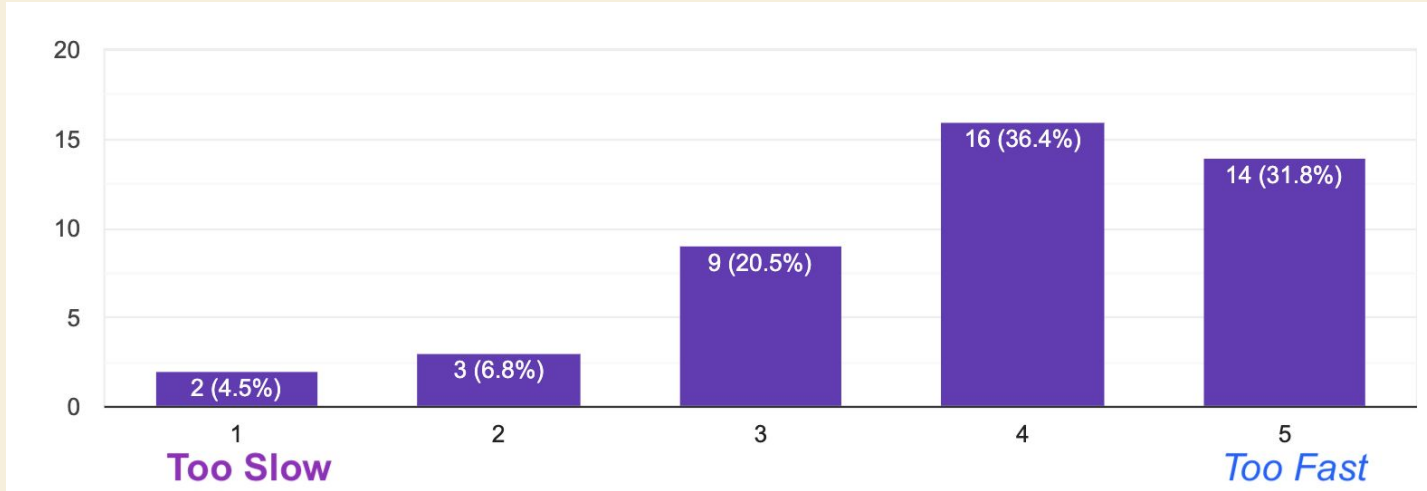


# Optimization and Automatic Differentiation

Mathematics for Machine Learning

Lecturer: Matthew Wicker

# Logistics!



Thank you all so much for your feedback! It is invaluable to me and I take it seriously.

Some say I use terms they don't know, so that will be the question after today's lecture!

# Errata

So there were some errors especially in the last section of last week's notes. So we are going over that again today to clear things up.

Thank you so much to all who reported errata!!

New reward system for this! If you email me and point out a valid (yet unreported) mathematical error in the lecture notes, you get 1 ticket in the end of term raffle!

Prizes will include: A Imperial hoodie, a Raspberry Pi Model B, & Amazon Gift Card

# Defining some terminology

$$\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}$$

column  
vectors  
by  
default

$$\mathbf{x}, \theta \in \mathbb{R}^{n \times 1}$$

## Defining some terminology

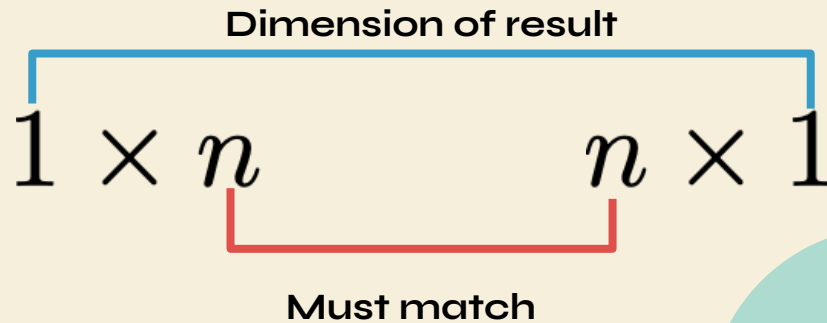
$$\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}, \quad \theta^\top = [\theta_1, \dots, \theta_n]$$

$$n \times 1$$

$$1 \times n$$

# Defining some terminology

$$\theta^{\top} = [\theta_1, \dots, \theta_n] \quad \mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}$$



# Writing out our model

$$\Theta^T X : 1 \times n \times n /$$

$$\begin{bmatrix} (- \dots) \\ (- \dots) \\ (- \dots) \end{bmatrix}$$

$$1 \times n$$

$$\begin{bmatrix} \mathbf{x}^T(1), \\ \mathbf{x}^T(2), \\ \vdots, \\ \mathbf{x}^T(N) \end{bmatrix}$$

$$\begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$n \times 1$$

# Writing out our model

$X\theta: N \times n \times n$

$$\begin{bmatrix} \hat{y}^{(1)} \\ \hat{y}^{(2)} \\ \dots \\ \hat{y}^{(N)} \end{bmatrix} = \begin{bmatrix} x_1^{(1)} & \dots & x_n^{(1)} \\ x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \ddots & \vdots \\ x_1^{(N)} & \dots & x_n^{(N)} \end{bmatrix} \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

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

$x_1^{(1)}\theta_1 + x_2^{(1)}\theta_2 + \dots$   
 $x_1^{(2)}\theta_1 + x_2^{(2)}\theta_2 + \dots$   
 $\vdots$   
 $\vdots$   
 $\vdots$



# Writing out our model

The diagram illustrates the matrix equation  $\hat{\mathbf{y}} = \mathbf{X}\boldsymbol{\theta}$ . The equation is enclosed in a red hand-drawn box. Three teal arrows point from the dimensions below to the terms in the equation: from  $N \times 1$  to  $\hat{\mathbf{y}}$ , from  $N \times n$  to  $\mathbf{X}$ , and from  $n \times 1$  to  $\boldsymbol{\theta}$ .

$$\hat{\mathbf{y}} = \mathbf{X}\boldsymbol{\theta}$$

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

# Writing out our model

$$\hat{\mathbf{y}} = \mathbf{X}\theta$$

*actual - prediction*

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{N} \sum_i (\mathbf{y}_i - \hat{\mathbf{y}}_i)^2$$

$$\begin{bmatrix} - \\ - \\ - \\ - \end{bmatrix} \begin{bmatrix} - \\ - \\ - \\ - \end{bmatrix}$$

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{N} \sum_i (\mathbf{y}_i - \mathbf{X}_{i,j} \theta_j)^2$$

*summing over  $j$  implicitly*

## Writing out our model

$$\hat{\mathbf{y}} = \mathbf{X}\boldsymbol{\theta}$$

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{N} \sum_i (\mathbf{y}_i - \mathbf{X}_{i,j} \theta_j)^2$$

$$\frac{1}{N} (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})^\top (\mathbf{X}\boldsymbol{\theta} - \mathbf{y})$$

## Expanding the loss

$$\begin{aligned} \frac{1}{N} (\mathbf{X}\theta - \mathbf{y})^\top (\mathbf{X}\theta - \mathbf{y}) &= \frac{1}{N} ((\mathbf{X}\theta)^\top - \mathbf{y}^\top) (\mathbf{X}\theta - \mathbf{y}) \\ &= \frac{1}{N} (\theta^\top \mathbf{X}^\top - \mathbf{y}^\top) (\mathbf{X}\theta - \mathbf{y}) = \frac{1}{N} (\theta^\top \mathbf{X}^\top \mathbf{X} \theta - \theta^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \theta + \mathbf{y}^\top \mathbf{y}) \\ &= \frac{1}{N} \left( \theta^\top (\mathbf{X}^\top \mathbf{X}) \theta - \theta^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} \theta + \mathbf{y}^\top \mathbf{y} \right) \end{aligned}$$

# Expanding the loss

$$\frac{1}{N} (\mathbf{X}\theta - \mathbf{y})^\top (\mathbf{X}\theta - \mathbf{y})$$

$$\underset{N \times 1}{\mathbf{y}} = \underset{N \times m}{\mathbf{X}} \underset{m \times 1}{\theta}$$

$$= \left( \theta^\top (\mathbf{X}^\top \mathbf{X}) \theta - \overset{(\mathbf{X}\theta)^\top \mathbf{y}}{\underset{\substack{(k \times m) \quad (m \times N) \\ = 1 \times 1}}{\theta^\top \mathbf{X}^\top \mathbf{y}}} - \underset{\substack{(1 \times N) \quad (N \times m) \quad (m \times 1) \\ = 1 \times 1}}{\mathbf{y}^\top \mathbf{X} \theta} + \mathbf{y}^\top \mathbf{y} \right)$$

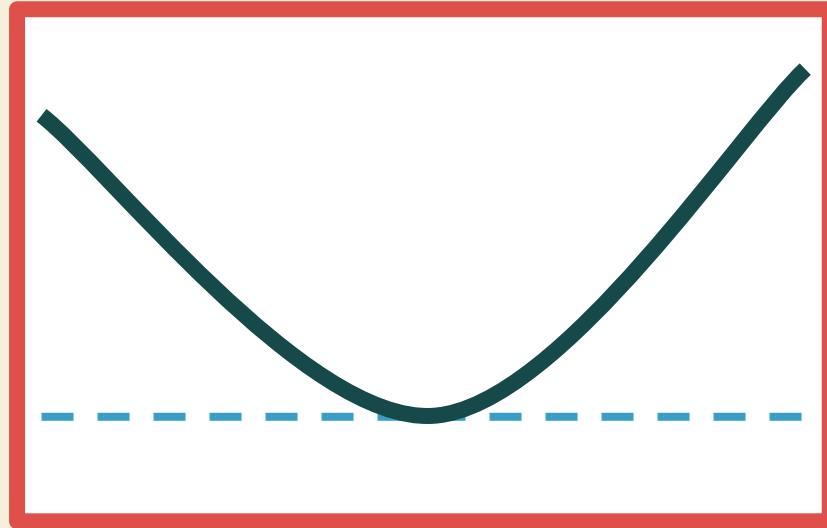
↖ scalar ↗

$$= \left( \theta^\top (\mathbf{X}^\top \mathbf{X}) \theta - 2(\underset{\substack{(1 \times N) \quad (N \times m) \quad (m \times 1) \\ = 1 \times 1}}{\mathbf{y}^\top \mathbf{X} \theta}) + \mathbf{y}^\top \mathbf{y} \right)$$

Why is this step true? Take a second to convince yourself that this is correct

# Using Calculus in Optimization

For a quadratic functions (formally defined in lecture 4), we know that the optimal function value will have gradient zero.



# Using Calculus in Optimization

For a quadratic functions (formally defined in lecture 4), we know that the optimal function value will have gradient zero.

$$\left( \theta^\top (\mathbf{X}^\top \mathbf{X}) \theta - 2(\mathbf{y}^\top \mathbf{X} \theta) + \mathbf{y}^\top \mathbf{y} \right)$$

We don't want to differentiate all of these by hand, so we will write out (and here prove) a few differentiation rules

# Proving some vector calc identities

$$\nabla_{\theta} \mathbf{c}^{\top} \theta = \mathbf{c}^{\top}$$

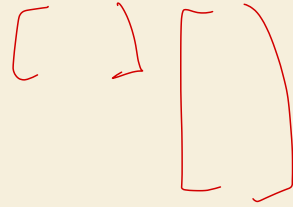
*vector columns  
by default*



# Proving some vector calc identities

Start by expanding into  
einstein or summation  
notation

$$\mathbf{c}^\top \boldsymbol{\theta} = \sum_j \mathbf{c}_j \theta_j$$



# Proving some vector calc identities

Start by expanding into  
einstein or summation  
notation

$$\mathbf{c}^\top \theta = \sum_j \mathbf{c}_j \theta_j$$

Now reason about what  
each of the partial  
derivatives of the einstein  
notation are

$$\frac{\partial \mathbf{c}^\top \theta}{\partial \theta_j} = \mathbf{c}_j$$

# Proving some vector calc identities

Start by expanding into einstein or summation notation

$$\mathbf{c}^\top \theta = \sum_j \mathbf{c}_j \theta_j$$

Now reason about what each of the partial derivatives of the einstein notation are

$$\frac{\partial \mathbf{c}^\top \theta}{\partial \theta_j} = \mathbf{c}_j$$

Finally put this from einstein notation back to vector notation

$$\nabla_\theta \mathbf{c}^\top \theta = \mathbf{c}^\top$$

# Proving some vector calc identities

$$\nabla_{\theta}(\overbrace{\theta^{\top} \mathbf{A} \theta}^{\text{scalar}}) = \mathbf{A} \theta + \mathbf{A}^{\top} \theta$$

$$\begin{matrix} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \\ 1 \times n \end{matrix} \begin{matrix} \begin{bmatrix} \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix} \\ n \times m \end{matrix} \begin{matrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \\ m \times 1 \end{matrix}$$

$$\text{row } i: \quad \text{col}(A) \times \theta$$

# Proving some vector calc identities

Start by expanding into  
einstein or summation  
notation

$$\theta^\top \mathbf{A} \theta = \sum_i \sum_j \theta_i \theta_j \mathbf{A}_{i,j}$$

*order doesn't matter*

# Proving some vector calc identities

Start by expanding into einstein or summation notation

$$\theta^\top \mathbf{A} \theta = \sum_i \sum_j \theta_i \theta_j \mathbf{A}_{i,j}$$

Now reason about what each of the partial derivatives of the einstein notation are

$$\frac{\partial \theta^\top \mathbf{A} \theta}{\partial \theta_k} = \sum_i \theta_i \mathbf{A}_{i,k} + \sum_j \mathbf{A}_{k,j} \theta_j$$

Handwritten notes in red:

- For the first term  $\sum_i \theta_i \mathbf{A}_{i,k}$ , an arrow points from  $k=j$  to the index  $k$  in  $\mathbf{A}_{i,k}$ .
- For the second term  $\sum_j \mathbf{A}_{k,j} \theta_j$ , an arrow points from  $k=i$  to the index  $k$  in  $\mathbf{A}_{k,j}$ .
- A handwritten definition for the Kronecker delta is shown:  $\frac{\partial (\theta_i \theta_j)}{\partial \theta_k} = \begin{cases} \theta_j & , k=i \\ \theta_i & , k=j \\ 0 & \text{else} \end{cases}$

# Proving some vector calc identities

Start by expanding into einstein or summation notation

$$\theta^\top \mathbf{A} \theta = \sum_i \sum_j \theta_i \theta_j \mathbf{A}_{i,j}$$

Now reason about what each of the partial derivatives of the einstein notation are

$$\frac{\partial \theta^\top \mathbf{A} \theta}{\partial \theta_k} = \boxed{\sum_i \theta_i \mathbf{A}_{i,k}} + \boxed{\sum_j \mathbf{A}_{k,j} \theta_j}$$

Finally put this from einstein notation back to vector notation

# Proving some vector calc identities

Start by expanding into  
einstein or summation  
notation

$$\theta^\top \mathbf{A} \theta = \sum_i \sum_j \theta_i \theta_j \mathbf{A}_{i,j}$$

Now reason about what  
each of the partial  
derivatives of the einstein  
notation are

$$\frac{\partial \theta^\top \mathbf{A} \theta}{\partial \theta_k} = \sum_i \underbrace{\theta_i \mathbf{A}_{i,k}}_{k^{\text{th}} \text{ element of } \mathbf{A} \theta} + \sum_j \underbrace{\mathbf{A}_{k,j} \theta_j}_{k^{\text{th}} \text{ element of } \mathbf{A}^\top \theta}$$

Finally put this from  
einstein notation back to  
vector notation

$$\nabla_\theta (\theta^\top \mathbf{A} \theta) = \mathbf{A} \theta + \mathbf{A}^\top \theta$$

(if  $\mathbf{A}$  is symmetric then it's  $2\mathbf{A}\theta$ )  
:



# Optimizing our parameter

$$\nabla_{\theta}(c^T \theta) = c^T$$

$$\nabla_{\theta}(\theta^T A \theta) = \theta^T (A + A^T)$$

$$(y - Ax)^T (y - Ax)$$

$$y^T A^T A x - y^T A x - x^T A^T y$$

$$\mathcal{L} = \frac{1}{N} (\mathbf{X}\theta - \mathbf{y})^T (\mathbf{X}\theta - \mathbf{y})$$

$$\theta^T A^T x \theta - 2y^T x \theta + y^T y \xrightarrow{\frac{\partial}{\partial \theta}} \theta^T (x A^T + (A x)^T) - 2y^T x = 0$$

$$: 2\theta^T A^T x - 2y^T x = 0$$

$$\theta^T = y^T A (A^T x)^{-1}$$

$$\theta = (x^T x)^{-1} x^T y$$

$$\nabla_{\theta} \mathcal{L} = \frac{2}{N} \left( (\mathbf{X}^T \mathbf{X}) \theta - \mathbf{X}^T \mathbf{y} \right)$$

$$(A x)^T = x^T A^T$$

$$\mathcal{L} = \frac{1}{N} \left( \theta^T (\mathbf{X}^T \mathbf{X}) \theta - 2(\mathbf{y}^T \mathbf{X} \theta) + \mathbf{y}^T \mathbf{y} \right)$$

$$\begin{bmatrix} 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$A = x x^T$$

$$\theta^T A \theta$$

$$- 2 \underbrace{(x^T y)^T}_{= c^T \theta} \theta + y^T y$$

$$\nabla_{\theta} \mathcal{L} = A \theta + A^T \theta - 2 x y = x^T x \theta + x^T x \theta - 2 x^T y = 2 x^T x \theta - 2 x^T y$$

## Solving for the parameter

set grad of loss to 0:

$$0 = \frac{2}{N} \left( (\mathbf{X}^\top \mathbf{X})\theta - \mathbf{X}^\top \mathbf{y} \right)$$

$$= (\mathbf{X}^\top \mathbf{X})\theta - \mathbf{X}^\top \mathbf{y}$$

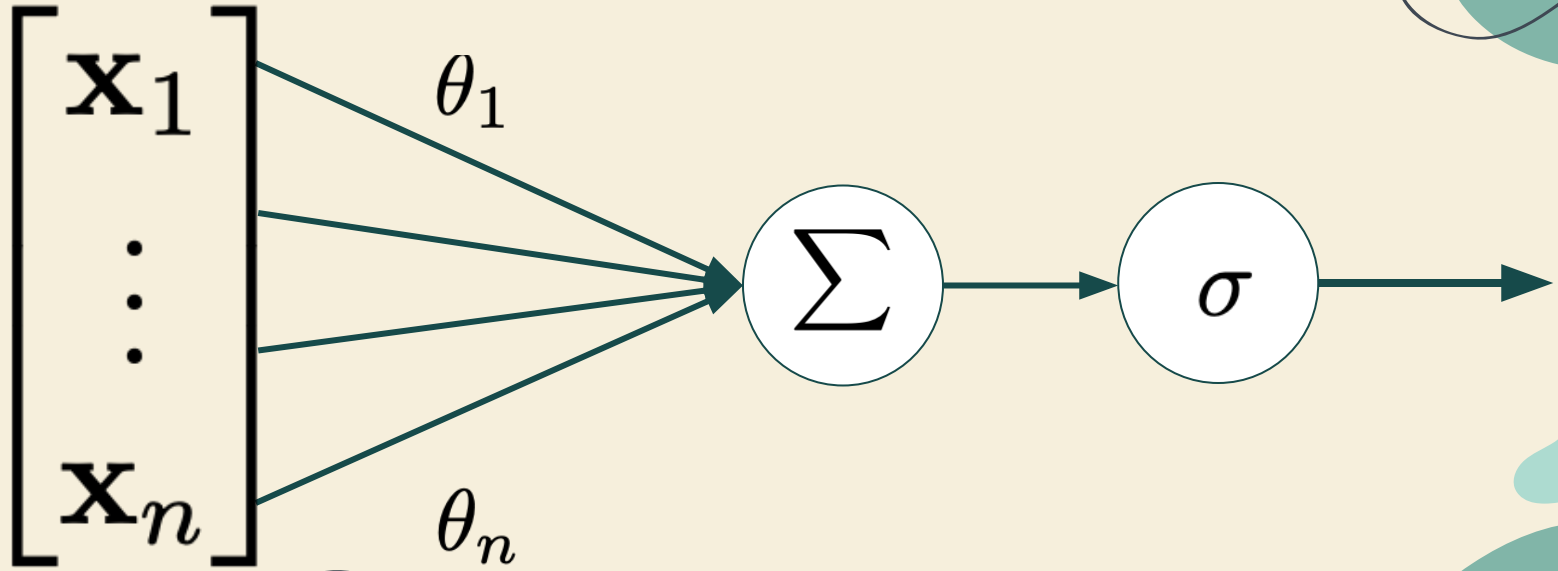
$$(\mathbf{X}^\top \mathbf{X})\theta = \mathbf{X}^\top \mathbf{y}$$

$$\theta = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

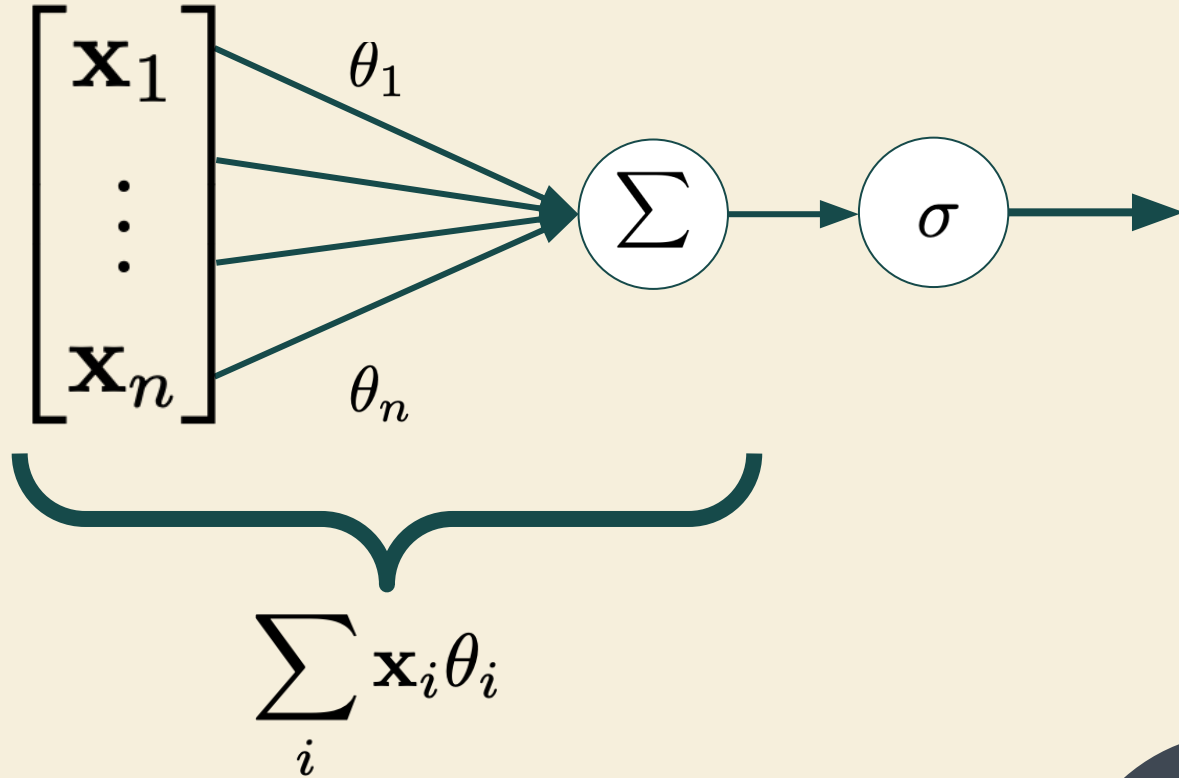
# Learning Objectives

- Understand the history of deep learning
- Formulation of the multi-layer perceptron
- Deep Networks and the need for auto-diff
- Forward-mode automatic differentiation
- Reverse-mode automatic differentiation

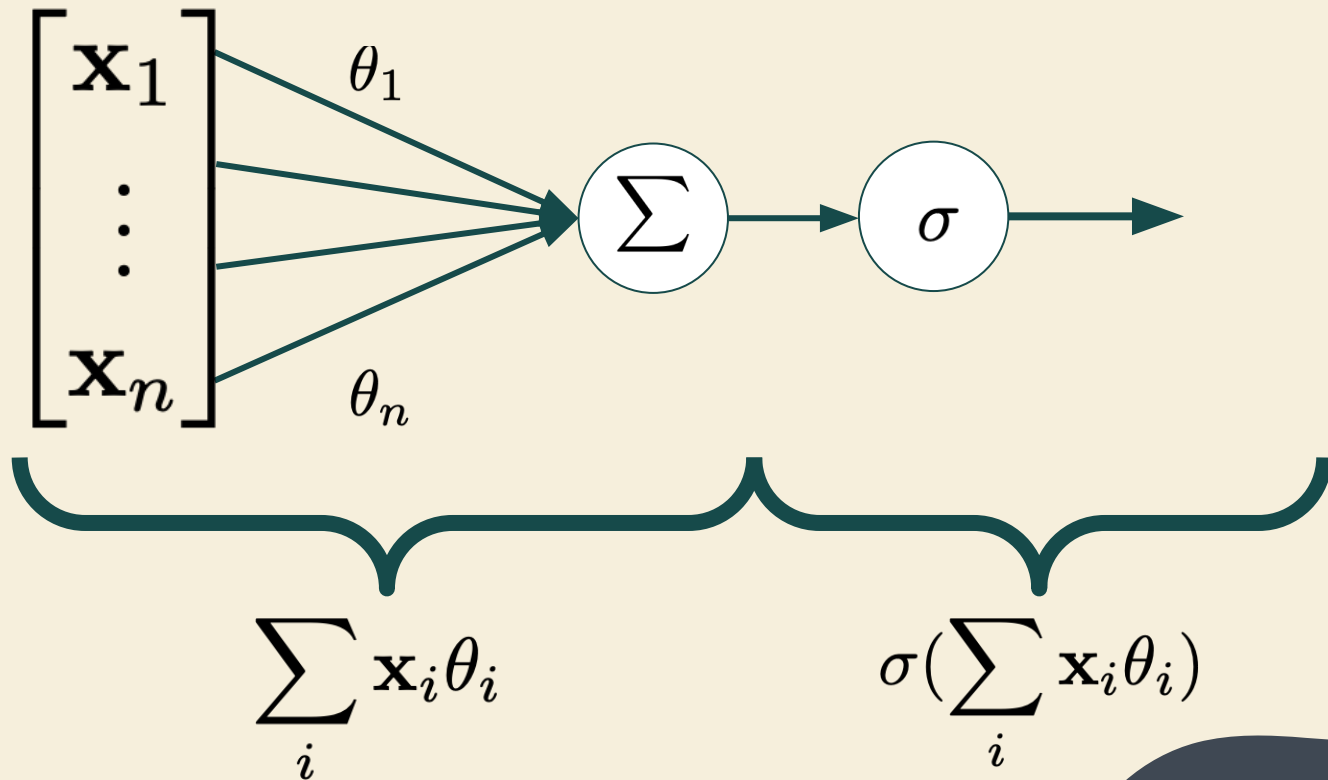
# Formulating the multilayer perceptron



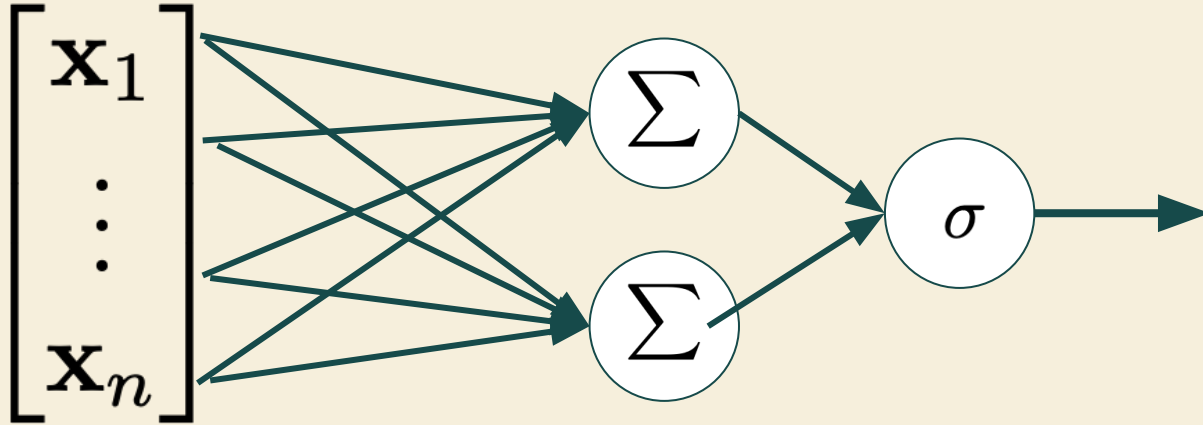
# Formulating the MLP



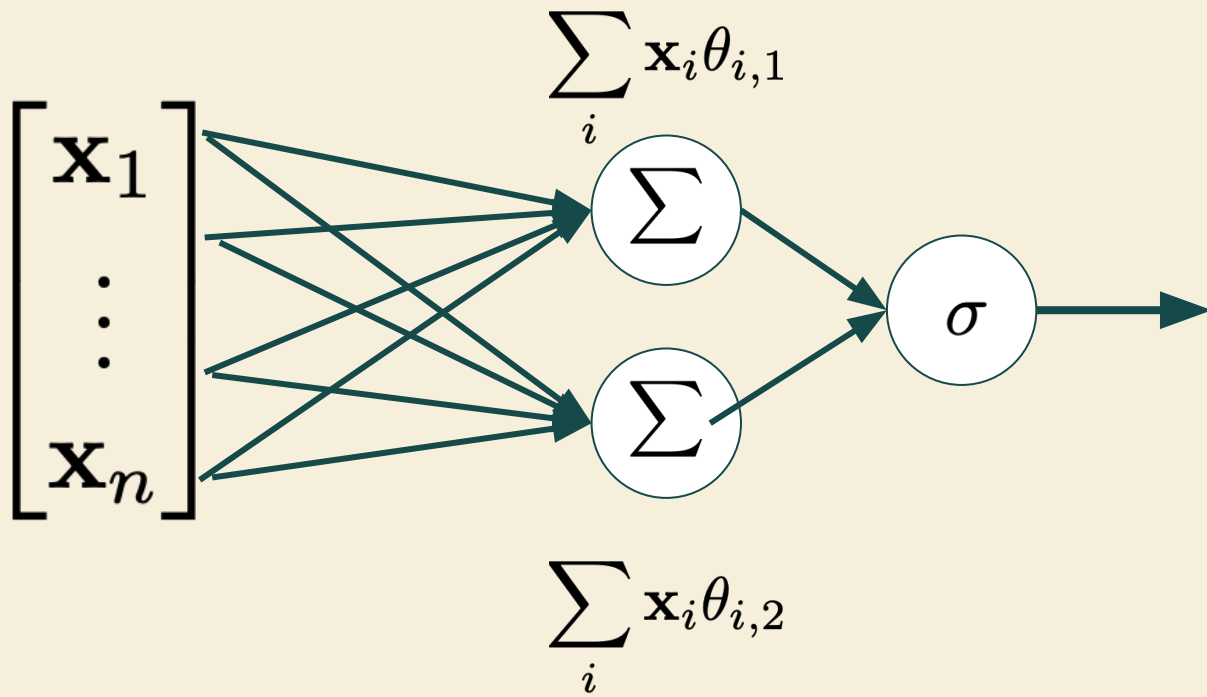
# Formulating the MLP



# Formulating the MLP

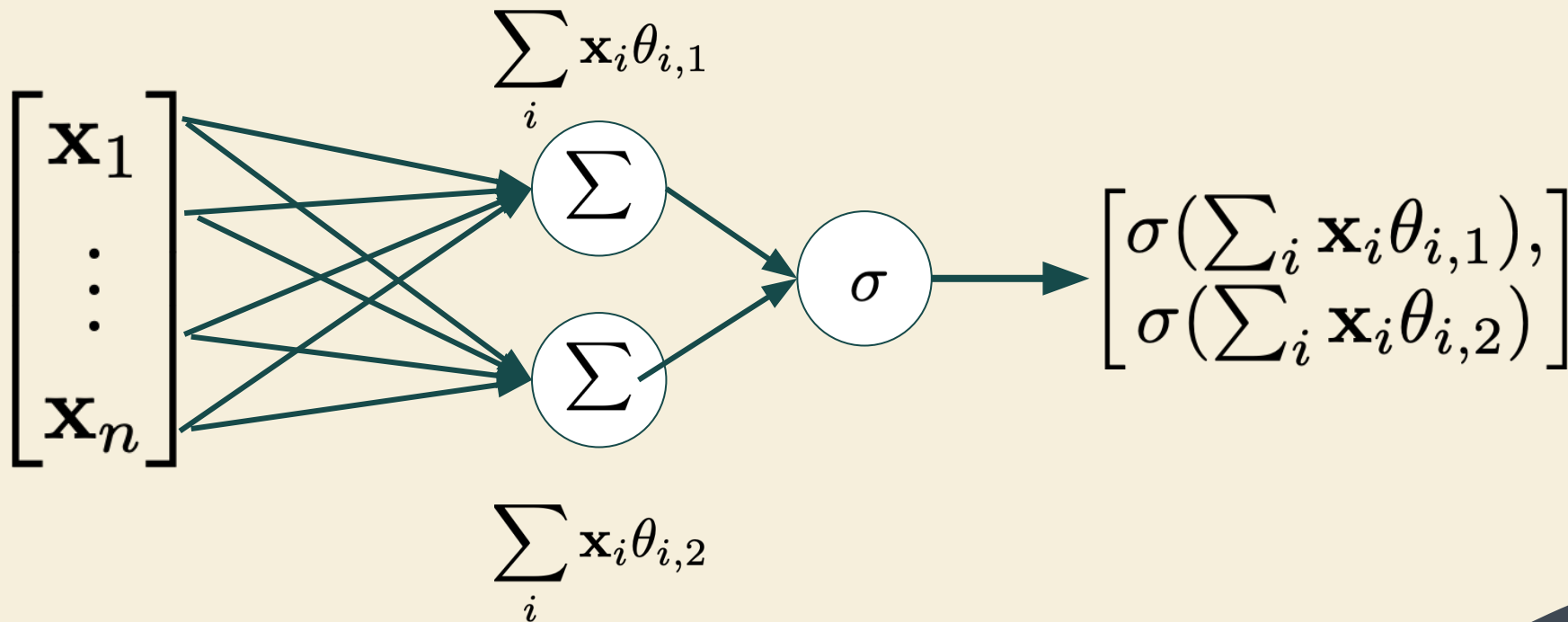


# Formulating the MLP

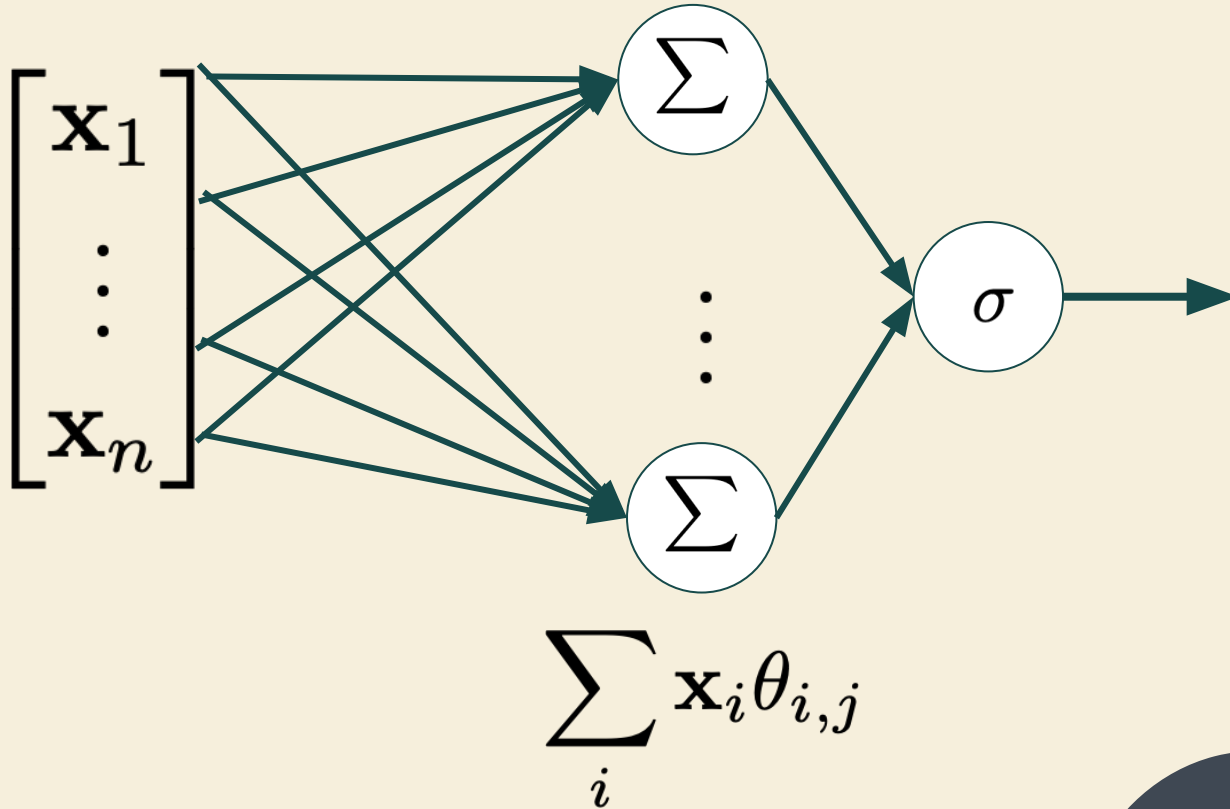




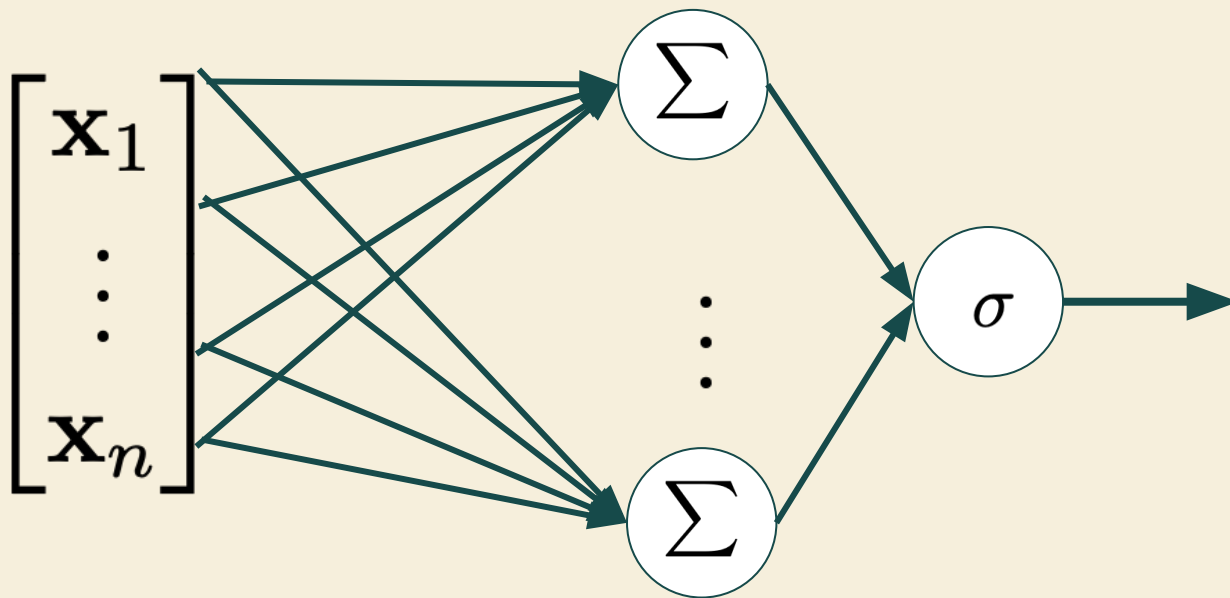
# Formulating the MLP



# Formulating the MLP



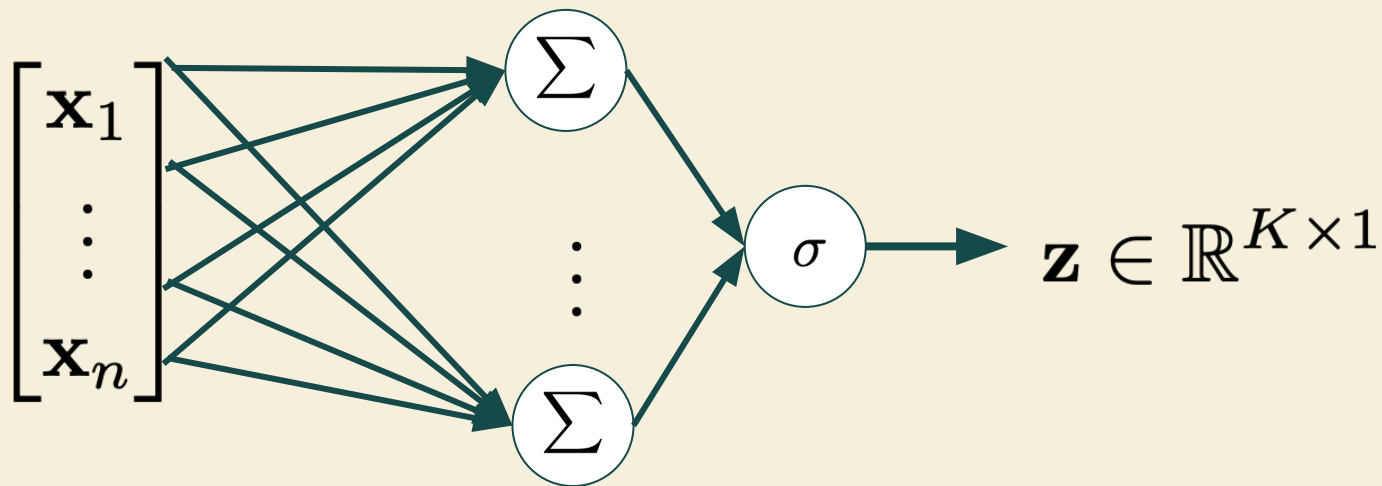
# Formulating the MLP



Matrix-vector  
product (just  
need to get  
shapes right)

$$\sum_i \mathbf{x}_i \theta_{i,j}$$

# Formulating the MLP



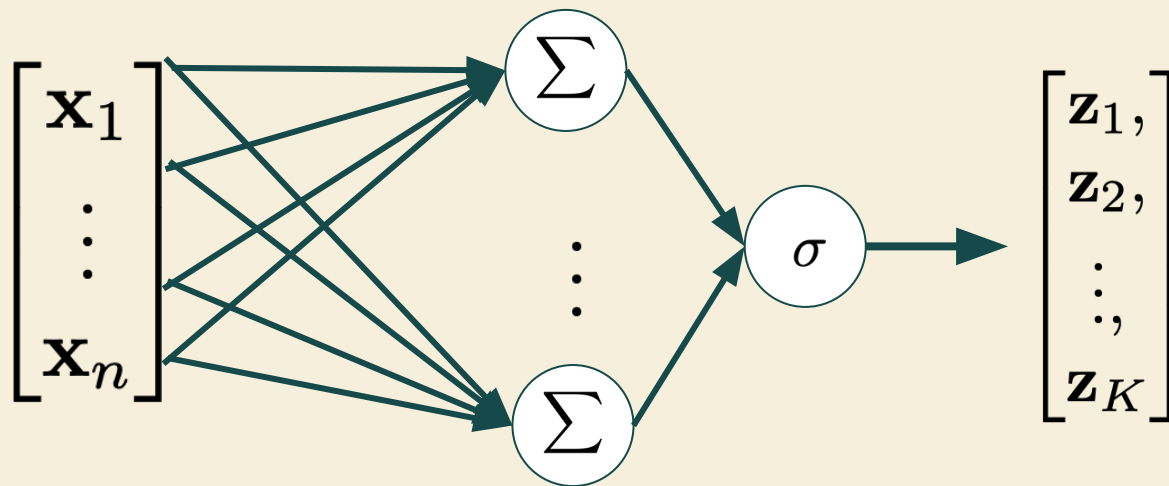
Matrix-vector  
product

$$\sum_i \mathbf{W}_{i,j} \mathbf{x}_i$$

$K \times n$        $n \times 1$

$\rightarrow K \times 1$

# Formulating the MLP

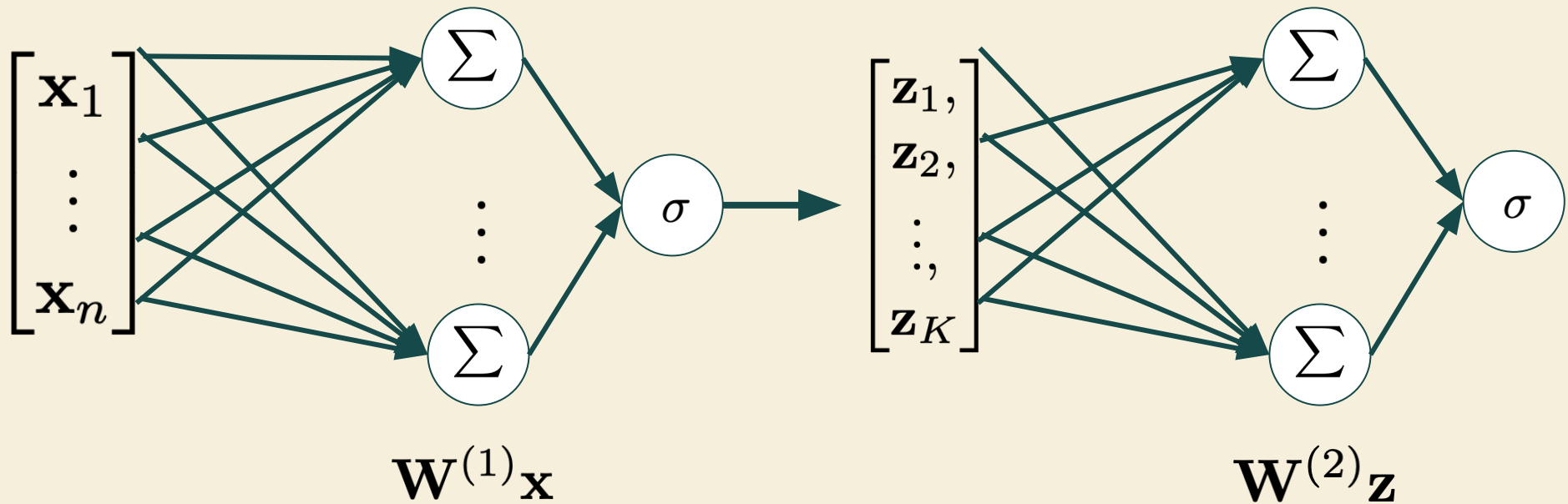


Matrix-vector  
product

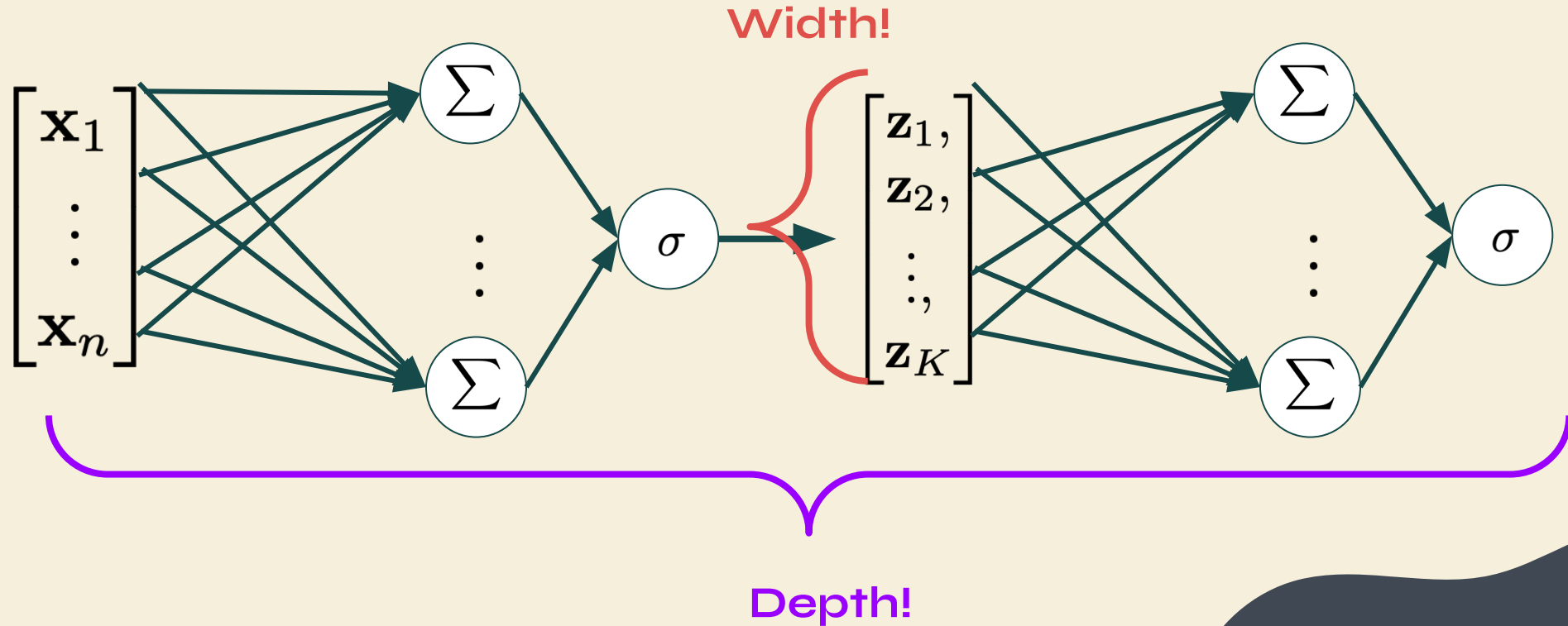
$$\sum_i \mathbf{W}_{i,j} \mathbf{x}_i$$

$K \times n$        $n \times 1$

# Formulating the MLP



# Formulating the MLP



# Mathematical abstraction

$$\mathbf{z}^{(1)} = \mathbf{x}$$

$$\mathbf{z}^{(i+1)} = \sigma(\mathbf{W}^{(i)} \mathbf{z}^{(i)} + \mathbf{b}^{(i)})$$



# Learning in a 1-layer NN

$$\mathbf{z}^{(1)} = \mathbf{x}$$

$$\mathbf{z}^{(i+1)} = \sigma(\mathbf{W}^{(i)} \mathbf{z}^{(i)} + \mathbf{b}^{(i)})$$

# Learning in fully connected networks

$$\hat{y} = \sigma(\mathbf{W}\mathbf{x})$$

$$\mathcal{L} = (\sigma(\mathbf{W}\mathbf{x}) - \mathbf{y})^2$$



Activation function is very generally defined and so we don't have anything analytically tractable in general

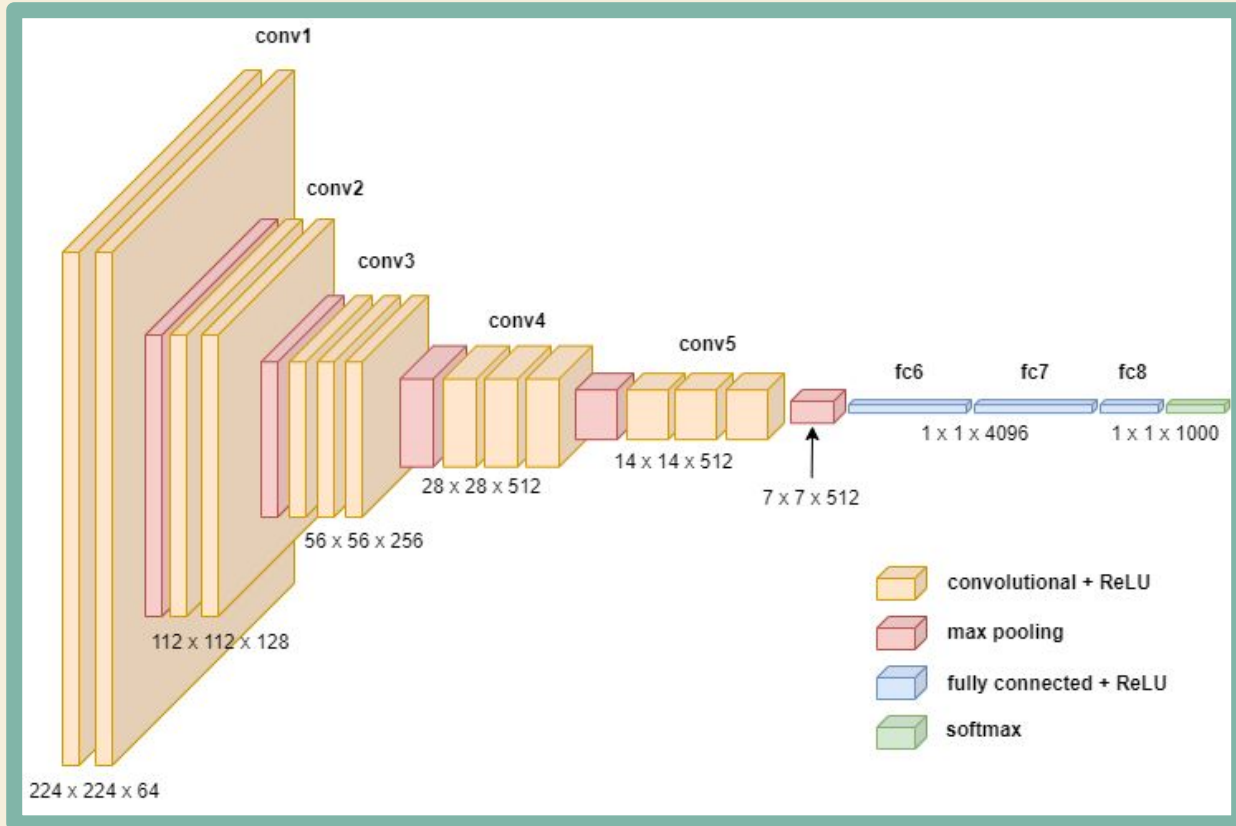
# Learning in fully connected networks

$$\log(1 + e^{wx+b})$$

$$\frac{e^{b_1+b_2+w_1x+w_2 \log[1 + e^{b_1+w_1x}]} w_2 x}{(1 + e^{b_1+w_1x}) \left(1 + e^{b_2+w_2 \log[1 + e^{b_1+w_1x}]}\right)}$$

Activation function is very generally defined and so we don't have anything analytically tractable in general

# Learning in more complex models



The background of the slide features abstract, wavy shapes in teal and orange. A large teal shape occupies the upper right and bottom right areas, while an orange shape fills the upper left and bottom left areas. The two colors meet in a diagonal line across the center. There are also thin, wavy lines in the opposite colors: a thin orange line in the top right and a thin teal line in the bottom left.

# Break time!

(save questions for the end of lecture)



# Automatic differentiation

We want to be able to differentiate extremely complex functions without needing to symbolically differentiate.

We call this numerical differentiation

# Automatic differentiation

$$f : \mathbb{R}^n \rightarrow \mathbb{R} \quad \frac{\partial f}{\partial x_1} = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(x)}{h},$$

$$\frac{\partial f}{\partial x_2} = \lim_{h \rightarrow 0} \frac{f(x_1, x_2 + h, \dots, x_n) - f(x)}{h},$$

$\vdots$

$$\frac{\partial f}{\partial x_n} = \lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, x_n + h) - f(x)}{h}$$

What about the  
finite difference  
method?

# Automatic differentiation

$$f : \mathbb{R}^n \rightarrow \mathbb{R} \quad \frac{\partial f}{\partial x_1} = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(x)}{h},$$

$$\frac{\partial f}{\partial x_2} = \lim_{h \rightarrow 0} \frac{f(x_1, x_2 + h, \dots, x_n) - f(x)}{h},$$

$\vdots$

$$\frac{\partial f}{\partial x_n} = \lim_{h \rightarrow 0} \frac{f(x_1, x_2, \dots, x_n + h) - f(x)}{h}$$

What about the  
finite difference  
method?

Too slow, too  
approximate :(






# Automatic differentiation

We want to be able to differentiate extremely complex functions without needing to symbolically differentiate.

Let's come up with a systematic way of computing numerical derivatives!



# Forward Mode Auto. Diff.

Key idea: along with computing the forward pass values, also keep track of the derivative. We will see this from a programming perspective and a *mathematical one*.

# Forward Mode Auto. Diff.

Consider the function:

$$f(\mathbf{x}) = h(g(\mathbf{x}))$$

# Forward Mode Auto. Diff.

Consider the function:

$$f(\mathbf{x}) = h(g(\mathbf{x}))$$

$$f(\mathbf{x}) = (h \circ g)(\mathbf{x}) = h(g(\mathbf{x}))$$

Aside: you will also see function composition this way

# Forward Mode Auto. Diff.

Consider the function:

$$f(\mathbf{x}) = h(g(x))$$

Chain rule tells us that:

$$\mathbf{J}_f = \mathbf{J}_{h \circ g} = \mathbf{J}_h(g(\mathbf{x})) \mathbf{J}_g(x)$$

*Jacobian*

*diff outside, diff inside*

$f = h(g(a(x)))$

$J_f = J_h(g(a(x))) J_g(a(x)) J_a(x)$

compute  $a(x) = y_1$ , do  $J_a(x)$   
compute  $g(y_1) = y_2$ , do  $J_g(y_1)$   
compute  $h(y_2) = y_3$ , do  $J_h(y_2)$

# Forward Mode Auto. Diff.

Generalizing the above function we have:

$$f = f^{(L)} \circ f^{(L-1)} \circ \dots \circ f^1$$

*$f^L(f^{L-1}(\dots))$*

Which in turn gives us the Jacobians:

$$\mathbf{J} = \mathbf{J}^{(L)} \cdot \mathbf{J}^{(L-1)} \cdot \dots \cdot \mathbf{J}^{(1)}$$

*diff of outside  $\times$  diff of inside*

# Forward Mode Auto. Diff.

To compute the gradient-input product we can use the products of Jacobians formula to see:

*must both sides by  $x$*

$$\mathbf{J}^{\downarrow} \mathbf{x} = \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(2)} \underbrace{(\mathbf{J}^{(1)} \mathbf{x}^{\downarrow})}$$

Compute the derivative of  
the first operation first!

# Forward Mode Auto. Diff.

To compute the gradient-input product we can use the products of Jacobians formula to see:

$$\begin{aligned}\mathbf{J}\mathbf{x} &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(2)} (\mathbf{J}^{(1)} \mathbf{x}) \\ &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(3)} (\underbrace{\mathbf{J}^{(2)} \mathbf{x}^{(1)}}_{\text{Jacobian wrt } \mathbf{x}^{(1)} ?})\end{aligned}$$

*Handwritten notes:*  
-  $\mathbf{x}^{(1)}$  (above  $\mathbf{J}^{(1)} \mathbf{x}$ )  
-  $\mathbf{J}^{(1)}$  = Jacobian wrt  $\mathbf{x}$  ?  
- Jacobian wrt  $\mathbf{x}^{(1)}$  ? (pointing to the underbrace)

$$\begin{aligned}\bar{\mathbf{J}}_x &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(3)} \mathbf{J}^{(2)} (\mathbf{J}^{(1)} x) \\ &= \dots \bar{\mathbf{J}}^{(3)} \mathbf{J}^{(2)}\end{aligned}$$

Use that result to compute the next Jacobian product



*Keep track of how  $J^n$  w.r.t input for which you calc. the den*

## Forward Mode Auto. Diff.

$$\begin{aligned}\mathbf{J}\mathbf{x} &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(2)} (\mathbf{J}^{(1)} \mathbf{x}) \\ &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(3)} (\mathbf{J}^{(2)} \mathbf{x}^{(1)}) \\ &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(4)} (\mathbf{J}^{(3)} \mathbf{x}^{(2)}) \\ &\dots \\ &= \mathbf{J}^L \mathbf{x}^{(L-1)}\end{aligned}$$

Continue until we have just  
the final Jacobian

# Programming Perspective: Operator Overloading

```
class AdVar:
    ...

    We want to keep track of both the value and
    the derivative of our variable, so it will
    hold two values one representing each
    ...

    def __init__(self, **kwargs):
        self.val = 0
        self.der = 1

        # re-assign these default values
        if 'val' in kwargs:
            self.val = kwargs['val']
        if 'der' in kwargs:
            self.der = kwargs['der']
```

↑  
allows operators like  $+$ ,  $-$ ,  $*$  etc. to  
have diff. meanings based on context

# Programming Perspective: Operator Overloading

```
class AdVar:
```

```
'''
```

We want to keep track of both the value and the derivative of our variable, so it will hold two values one representing each

```
'''
```

```
def __init__(self,**kwargs):
```

```
    self.val = 0
```

```
    self.der = 1
```

```
    # re-assign these default values
```

```
    if 'val' in kwargs:
```

```
        self.val = kwargs['val']
```

```
    if 'der' in kwargs:
```

```
        self.der = kwargs['der']
```

```
def add(a,b):
```

```
    # Create output evaluation and derivative object
```

```
    c = AdVar()
```

```
    # switch to determine if a or b is a constant
```

```
    if type(a) != AdVar:
```

```
        c.val = a + b.val
```

```
        c.der = b.der
```

```
    elif type(b) != AdVar:
```

```
        c.val = a.val + b
```

```
        c.der = a.der
```

```
    else:
```

```
        c.val = a.val + b.val
```

```
        c.der = a.der + b.der
```

```
    return c
```

*keep track of derivative*

When we update our value, we must also update our derivative!

## Consider the simple neural network:

$$\mathbf{z}^{(1)} = \mathbf{x} \quad \rightarrow \text{do } J_{z^{(1)}}(x)$$

$$\mathbf{z}^{(2)} = \tanh(\mathbf{W}^{(1)} \mathbf{z}^{(1)}) \rightarrow \text{do } J_{z^{(1)}}(z^{(1)})$$

$$\mathbf{z}^{(3)} = \mathbf{W}^{(2)} \mathbf{z}^{(2)} \quad \rightarrow \text{do } J_{z^{(1)}}(z^{(2)})$$

## Consider the simple neural network:

$$\mathbf{z}^{(1)} = \mathbf{x}$$

$$\mathbf{z}^{(2)} = \tanh(\mathbf{W}^{(1)} \mathbf{z}^{(1)})$$

$$\mathbf{z}^{(3)} = \mathbf{W}^{(2)} \mathbf{z}^{(2)}$$

We start by considering all of the operations  
that are needed to get us from the input to the  
loss function

## Consider the simple neural network:

$$\mathbf{z}^{(2)} = \text{Matmul}\left(\mathbf{W}^{(1)}, \tanh\left(\text{Matmul}(\mathbf{W}^{(0)} \mathbf{z}^{(0)})\right)\right)$$

$$\mathcal{L} = (\mathbf{y} - \mathbf{z}^{(2)})^2$$

We start by considering all of the operations that are needed to get us from the input to the loss function

(Sorry I switched indexing schemes)

*starts at  $z^{(0)}$  here*

## Consider the simple neural network:

$$\mathbf{z}^{(2)} = \text{Matmul}\left(\mathbf{W}^{(1)}, \tanh\left(\text{Matmul}(\mathbf{W}^{(0)} \mathbf{z}^{(0)})\right)\right)$$

$$\mathcal{L} = (\mathbf{y} - \mathbf{z}^{(2)})^2$$

Now, let's write out all of the derivative rules we need!

## Consider the simple neural network:

$$\mathbf{z}^{(2)} = \text{Matmul}(\mathbf{W}^{(1)}, \tanh(\text{Matmul}(\mathbf{W}^{(0)} \mathbf{z}^{(0)})))$$

*Handwritten red annotations:*  
- Under  $\mathbf{z}^{(2)}$ :  $g(w)$  and  $z^{(2)}$  with an arrow pointing to  $\mathbf{z}^{(2)}$ .  
- Under  $\tanh$ :  $g(w)$  and  $z^{(1)}$  with an arrow pointing to the  $\tanh$  function.  
- Under  $\mathbf{z}^{(0)}$ :  $x$  with an arrow pointing to  $\mathbf{z}^{(0)}$ .

$$\mathcal{L} = (\mathbf{y} - \mathbf{z}^{(2)})^2$$

Operation	Value Update	Derivative Update
Addition of a constant $c$	$g(w) + c$	$\frac{d}{dw}(g(w) + c) = \frac{d}{dw}g(w)$
Multiplication by a constant $c$	$cg(w)$	$\frac{d}{dw}(cg(w)) = c \frac{d}{dw}g(w)$
Raising to a power $n$	$g(w)^n$	$\frac{d}{dw}(g(w)^n) = n(g(w)^{n-1}) \frac{d}{dw}g(w)$
Applying Tanh	$\tanh(g(w))$	$\frac{d}{dw}(\tanh(g(w))) = 1 - \tanh^2(g(w)) \frac{d}{dw}g(w)$



## Consider the simple neural network:

Operation	Value Update	Derivative Update
Addition of a constant $c$	$g(w) + c$	$\frac{d}{dw}(g(w) + c) = \frac{d}{dw}g(w)$
Multiplication by a constant $c$	$cg(w)$	$\frac{d}{dw}(cg(w)) = c\frac{d}{dw}g(w)$
Raising to a power $n$	$g(w)^n$	$\frac{d}{dw}(g(w)^n) = n(g(w)^{n-1})\frac{d}{dw}g(w)$
Applying Tanh	$\tanh(g(w))$	$\frac{d}{dw}(\tanh(g(w))) = 1 - \tanh^2(g(w))\frac{d}{dw}$

Now lets go through each operation and write out how we update the value and how we update the derivative

## Consider the simple neural network:

$$\mathbf{z}^{(2)} = \text{Matmul}(\mathbf{W}^{(1)}, \tanh(\text{Matmul}(\mathbf{W}^{(0)} \mathbf{z}^{(0)})))$$

$$\mathcal{L} = (\mathbf{y} - \mathbf{z}^{(2)})^2$$

Value  
 $x_0 = x$   *$z^{(0)} = 1$*

Derivative  
 $d_0 = 1$  *deriv of input = 1*

*diff wrt  $x$  ( $x$  is input)*

## Consider the simple neural network:

$$\mathbf{z}^{(2)} = \text{Matmul}(\mathbf{W}^{(1)}, \tanh(\text{Matmul}(\mathbf{W}^{(0)} \mathbf{z}^{(0)})))$$

$$\mathcal{L} = (\mathbf{y} - \mathbf{z}^{(2)})^2$$

$$[ ] [ ] = [ ]$$

$$Ax \rightarrow A_i = \sum_j a_{ij} x_j \quad \frac{\partial}{\partial a} = \sum_j a_{ij} \delta_j$$

$$= \sum a_{iq}$$

Value	Derivative
$x_0 = x$	$d_0 = 1$
$x_1 = W^{(0)} x_0$	$d_1 = W^{(0)\top} d_0$

d. w. w.r.t  $x_0$   
 $\longrightarrow$

d. w.  $x_0 = d_0 (=1)$

## Consider the simple neural network:

$$\mathbf{z}^{(2)} = \text{Matmul}\left(\mathbf{W}^{(1)}, \tanh(\text{Matmul}(\mathbf{W}^{(0)} \mathbf{z}^{(0)}))\right)$$

$$\mathcal{L} = (\mathbf{y} - \mathbf{z}^{(2)})^2$$

Value	Derivative
$x_0 = x$	$d_0 = \mathbf{1}$
$x_1 = W^{(0)} x_0$	$d_1 = W^{(0)\top} d_0$
$x_2 = \tanh(x_1)$	$d_2 = 1 - \tanh^2(x_1) d_1$

*d. w. r. t  $x_1$*

# Consider the simple neural network:

$$\mathbf{z}^{(2)} = \text{Matmul}(\mathbf{W}^{(1)}, \tanh(\text{Matmul}(\mathbf{W}^{(0)} \mathbf{z}^{(0)})))$$

$$\mathcal{L} = (\mathbf{y} - \mathbf{z}^{(2)})^2$$

Operation	Value Update	Derivative Update
Addition of a constant $c$	$g(w) + c$	$\frac{d}{dw}(g(w) + c) = \frac{d}{dw}g(w)$
Multiplication by a constant $c$	$cg(w)$	$\frac{d}{dw}(cg(w)) = c \frac{d}{dw}g(w)$
Raising to a power $n$	$g(w)^n$	$\frac{d}{dw}(g(w)^n) = n(g(w)^{n-1}) \frac{d}{dw}g(w)$
Applying Tanh	$\tanh(g(w))$	$\frac{d}{dw}(\tanh(g(w))) = 1 - \tanh^2(g(w)) \frac{d}{dw}g(w)$

Value

$$x_0 = x$$

$$x_1 = W^{(0)} x_0$$

$$x_2 = \tanh(x_1)$$

$$\stackrel{=: z^{(2)}}{\rightarrow} x_3 = W^{(1)} x_2$$

$$x_4 = (x_3 - y) \quad y - x_3$$

$$\stackrel{=: \text{loss}}{\rightarrow} x_5 = (x_4)^2$$

Derivative

$$d_0 = 1 \quad \text{transposed as } d_0 \text{ is a row vector by convention}$$

$$\frac{\partial x_1}{\partial x_0} : d_1 = W^{(0)\top} d_0 \quad W^{(0)\top} \frac{\partial}{\partial x_1}(x_1) = W^{(0)\top} d_0$$

$$\frac{\partial x_2}{\partial x_1} : d_2 = 1 - \tanh^2(x_1) d_1 \quad \frac{\partial}{\partial x_1}(\tanh(x_1)) = 1 - \tanh^2(x_1)$$

$$d_3 = W^{(1)\top} d_2$$

$$\frac{\partial x_4}{\partial x_3} : d_4 = d_3$$

$$\frac{\partial x_5}{\partial x_4} : d_5 = -2(x_4) d_4$$

$$2x_4 \frac{\partial}{\partial x_4}(x_4)$$

intermediate variables  $\rightarrow$

# Consider the simple neural network:

$x_3 - y$  is -ve. b/c this -ve would come out in  $d_4$ .  $x_4$  should be 1.  $y - x_3 \rightarrow d_4 = -d_3$  so  $d_4 = 2d_3$ , where  $d_4 = -d_3$

$$\mathbf{z}^{(2)} = \text{Matmul}(\mathbf{W}^{(1)}, \tanh(\text{Matmul}(\mathbf{W}^{(0)} \mathbf{z}^{(0)})))$$

$$\mathcal{L} = (\mathbf{y} - \mathbf{z}^{(2)})^2$$

$$x = 1, W^{(0)} = 2, W^{(1)} = 1.2$$

diff w.r.t  $x$

Value	Derivative
$x_0 = x$	$d_0 = 1$
$x_1 = W^{(0)} x_0$	$d_1 = W^{(0)\top} d_0$
$x_2 = \tanh(x_1)$	$d_2 = 1 - \tanh^2(x_1) d_1$
$x_3 = W^{(1)} x_2$	$d_3 = W^{(1)\top} d_2$
$x_4 = (x_3 - y)$	$d_4 = d_3$
$x_5 = (x_4)^2$	$d_4 = -2(x_4) d_4$

$$x_0 = 1 \quad x_2^{(0)}$$

$$d_0 = 1$$

# Consider the simple neural network:

$$\mathbf{z}^{(2)} = \text{Matmul}\left(\mathbf{W}^{(1)}, \tanh\left(\text{Matmul}(\mathbf{W}^{(0)} \mathbf{z}^{(0)})\right)\right)$$

$$\mathcal{L} = (\mathbf{y} - \mathbf{z}^{(2)})^2$$

$$x = 1, W^{(0)} = 2, W^{(1)} = 1.2$$

deriv. w.r.t.  $x$  at  $x=1$

$$x_0 = 1$$

$$x_1 = 1 * 2 = 2$$

$$d_0 = 1$$

$$d_1 = 2 * d_0 = 2$$

Value	Derivative
$x_0 = x$	$d_0 = 1$
$x_1 = W^{(0)} x_0$	$d_1 = W^{(0)\top} d_0$
$x_2 = \tanh(x_1)$	$d_2 = 1 - \tanh^2(x_1) d_1$
$x_3 = W^{(1)} x_2$	$d_3 = W^{(1)\top} d_2$
$x_4 = (x_3 - y)$	$d_4 = d_3$
$x_5 = (x_4)^2$	$d_4 = -2(x_4) d_4$

## Consider the simple neural network:

$$\mathbf{z}^{(2)} = \text{Matmul}(\mathbf{W}^{(1)}, \tanh(\text{Matmul}(\mathbf{W}^{(0)} \mathbf{z}^{(0)})))$$

$$\mathcal{L} = (\mathbf{y} - \mathbf{z}^{(2)})^2$$

$$x = 1, W^{(0)} = 2, W^{(1)} = 1.2$$

$$x_0 = 1$$

$$x_1 = 1 * 2 = 2$$

$$x_2 = \tanh(2) \approx 0.964$$

$$d_0 = 1$$

$$d_1 = 2 * d_0 = 2$$

$$d_2 = (1 - \tanh^2(2)) d_1 \approx 0.14$$

Value

$$x_0 = x$$

$$x_1 = W^{(0)} x_0$$

$$x_2 = \tanh(x_1)$$

$$x_3 = W^{(1)} x_2$$

$$x_4 = (x_3 - y)$$

$$x_5 = (x_4)^2$$

Derivative

$$d_0 = 1$$

$$d_1 = W^{(0)\top} d_0$$

$$d_2 = 1 - \tanh^2(x_1) d_1$$

$$d_3 = W^{(1)\top} d_2$$

$$d_4 = d_3$$

$$d_5 = -2(x_4) d_4$$



# Consider the simple neural network:

$$\mathbf{z}^{(2)} = \text{Matmul}(\mathbf{W}^{(1)}, \tanh(\text{Matmul}(\mathbf{W}^{(0)} \mathbf{z}^{(0)})))$$

$$\mathcal{L} = (\mathbf{y} - \mathbf{z}^{(2)})^2$$

$$x = 1, W^{(0)} = 2, W^{(1)} = 1.2$$

For  $x_4, x_5, d_4, d_5$ , alternate version which notes we solve in 'or's' - i.e.  $y - z^{(2)}$  rather than  $z^{(2)} - y$

$$x_0 = 1$$

$$x_1 = 1 * 2 = 2$$

$$x_2 = \tanh(2) \approx 0.964$$

$$x_3 = 0.964 * 1.2 \approx 1.15$$

$$(x_3 - y) \rightarrow x_4 = (1.15 - 2) \approx -0.85 \text{ or } y - x_3 = 0.85$$

$$x_5 = (-0.85)^2 = 0.722 \text{ or } x_5 = |x_4|^2$$

$$d_0 = 1$$

$$d_1 = 2 * d_0 = 2$$

$$d_2 = (1 - \tanh^2(2))d_1 \approx 0.14$$

$$d_3 = 1.2 * d_2 \approx 0.169$$

$$d_4 = d_3 = 0.169 \text{ or } d_4 = -d_3 = -0.169$$

$$d_5 = 2(-0.85)d_4 = -0.283 \text{ or } d_5 = 2(x_4)d_4 = 2(-0.85)(0.169) = -0.283$$

Value

Derivative

$$x_0 = x$$

$$d_0 = 1$$

$$x_1 = W^{(0)} x_0$$

$$d_1 = W^{(0)\top} d_0$$

$$x_2 = \tanh(x_1)$$

$$d_2 = 1 - \tanh^2(x_1) d_1$$

$$x_3 = W^{(1)} x_2$$

$$d_3 = W^{(1)\top} d_2$$

$$x_4 = (x_3 - y)$$

$$d_4 = d_3$$

$$x_5 = (x_4)^2$$

$$d_4 = -2(x_4) d_4$$

# Taking a step back

A single pass of forward mode automatic differentiation computes the "directional derivative" with respect to  $\mathbf{x}$ . By computing it with respect to  $\mathbf{e}_i$ , we can get one column of the Jacobian

$$\begin{aligned}\mathbf{J}\mathbf{x} &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(2)} (\mathbf{J}^{(1)} \mathbf{x}) \\ &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(3)} (\mathbf{J}^{(2)} \mathbf{x}^{(1)}) \\ &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(4)} (\mathbf{J}^{(3)} \mathbf{x}^{(2)}) \\ &\dots \\ &= \mathbf{J}^L \mathbf{x}^{(L-1)}\end{aligned}$$

Jacobian:

$$\begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}$$

# Taking a step back

To get the full Jacobian we would need to compute as many forward mode sweeps as we have columns

But! In ML problems we have MANY columns and few rows, so this seems inefficient.

$$\begin{aligned}\mathbf{J}\mathbf{x} &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(2)} (\mathbf{J}^{(1)} \mathbf{x}) \\ &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(3)} (\mathbf{J}^{(2)} \mathbf{x}^{(1)}) \\ &= \mathbf{J}^{(L)} \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(4)} (\mathbf{J}^{(3)} \mathbf{x}^{(2)}) \\ &\dots \\ &= \mathbf{J}^L \mathbf{x}^{(L-1)}\end{aligned}$$

*ML usually  
many inputs,  
few outputs so  
inefficient*

# Reverse Mode Auto. Diff.

Reverse mode is also known as adjoint mode because we are interested in computing the adjoint derivatives:

$$\bar{x} = \frac{\partial z}{\partial x}$$

*grad of final output  
w.r.t that variable*

# Reverse Mode Auto. Diff.

Reverse mode is also known as adjoint mode because we are interested in computing the adjoint derivatives:

$$\mathbf{J}_x = (\mathbf{y} \mathbf{J}^{(L)}) \mathbf{J}^{(L-1)} \dots \mathbf{J}^{(1)}$$

# Reverse Mode Auto. Diff.

Reverse mode is also known as adjoint mode because we are interested in computing the adjoint derivatives:

$$\begin{aligned}\mathbf{J}\mathbf{x} &= (\mathbf{y}\mathbf{J}^{(L)})\mathbf{J}^{(L-1)} \dots \mathbf{J}^{(1)} \\ &= (\bar{\mathbf{y}}^{(L)}\mathbf{J}^{(L-1)})\mathbf{J}^{(L-2)} \dots \mathbf{J}^{(1)}\end{aligned}$$

# Reverse Mode Auto. Diff.

Reverse mode is also known as adjoint mode because we are interested in computing the adjoint derivatives:

$$\begin{aligned}\mathbf{J}\mathbf{x} &= (\mathbf{y}\mathbf{J}^{(L)})\mathbf{J}^{(L-1)} \dots \mathbf{J}^{(1)} \\ &= (\bar{\mathbf{y}}^{(L)}\mathbf{J}^{(L-1)})\mathbf{J}^{(L-2)} \dots \mathbf{J}^{(1)} \\ &\dots \\ &= \bar{\mathbf{y}}^{(2)}\mathbf{J}^{(1)}\end{aligned}$$

# Key Observations

$$\mathbf{J}\mathbf{x} = (\mathbf{y}\mathbf{J}^{(L)})\mathbf{J}^{(L-1)} \dots \mathbf{J}^{(1)}$$



This is the initial co-tangent (output) that we will compute all our adjoints with respect to! But how did we get it?  
-> we have to compute  $f(\mathbf{x})$  first



# Key Observations

$$\mathbf{J}\mathbf{x} = (\mathbf{y}\mathbf{J}^{(L)})\mathbf{J}^{(L-1)} \dots \mathbf{J}^{(1)}$$



Similarly each adjoint requires us to use the intermediate computations that got us to  $\mathbf{y}$

This is the initial co-tangent (output) that we will compute all our adjoints with respect to! But how did we get it?  
-> we have to compute  $f(\mathbf{x})$  first

# Key Observations

- We need to complete a complete forward pass to get our initial cotangent
- We need to store all of our intermediate computations to compute the adjoints we need

*Note: we have a discussion of reverse mode for NNs in the notes, but won't have time in lecture to go through that.*

# Computational Graphs

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2))$$

Neural networks are actually a particularly easy/nice case for adjoint/reverse mode auto. diff. so let's look at something a little harder!

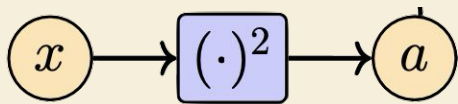
# Computational Graphs

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2))$$

1. Identify the elementary operations and look for the ones that occur most and for shared structure
2. We will have variable nodes leading to operations.
3. Start with the most common operations and then build the graph outwards

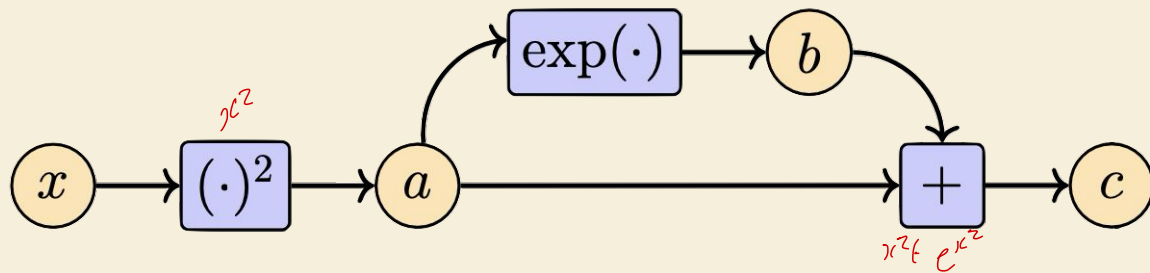
# Computational Graphs

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2))$$



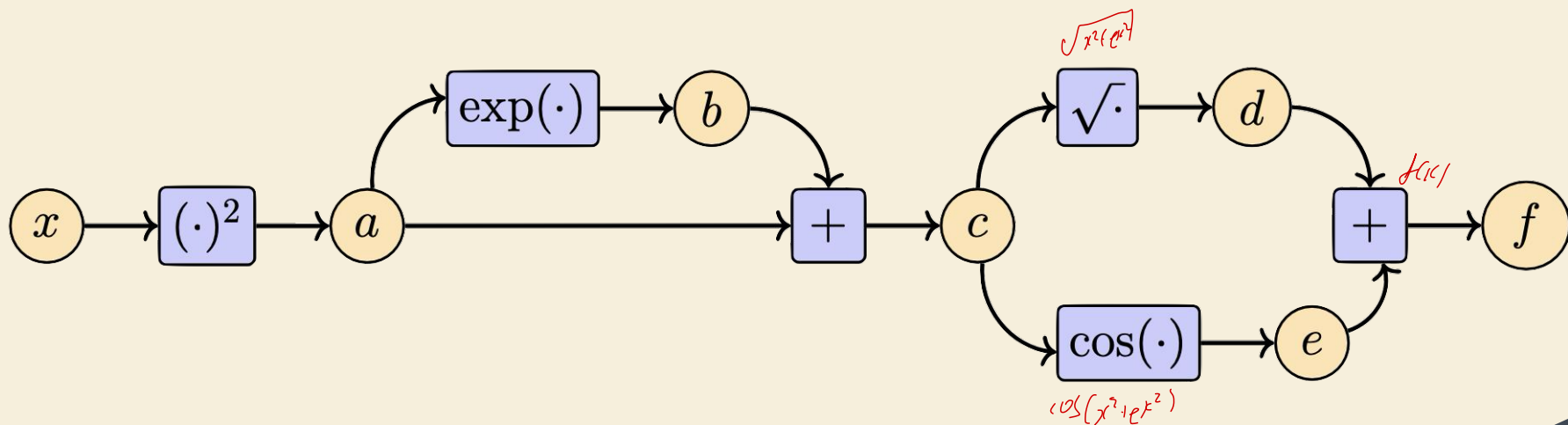
# Computational Graphs

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2))$$



# Computational Graphs

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2))$$



# Computational Graphs

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2))$$

$$b = e^a$$

$$c = a + b$$

$$\frac{\partial a}{\partial x} = 2x$$

$$\frac{\partial b}{\partial a} = \exp(a)$$

$$\frac{\partial c}{\partial a} = 1 = \frac{\partial c}{\partial b}$$

$$d = \sqrt{c} \quad \frac{\partial d}{\partial c} = \frac{1}{2\sqrt{c}}$$

$$e = \cos(c) \quad \frac{\partial e}{\partial c} = -\sin(c)$$

$$f = d + e \quad \frac{\partial f}{\partial d} = 1 = \frac{\partial f}{\partial e}$$



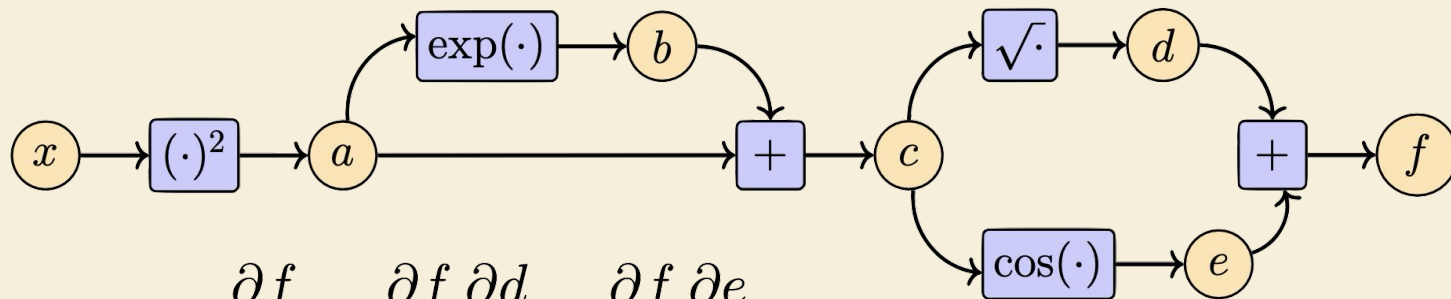
# Computational Graphs

$$f(x) = \sqrt{x^2 + \exp(x^2)} + \cos(x^2 + \exp(x^2))$$

To compute the derivative we want, we have to back propagate through all of the nodes in the graph using the chain rule for each incoming node and summing them together

$$\frac{\partial f}{\partial x_i} = \sum_{x_j: x_i \in \text{Pa}(x_j)} \frac{\partial f}{\partial x_j} \frac{\partial x_j}{\partial x_i} = \sum_{x_j: x_i \in \text{Pa}(x_j)} \frac{\partial f}{\partial x_j} \frac{\partial g_j}{\partial x_i}$$

# Computational Graphs



$$\frac{\partial f}{\partial c} = \frac{\partial f}{\partial d} \frac{\partial d}{\partial c} + \frac{\partial f}{\partial e} \frac{\partial e}{\partial c}$$

$$\frac{\partial f}{\partial b} = \frac{\partial f}{\partial c} \frac{\partial c}{\partial b}$$

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial b} \frac{\partial b}{\partial a} + \frac{\partial f}{\partial c} \frac{\partial c}{\partial a}$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial a} \frac{\partial a}{\partial x}$$



## Next Lecture:

*When and why does gradient descent  
give us a good parameter value?*

