theta} - \mathbf{y}$ ,  $g(\mathbf{a}) = \|\mathbf{a}\|_2^2$  and  $r = g(\mathbf{a}) = f(\boldsymbol{\theta})$ .

$$\begin{aligned}\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0)^\top &= \left. \frac{\partial r}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_0} = \left. \frac{\partial r}{\partial \mathbf{a}} \right|_{\mathbf{W}\boldsymbol{\theta}_0 - \mathbf{y}} \times \left. \frac{\partial \mathbf{a}}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_0} \\ &= \nabla_{\mathbf{a}} g(\mathbf{W}\boldsymbol{\theta}_0 - \mathbf{y})^\top \times \mathbf{W}\end{aligned}$$

## Exercise: *Least Squares*

- Calculate the gradient of  $f(\boldsymbol{\theta}) = \|\mathbf{W}\boldsymbol{\theta} - \mathbf{y}\|_2^2$  where  $\mathbf{W} \in \mathbb{R}^{N \times D}$

**Hint:** Let  $\mathbf{a} = \mathbf{W}\boldsymbol{\theta} - \mathbf{y}$ ,  $g(\mathbf{a}) = \|\mathbf{a}\|_2^2$  and  $r = g(\mathbf{a}) = f(\boldsymbol{\theta})$ .

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0)^\top = \left. \frac{\partial r}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_0} = \left. \frac{\partial r}{\partial \mathbf{a}} \right|_{\mathbf{W}\boldsymbol{\theta}_0 - \mathbf{y}} \times \left. \frac{\partial \mathbf{a}}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_0}$$

$$= \nabla_{\mathbf{a}} g(\mathbf{W}\boldsymbol{\theta}_0 - \mathbf{y})^\top \times \mathbf{W}$$

$$\text{We have: } \nabla_{\mathbf{a}} g(\mathbf{a}_0) = 2\mathbf{a}_0$$



## Exercise: *Least Squares*

- Calculate the gradient of  $f(\boldsymbol{\theta}) = \|\mathbf{W}\boldsymbol{\theta} - \mathbf{y}\|_2^2$  where  $\mathbf{W} \in \mathbb{R}^{N \times D}$

**Hint:** Let  $\mathbf{a} = \mathbf{W}\boldsymbol{\theta} - \mathbf{y}$ ,  $g(\mathbf{a}) = \|\mathbf{a}\|_2^2$  and  $r = g(\mathbf{a}) = f(\boldsymbol{\theta})$ .

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0)^\top &= \left. \frac{\partial r}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_0} = \left. \frac{\partial r}{\partial \mathbf{a}} \right|_{\mathbf{W}\boldsymbol{\theta}_0 - \mathbf{y}} \times \left. \frac{\partial \mathbf{a}}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_0} \\ &= \nabla_{\mathbf{a}} g(\mathbf{W}\boldsymbol{\theta}_0 - \mathbf{y})^\top \times \mathbf{W} && \text{We have: } \nabla_{\mathbf{a}} g(\mathbf{a}_0) = 2\mathbf{a}_0 \\ &= 2(\mathbf{W}\boldsymbol{\theta}_0 - \mathbf{y})^\top \mathbf{W} \end{aligned}$$

## Exercise: *Least Squares*

- Calculate the gradient of  $f(\boldsymbol{\theta}) = \|\mathbf{W}\boldsymbol{\theta} - \mathbf{y}\|_2^2$  where  $\mathbf{W} \in \mathbb{R}^{N \times D}$

**Hint:** Let  $\mathbf{a} = \mathbf{W}\boldsymbol{\theta} - \mathbf{y}$ ,  $g(\mathbf{a}) = \|\mathbf{a}\|_2^2$  and  $r = g(\mathbf{a}) = f(\boldsymbol{\theta})$ .

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0)^\top &= \left. \frac{\partial r}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_0} = \left. \frac{\partial r}{\partial \mathbf{a}} \right|_{\mathbf{W}\boldsymbol{\theta}_0 - \mathbf{y}} \times \left. \frac{\partial \mathbf{a}}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}_0} \\ &= \nabla_{\mathbf{a}} g(\mathbf{W}\boldsymbol{\theta}_0 - \mathbf{y})^\top \times \mathbf{W} && \text{We have: } \nabla_{\mathbf{a}} g(\mathbf{a}_0) = 2\mathbf{a}_0 \\ &= 2(\mathbf{W}\boldsymbol{\theta}_0 - \mathbf{y})^\top \mathbf{W} \end{aligned}$$

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta}_0) = 2\mathbf{W}^\top (\mathbf{W}\boldsymbol{\theta}_0 - \mathbf{y})$$