



**Matrix differentiation. Automatic  
differentiation. Checkpointintg.**

**Daniil Merkulov**

Applied Math for Data Science. Sberuniversity.

# Matrix calculus

## Gradient

Let  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ , then vector, which contains all first-order partial derivatives:

$$\nabla f(x) = \frac{df}{dx} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$

# Gradient

Let  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ , then vector, which contains all first-order partial derivatives:

$$\nabla f(x) = \frac{df}{dx} = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$

named gradient of  $f(x)$ . This vector indicates the direction of the steepest ascent. Thus, vector  $-\nabla f(x)$  means the direction of the steepest descent of the function in the point. Moreover, the gradient vector is always orthogonal to the contour line in the point.

## Example

For the function  $f(x, y) = x^2 + y^2$ , the gradient is:

$$\nabla f(x, y) = \begin{bmatrix} 2x \\ 2y \end{bmatrix}$$

This gradient points in the direction of the steepest ascent of the function.

## Question

How does the magnitude of the gradient relate to the steepness of the function?

## Hessian

Let  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ , then matrix, containing all the second order partial derivatives:

$$f''(x) = \nabla^2 f(x) = \frac{\partial^2 f}{\partial x_i \partial x_j} = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix}$$

## Hessian

Let  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ , then matrix, containing all the second order partial derivatives:

$$f''(x) = \nabla^2 f(x) = \frac{\partial^2 f}{\partial x_i \partial x_j} = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix}$$

In fact, Hessian could be a tensor in such a way:  $(f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m)$  is just 3d tensor, every slice is just hessian of corresponding scalar function  $(\nabla^2 f_1(x), \dots, \nabla^2 f_m(x))$ .

### Example

For the function  $f(x, y) = x^2 + y^2$ , the Hessian is:

$$H_f(x, y) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

This matrix provides information about the curvature of the function in different directions.

### Question

How can the Hessian matrix be used to determine the concavity or convexity of a function?

## Schwartz theorem

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a function. If the mixed partial derivatives  $\frac{\partial^2 f}{\partial x_i \partial x_j}$  and  $\frac{\partial^2 f}{\partial x_j \partial x_i}$  are both continuous on an open set containing a point  $a$ , then they are equal at the point  $a$ . That is,

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(a) = \frac{\partial^2 f}{\partial x_j \partial x_i}(a)$$

## Schwartz theorem

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a function. If the mixed partial derivatives  $\frac{\partial^2 f}{\partial x_i \partial x_j}$  and  $\frac{\partial^2 f}{\partial x_j \partial x_i}$  are both continuous on an open set containing a point  $a$ , then they are equal at the point  $a$ . That is,

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(a) = \frac{\partial^2 f}{\partial x_j \partial x_i}(a)$$

Given the Schwartz theorem, if the mixed partials are continuous on an open set, the Hessian matrix is symmetric. That means the entries above the main diagonal mirror those below the main diagonal:

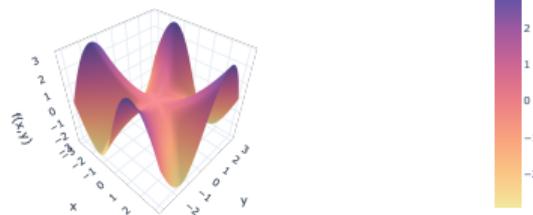
$$\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{\partial^2 f}{\partial x_j \partial x_i} \quad \nabla^2 f(x) = (\nabla^2 f(x))^T$$

This symmetry simplifies computations and analysis involving the Hessian matrix in various applications, particularly in optimization.

## Schwartz counterexample

$$f(x, y) = \begin{cases} \frac{xy(x^2 - y^2)}{x^2 + y^2} & \text{for } (x, y) \neq (0, 0), \\ 0 & \text{for } (x, y) = (0, 0). \end{cases}$$

### Counterexample ♣



One can verify, that  $\frac{\partial^2 f}{\partial x \partial y}(0, 0) \neq \frac{\partial^2 f}{\partial y \partial x}(0, 0)$ , although the mixed partial derivatives do exist, and at every other point the symmetry does hold.

# Jacobian

The extension of the gradient of multidimensional  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is the following matrix:

$$J_f = f'(x) = \frac{df}{dx^T} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix}$$

This matrix provides information about the rate of change of the function with respect to its inputs.

## Question

Can we somehow connect those three definitions above (gradient, jacobian, and hessian) using a single correct statement?

## Example

For the function

$$f(x, y) = \begin{bmatrix} x + y \\ x - y \end{bmatrix},$$

the Jacobian is:

$$J_f(x, y) = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

## Question

How does the Jacobian matrix relate to the gradient for scalar-valued functions?

## Summary

$$f(x) : X \rightarrow Y; \quad \frac{\partial f(x)}{\partial x} \in G$$

$X$	$Y$	$G$	Name
$\mathbb{R}$	$\mathbb{R}$	$\mathbb{R}$	$f'(x)$ (derivative)
$\mathbb{R}^n$	$\mathbb{R}$	$\mathbb{R}^n$	$\frac{\partial f}{\partial x_i}$ (gradient)
$\mathbb{R}^n$	$\mathbb{R}^m$	$\mathbb{R}^{n \times m}$	$\frac{\partial f_i}{\partial x_j}$ (jacobian)
$\mathbb{R}^{m \times n}$	$\mathbb{R}$	$\mathbb{R}^{m \times n}$	$\frac{\partial f}{\partial x_{ij}}$

# Differentials

## Theorem

Let  $x \in S$  be an interior point of the set  $S$ , and let  $D : U \rightarrow V$  be a linear operator. We say that the function  $f$  is differentiable at the point  $x$  with derivative  $D$  if for all sufficiently small  $h \in U$  the following decomposition holds:

$$f(x + h) = f(x) + D[h] + o(\|h\|)$$

If for any linear operator  $D : U \rightarrow V$  the function  $f$  is not differentiable at the point  $x$  with derivative  $D$ , then we say that  $f$  is not differentiable at the point  $x$ .

## Differentials

After obtaining the differential notation of  $df$  we can retrieve the gradient using the following formula:

$$df(x) = \langle \nabla f(x), dx \rangle$$

## Differentials

After obtaining the differential notation of  $df$  we can retrieve the gradient using the following formula:

$$df(x) = \langle \nabla f(x), dx \rangle$$

Then, if we have a differential of the above form and we need to calculate the second derivative of the matrix/vector function, we treat “old”  $dx$  as the constant  $dx_1$ , then calculate  $d(df) = d^2 f(x)$

$$d^2 f(x) = \langle \nabla^2 f(x)dx_1, dx \rangle = \langle H_f(x)dx_1, dx \rangle$$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X + Y) = dX + dY$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X + Y) = dX + dY$
- $d(X^T) = (dX)^T$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X + Y) = dX + dY$
- $d(X^T) = (dX)^T$
- $d(XY) = (dX)Y + X(dY)$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X + Y) = dX + dY$
- $d(X^T) = (dX)^T$
- $d(XY) = (dX)Y + X(dY)$
- $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X + Y) = dX + dY$
- $d(X^T) = (dX)^T$
- $d(XY) = (dX)Y + X(dY)$
- $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$
- $d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X + Y) = dX + dY$
- $d(X^T) = (dX)^T$
- $d(XY) = (dX)Y + X(dY)$
- $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$
- $d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$
- $d(\det X) = \det X \langle X^{-T}, dX \rangle$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X + Y) = dX + dY$
- $d(X^T) = (dX)^T$
- $d(XY) = (dX)Y + X(dY)$
- $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$
- $d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$
- $d(\det X) = \det X \langle X^{-T}, dX \rangle$
- $d(\text{tr } X) = \langle I, dX \rangle$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X + Y) = dX + dY$
- $d(X^T) = (dX)^T$
- $d(XY) = (dX)Y + X(dY)$
- $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$
- $d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$
- $d(\det X) = \det X \langle X^{-T}, dX \rangle$
- $d(\text{tr } X) = \langle I, dX \rangle$
- $df(g(x)) = \frac{df}{dg} \cdot dg(x)$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X + Y) = dX + dY$
- $d(X^T) = (dX)^T$
- $d(XY) = (dX)Y + X(dY)$
- $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$
- $d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$
- $d(\det X) = \det X \langle X^{-T}, dX \rangle$
- $d(\text{tr } X) = \langle I, dX \rangle$
- $df(g(x)) = \frac{df}{dg} \cdot dg(x)$
- $H = (J(\nabla f))^T$

## Differential properties

Let  $A$  and  $B$  be the constant matrices, while  $X$  and  $Y$  are the variables (or matrix functions).

- $dA = 0$
- $d(\alpha X) = \alpha(dX)$
- $d(AXB) = A(dX)B$
- $d(X + Y) = dX + dY$
- $d(X^T) = (dX)^T$
- $d(XY) = (dX)Y + X(dY)$
- $d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$
- $d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$
- $d(\det X) = \det X \langle X^{-T}, dX \rangle$
- $d(\text{tr } X) = \langle I, dX \rangle$
- $df(g(x)) = \frac{df}{dg} \cdot dg(x)$
- $H = (J(\nabla f))^T$
- $d(X^{-1}) = -X^{-1}(dX)X^{-1}$

## Matrix calculus. Example 1

### Example

Find  $df, \nabla f(x)$ , if  $f(x) = \langle x, Ax \rangle - b^T x + c$ .

## Matrix calculus. Example 2

### Example

Find  $df, \nabla f(x)$ , if  $f(x) = \ln\langle x, Ax \rangle$ .

1. It is essential for  $A$  to be positive definite, because it is a logarithm argument. So,  $A \in \mathbb{S}_{++}^n$ . Let's find the differential first:

$$\begin{aligned} df &= d(\ln\langle x, Ax \rangle) = \frac{d(\langle x, Ax \rangle)}{\langle x, Ax \rangle} = \frac{\langle dx, Ax \rangle + \langle x, d(Ax) \rangle}{\langle x, Ax \rangle} = \\ &= \frac{\langle Ax, dx \rangle + \langle x, Adx \rangle}{\langle x, Ax \rangle} = \frac{\langle Ax, dx \rangle + \langle A^T x, dx \rangle}{\langle x, Ax \rangle} = \frac{\langle (A + A^T)x, dx \rangle}{\langle x, Ax \rangle} \end{aligned}$$

## Matrix calculus. Example 2

### Example

Find  $df, \nabla f(x)$ , if  $f(x) = \ln\langle x, Ax \rangle$ .

1. It is essential for  $A$  to be positive definite, because it is a logarithm argument. So,  $A \in \mathbb{S}_{++}^n$ . Let's find the differential first:

$$\begin{aligned} df &= d(\ln\langle x, Ax \rangle) = \frac{d(\langle x, Ax \rangle)}{\langle x, Ax \rangle} = \frac{\langle dx, Ax \rangle + \langle x, d(Ax) \rangle}{\langle x, Ax \rangle} = \\ &= \frac{\langle Ax, dx \rangle + \langle x, Adx \rangle}{\langle x, Ax \rangle} = \frac{\langle Ax, dx \rangle + \langle A^T x, dx \rangle}{\langle x, Ax \rangle} = \frac{\langle (A + A^T)x, dx \rangle}{\langle x, Ax \rangle} \end{aligned}$$

2. Note, that our main goal is to derive the form  $df = \langle \cdot, dx \rangle$

$$df = \left\langle \frac{2Ax}{\langle x, Ax \rangle}, dx \right\rangle$$

Hence, the gradient is  $\nabla f(x) = \frac{2Ax}{\langle x, Ax \rangle}$

## Matrix calculus. Example 3

### Example

Find  $df, \nabla f(X)$ , if  $f(X) = \langle S, X \rangle - \log \det X$ .

## Automatic differentiation



@dpiponi@mathstodon.xyz  
@sigfpe

...

I think the first 40 years or so of automatic differentiation was largely people not using it because they didn't believe such an algorithm could possibly exist.

11:36 PM · Sep 17, 2019

---

9

26

159

13

↑

Figure 1: When you got the idea



Figure 2: This is not autograd

## Problem

Suppose we need to solve the following problem:

$$L(w) \rightarrow \min_{w \in \mathbb{R}^d}$$

## Problem

Suppose we need to solve the following problem:

$$L(w) \rightarrow \min_{w \in \mathbb{R}^d}$$

- Such problems typically arise in machine learning, when you need to find optimal hyperparameters  $w$  of an ML model (i.e. train a neural network).

## Problem

Suppose we need to solve the following problem:

$$L(w) \rightarrow \min_{w \in \mathbb{R}^d}$$

- Such problems typically arise in machine learning, when you need to find optimal hyperparameters  $w$  of an ML model (i.e. train a neural network).
- You may use a lot of algorithms to approach this problem, but given the modern size of the problem, where  $d$  could be dozens of billions it is very challenging to solve this problem without information about the gradients using zero-order optimization algorithms.

## Problem

Suppose we need to solve the following problem:

$$L(w) \rightarrow \min_{w \in \mathbb{R}^d}$$

- Such problems typically arise in machine learning, when you need to find optimal hyperparameters  $w$  of an ML model (i.e. train a neural network).
- You may use a lot of algorithms to approach this problem, but given the modern size of the problem, where  $d$  could be dozens of billions it is very challenging to solve this problem without information about the gradients using zero-order optimization algorithms.
- That is why it would be beneficial to be able to calculate the gradient vector  $\nabla_w L = \left( \frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_d} \right)^T$ .

## Problem

Suppose we need to solve the following problem:

$$L(w) \rightarrow \min_{w \in \mathbb{R}^d}$$

- Such problems typically arise in machine learning, when you need to find optimal hyperparameters  $w$  of an ML model (i.e. train a neural network).
- You may use a lot of algorithms to approach this problem, but given the modern size of the problem, where  $d$  could be dozens of billions it is very challenging to solve this problem without information about the gradients using zero-order optimization algorithms.
- That is why it would be beneficial to be able to calculate the gradient vector  $\nabla_w L = \left( \frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_d} \right)^T$ .
- Typically, first-order methods perform much better in huge-scale optimization, while second-order methods require too much memory.

## Example: multidimensional scaling

Suppose, we have a pairwise distance matrix for  $N$   $d$ -dimensional objects  $D \in \mathbb{R}^{N \times N}$ . Given this matrix, our goal is to recover the initial coordinates  $W_i \in \mathbb{R}^d$ ,  $i = 1, \dots, N$ .

## Example: multidimensional scaling

Suppose, we have a pairwise distance matrix for  $N$   $d$ -dimensional objects  $D \in \mathbb{R}^{N \times N}$ . Given this matrix, our goal is to recover the initial coordinates  $W_i \in \mathbb{R}^d$ ,  $i = 1, \dots, N$ .

$$L(W) = \sum_{i,j=1}^N (\|W_i - W_j\|_2^2 - D_{i,j})^2 \rightarrow \min_{W \in \mathbb{R}^{N \times d}}$$

## Example: multidimensional scaling

Suppose, we have a pairwise distance matrix for  $N$   $d$ -dimensional objects  $D \in \mathbb{R}^{N \times N}$ . Given this matrix, our goal is to recover the initial coordinates  $W_i \in \mathbb{R}^d$ ,  $i = 1, \dots, N$ .

$$L(W) = \sum_{i,j=1}^N (\|W_i - W_j\|_2^2 - D_{i,j})^2 \rightarrow \min_{W \in \mathbb{R}^{N \times d}}$$

Link to a nice visualization ♣, where one can see, that gradient-free methods handle this problem much slower, especially in higher dimensions.

### Question

Is it somehow connected with PCA?

## Example: Gradient Descent without gradient

Suppose we need to solve the following problem:

$$L(w) \rightarrow \min_{w \in \mathbb{R}^d}$$

## Example: Gradient Descent without gradient

Suppose we need to solve the following problem:

$$L(w) \rightarrow \min_{w \in \mathbb{R}^d}$$

with the Gradient Descent (GD) algorithm:

$$w_{k+1} = w_k - \alpha_k \nabla_w L(w_k)$$

## Example: Gradient Descent without gradient

Suppose we need to solve the following problem:

$$L(w) \rightarrow \min_{w \in \mathbb{R}^d}$$

with the Gradient Descent (GD) algorithm:

$$w_{k+1} = w_k - \alpha_k \nabla_w L(w_k)$$

Is it possible to replace  $\nabla_w L(w_k)$  using only zero-order information?

## Example: Gradient Descent without gradient

Suppose we need to solve the following problem:

$$L(w) \rightarrow \min_{w \in \mathbb{R}^d}$$

with the Gradient Descent (GD) algorithm:

$$w_{k+1} = w_k - \alpha_k \nabla_w L(w_k)$$

Is it possible to replace  $\nabla_w L(w_k)$  using only zero-order information?

Yes, but at a cost.

## Example: Gradient Descent without gradient

Suppose we need to solve the following problem:

$$L(w) \rightarrow \min_{w \in \mathbb{R}^d}$$

with the Gradient Descent (GD) algorithm:

$$w_{k+1} = w_k - \alpha_k \nabla_w L(w_k)$$

Is it possible to replace  $\nabla_w L(w_k)$  using only zero-order information?

Yes, but at a cost.

One can consider 2-point gradient estimator<sup>a</sup>  $G$ :

$$G = d \frac{L(w + \varepsilon v) - L(w - \varepsilon v)}{2\varepsilon} v,$$

where  $v$  is spherically symmetric.

---

<sup>a</sup>I suggest a nice presentation about gradient-free methods

## Example: Gradient Descent without gradient

Suppose we need to solve the following problem:

$$L(w) \rightarrow \min_{w \in \mathbb{R}^d}$$

with the Gradient Descent (GD) algorithm:

$$w_{k+1} = w_k - \alpha_k \nabla_w L(w_k)$$

Is it possible to replace  $\nabla_w L(w_k)$  using only zero-order information?

Yes, but at a cost.

One can consider 2-point gradient estimator<sup>a</sup>  $G$ :

$$G = d \frac{L(w + \varepsilon v) - L(w - \varepsilon v)}{2\varepsilon} v,$$

where  $v$  is spherically symmetric.

<sup>a</sup>I suggest a nice presentation about gradient-free methods

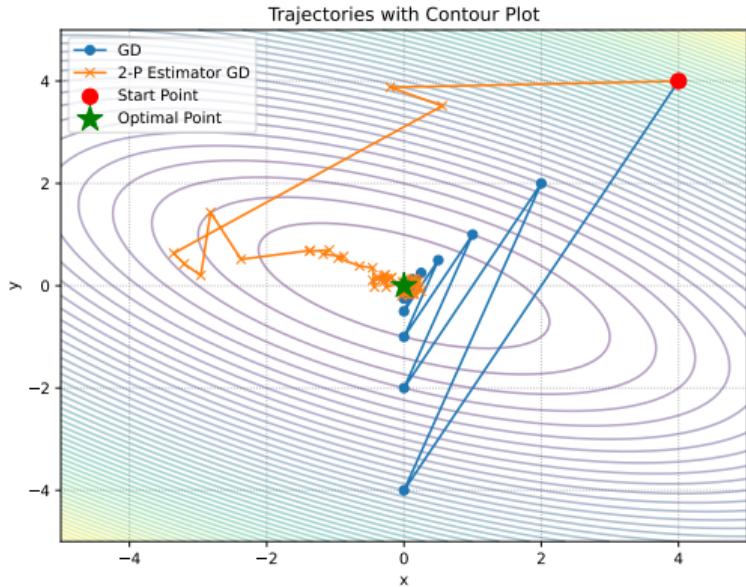


Figure 3: "Illustration of two-point estimator of Gradient Descent"

## Example: Gradient Descent without gradient

$$w_{k+1} = w_k - \alpha_k G$$

## Example: Gradient Descent without gradient

$$w_{k+1} = w_k - \alpha_k G$$

One can also consider the idea of finite differences:

$$G = \sum_{i=1}^d \frac{L(w + \varepsilon e_i) - L(w - \varepsilon e_i)}{2\varepsilon} e_i$$

Open In Colab ♣

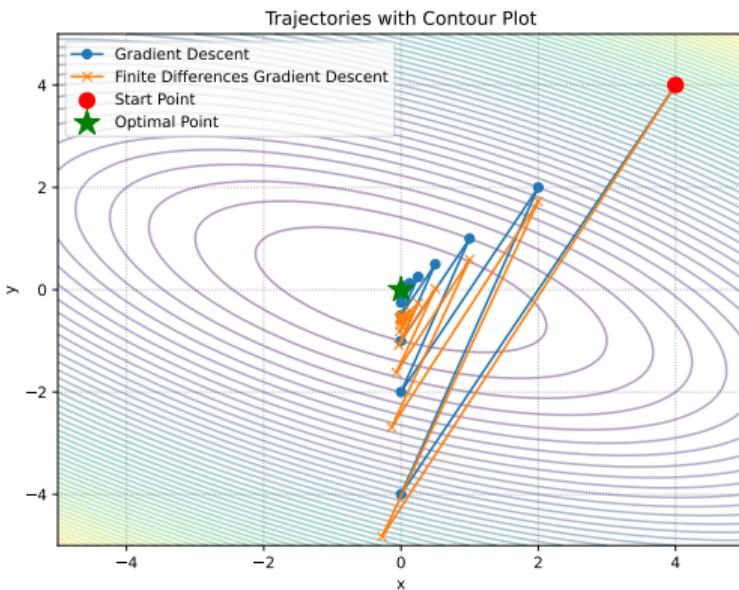


Figure 4: “Illustration of finite differences estimator of Gradient Descent”

## The curse of dimensionality for zero-order methods

$$\min_{x \in \mathbb{R}^n} f(x)$$

## The curse of dimensionality for zero-order methods

$$\min_{x \in \mathbb{R}^n} f(x)$$

GD:  $x_{k+1} = x_k - \alpha_k \nabla f(x_k)$

Zero order GD:  $x_{k+1} = x_k - \alpha_k G,$

where  $G$  is a 2-point or multi-point estimator of the gradient.

## The curse of dimensionality for zero-order methods

$$\min_{x \in \mathbb{R}^n} f(x)$$

GD:  $x_{k+1} = x_k - \alpha_k \nabla f(x_k)$

Zero order GD:  $x_{k+1} = x_k - \alpha_k G,$

where  $G$  is a 2-point or multi-point estimator of the gradient.

	$f(x)$ - smooth	$f(x)$ - smooth and convex	$f(x)$ - smooth and strongly convex
GD	$\ \nabla f(x_k)\ ^2 \approx \mathcal{O}\left(\frac{1}{k}\right)$	$f(x_k) - f^* \approx \mathcal{O}\left(\frac{1}{k}\right)$	$\ x_k - x^*\ ^2 \approx \mathcal{O}\left(\left(1 - \frac{\mu}{L}\right)^k\right)$
Zero order GD	$\ \nabla f(x_k)\ ^2 \approx \mathcal{O}\left(\frac{n}{k}\right)$	$f(x_k) - f^* \approx \mathcal{O}\left(\frac{n}{k}\right)$	$\ x_k - x^*\ ^2 \approx \mathcal{O}\left(\left(1 - \frac{\mu}{nL}\right)^k\right)$

## Finite differences

The naive approach to get approximate values of gradients is **Finite differences** approach. For each coordinate, one can calculate the partial derivative approximation:

$$\frac{\partial L}{\partial w_k}(w) \approx \frac{L(w + \varepsilon e_k) - L(w)}{\varepsilon}, \quad e_k = (0, \dots, \underset{k}{1}, \dots, 0)$$

---

<sup>1</sup>Linnainmaa S. The representation of the cumulative rounding error of an algorithm as a Taylor expansion of the local rounding errors. Master's Thesis (in Finnish), Univ. Helsinki, 1970.

## Finite differences

The naive approach to get approximate values of gradients is **Finite differences** approach. For each coordinate, one can calculate the partial derivative approximation:

$$\frac{\partial L}{\partial w_k}(w) \approx \frac{L(w + \varepsilon e_k) - L(w)}{\varepsilon}, \quad e_k = (0, \dots, \underset{k}{1}, \dots, 0)$$

### Question

If the time needed for one calculation of  $L(w)$  is  $T$ , what is the time needed for calculating  $\nabla_w L$  with this approach?

## Finite differences

The naive approach to get approximate values of gradients is **Finite differences** approach. For each coordinate, one can calculate the partial derivative approximation:

$$\frac{\partial L}{\partial w_k}(w) \approx \frac{L(w + \varepsilon e_k) - L(w)}{\varepsilon}, \quad e_k = (0, \dots, \underset{k}{1}, \dots, 0)$$

### Question

If the time needed for one calculation of  $L(w)$  is  $T$ , what is the time needed for calculating  $\nabla_w L$  with this approach?

**Answer**  $2dT$ , which is extremely long for the huge scale optimization. Moreover, this exact scheme is unstable, which means that you will have to choose between accuracy and stability.

## Finite differences

The naive approach to get approximate values of gradients is **Finite differences** approach. For each coordinate, one can calculate the partial derivative approximation:

$$\frac{\partial L}{\partial w_k}(w) \approx \frac{L(w + \varepsilon e_k) - L(w)}{\varepsilon}, \quad e_k = (0, \dots, \underset{k}{1}, \dots, 0)$$

### Question

If the time needed for one calculation of  $L(w)$  is  $T$ , what is the time needed for calculating  $\nabla_w L$  with this approach?

**Answer**  $2dT$ , which is extremely long for the huge scale optimization. Moreover, this exact scheme is unstable, which means that you will have to choose between accuracy and stability.

### Theorem

There is an algorithm to compute  $\nabla_w L$  in  $\mathcal{O}(T)$  operations.<sup>1</sup>

---

<sup>1</sup>Linnainmaa S. The representation of the cumulative rounding error of an algorithm as a Taylor expansion of the local rounding errors. Master's Thesis (in Finnish), Univ. Helsinki, 1970.

## Forward mode automatic differentiation

To dive deep into the idea of automatic differentiation we will consider a simple function for calculating derivatives:

$$L(w_1, w_2) = w_2 \log w_1 + \sqrt{w_2 \log w_1}$$

## Forward mode automatic differentiation

To dive deep into the idea of automatic differentiation we will consider a simple function for calculating derivatives:

$$L(w_1, w_2) = w_2 \log w_1 + \sqrt{w_2 \log w_1}$$

Let's draw a *computational graph* of this function:

$$L(w_1, w_2) = w_2 \log w_1 + \sqrt{w_2 \log w_1}$$

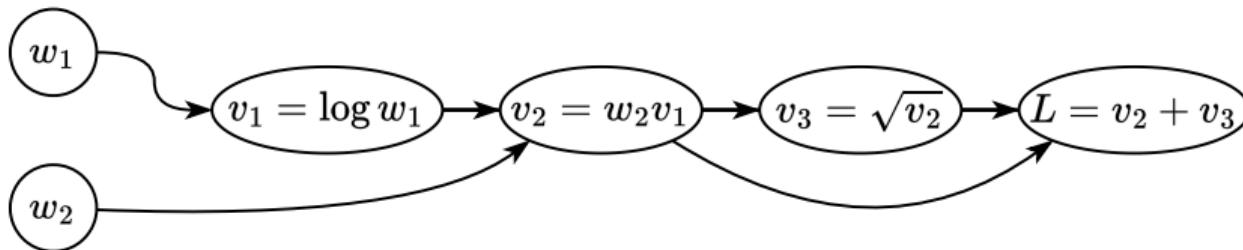


Figure 5: Illustration of computation graph of primitive arithmetic operations for the function  $L(w_1, w_2)$

## Forward mode automatic differentiation

To dive deep into the idea of automatic differentiation we will consider a simple function for calculating derivatives:

$$L(w_1, w_2) = w_2 \log w_1 + \sqrt{w_2 \log w_1}$$

Let's draw a *computational graph* of this function:

$$L(w_1, w_2) = w_2 \log w_1 + \sqrt{w_2 \log w_1}$$

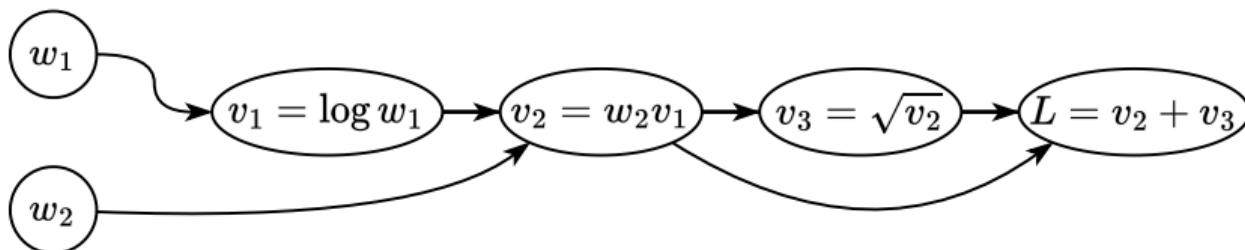


Figure 5: Illustration of computation graph of primitive arithmetic operations for the function  $L(w_1, w_2)$

Let's go from the beginning of the graph to the end and calculate the derivative  $\frac{\partial L}{\partial w_1}$ .

## Forward mode automatic differentiation

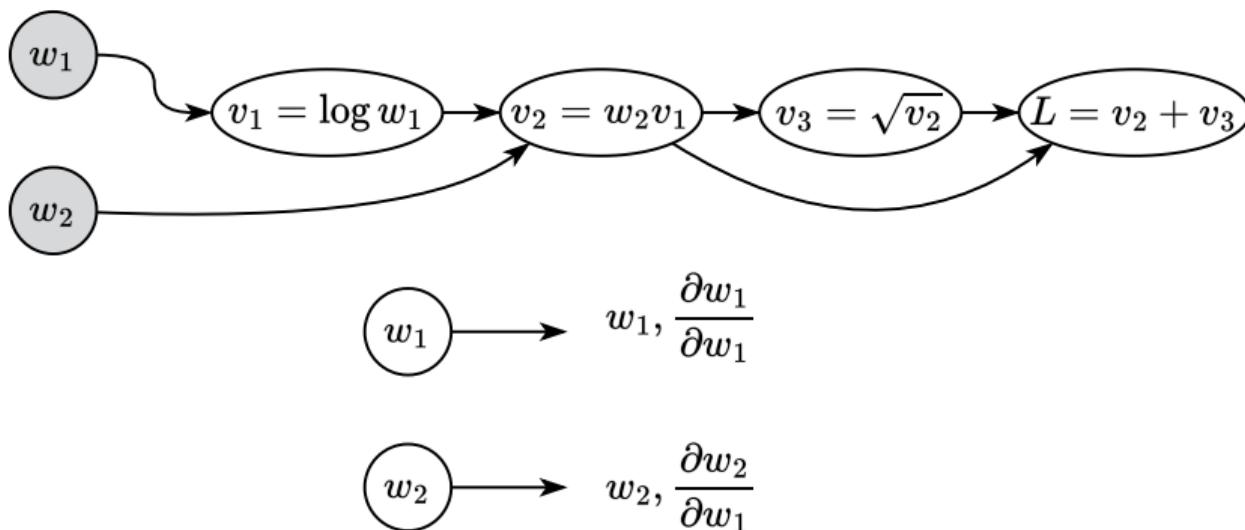


Figure 6: Illustration of forward mode automatic differentiation

## Function

$$w_1 = w_1, w_2 = w_2$$

## Forward mode automatic differentiation

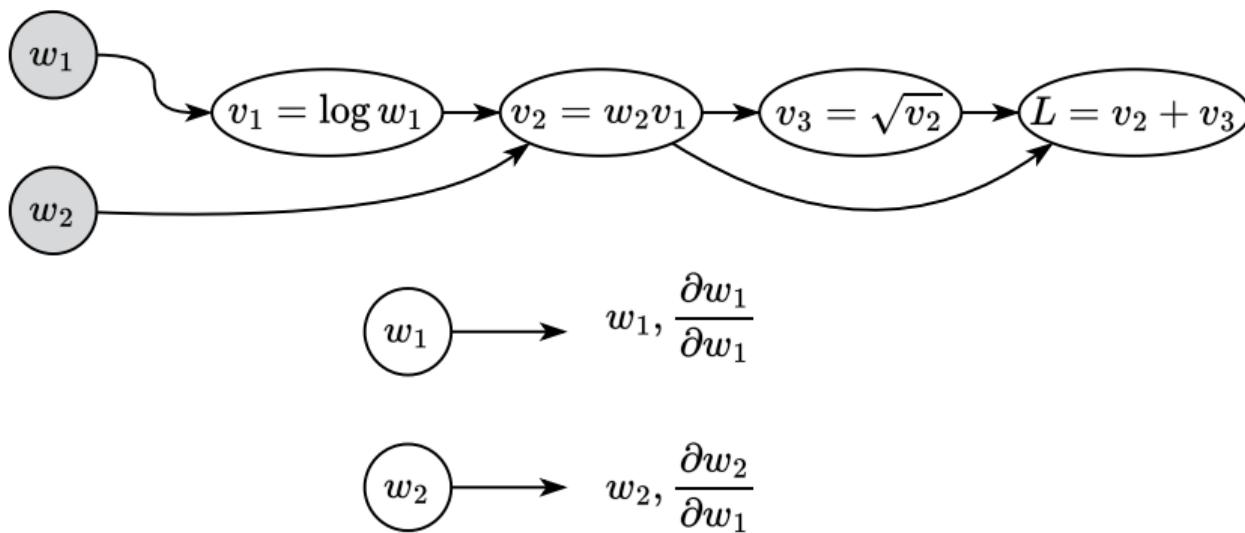


Figure 6: Illustration of forward mode automatic differentiation

### Function

$$w_1 = w_1, w_2 = w_2$$

### Derivative

$$\frac{\partial w_1}{\partial w_1} = 1, \frac{\partial w_2}{\partial w_1} = 0$$

## Forward mode automatic differentiation

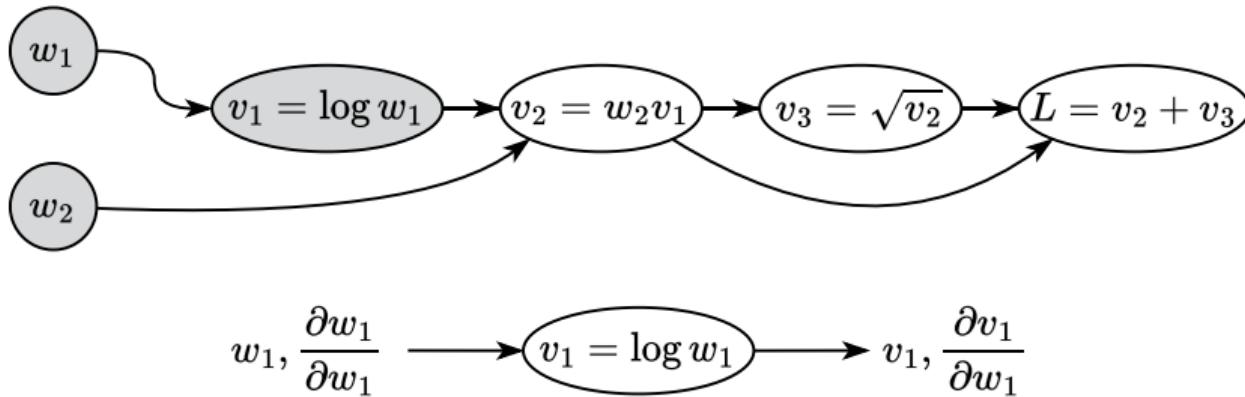


Figure 7: Illustration of forward mode automatic differentiation

## Forward mode automatic differentiation

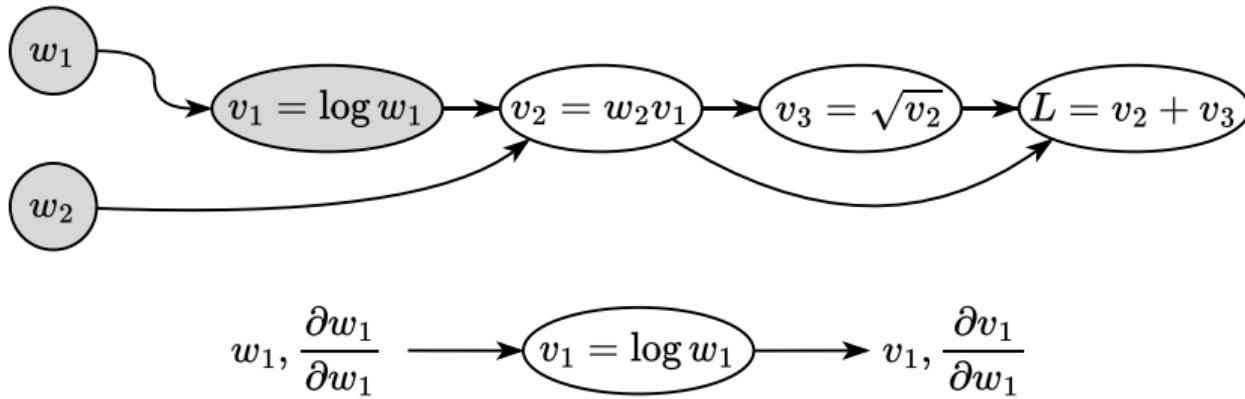


Figure 7: Illustration of forward mode automatic differentiation

### Function

$$v_1 = \log w_1$$

## Forward mode automatic differentiation

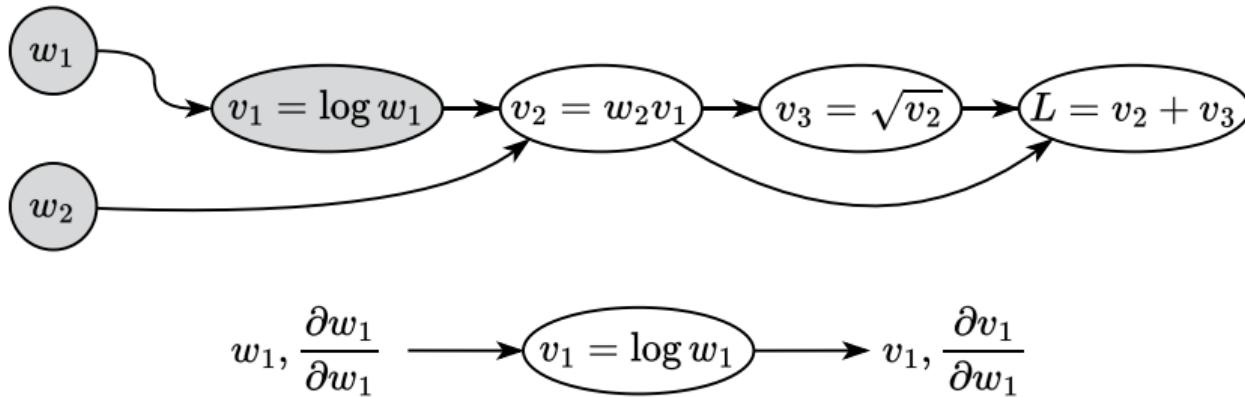


Figure 7: Illustration of forward mode automatic differentiation

### Function

$$v_1 = \log w_1$$

### Derivative

$$\frac{\partial v_1}{\partial w_1} = \frac{\partial v_1}{\partial w_1} \frac{\partial w_1}{\partial w_1} = \frac{1}{w_1} 1$$

## Forward mode automatic differentiation

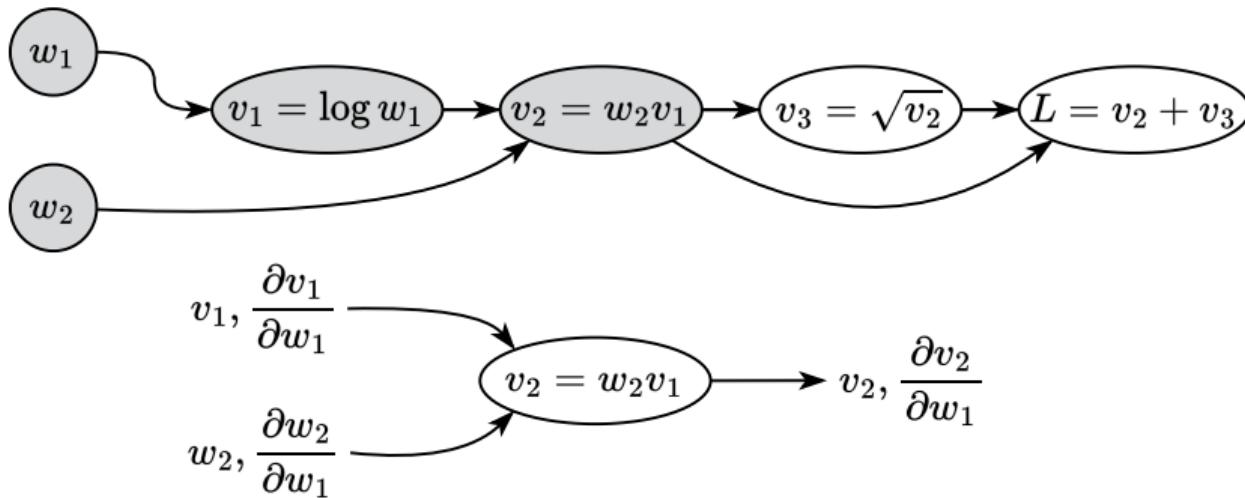


Figure 8: Illustration of forward mode automatic differentiation

## Forward mode automatic differentiation

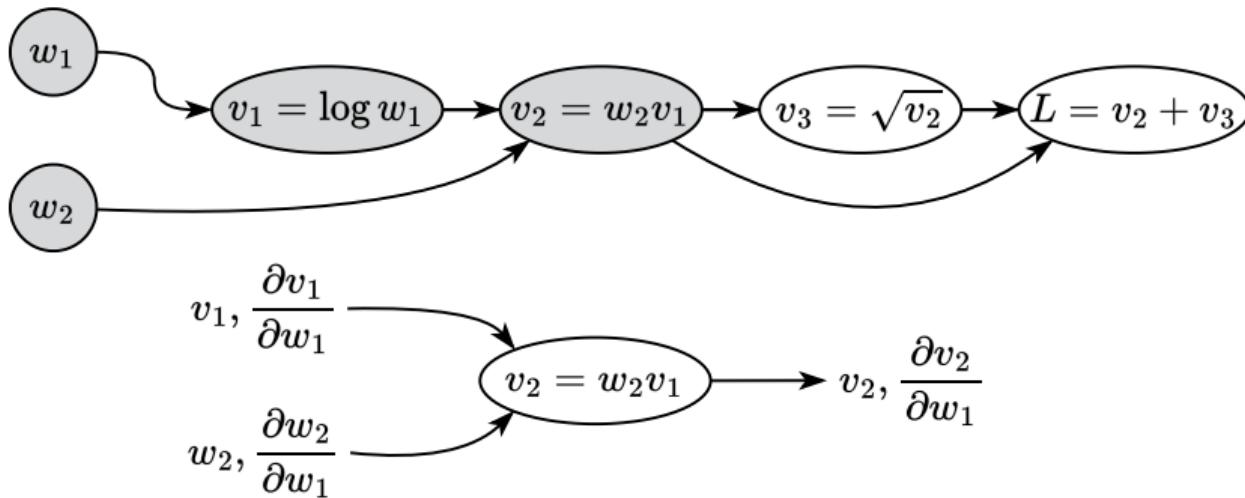


Figure 8: Illustration of forward mode automatic differentiation

### Function

$$v_2 = w_2 v_1$$

## Forward mode automatic differentiation

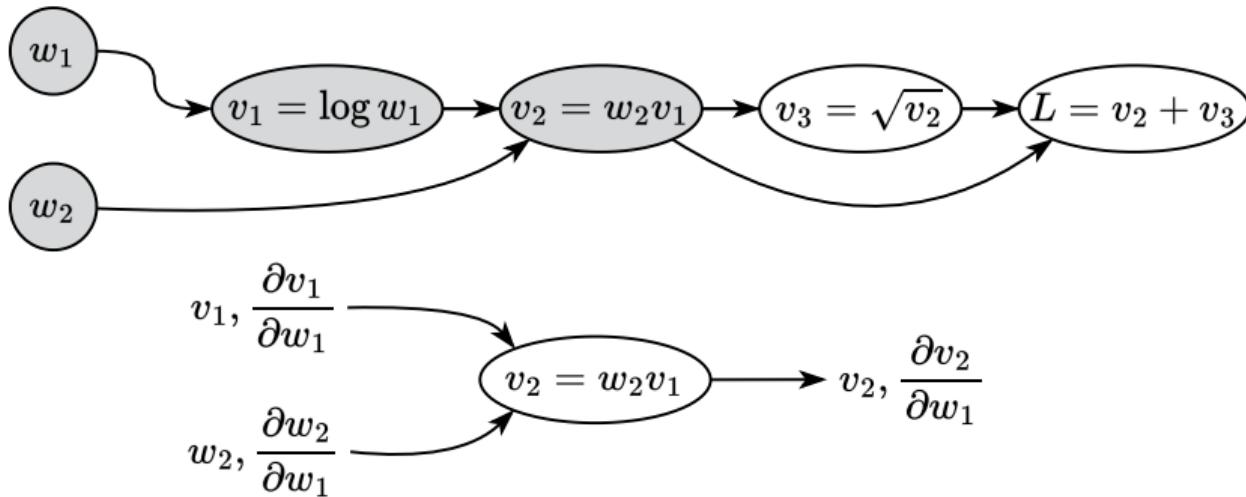


Figure 8: Illustration of forward mode automatic differentiation

Function

$$v_2 = w_2 v_1$$

Derivative

$$\frac{\partial v_2}{\partial w_1} = \frac{\partial v_2}{\partial v_1} \frac{\partial v_1}{\partial w_1} + \frac{\partial v_2}{\partial w_2} \frac{\partial w_2}{\partial w_1} = w_2 \frac{\partial v_1}{\partial w_1} + v_1 \frac{\partial w_2}{\partial w_1}$$

## Forward mode automatic differentiation

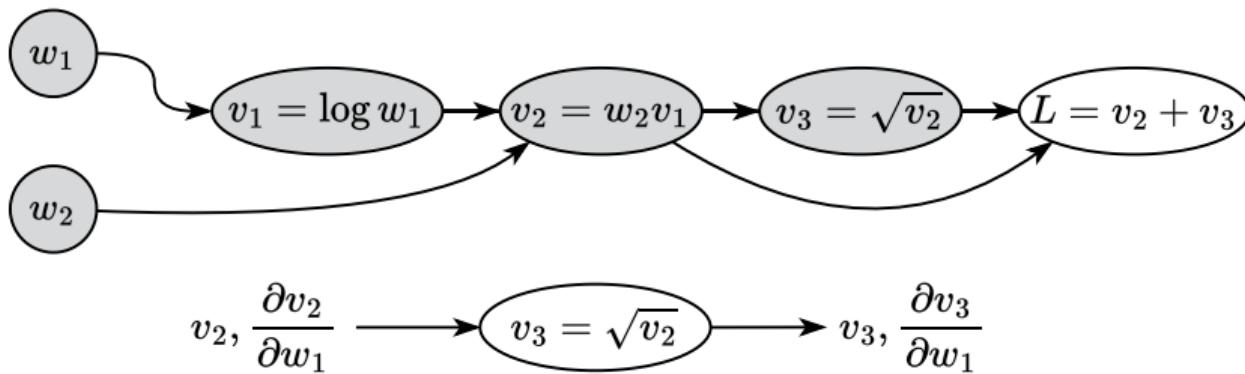


Figure 9: Illustration of forward mode automatic differentiation

## Forward mode automatic differentiation

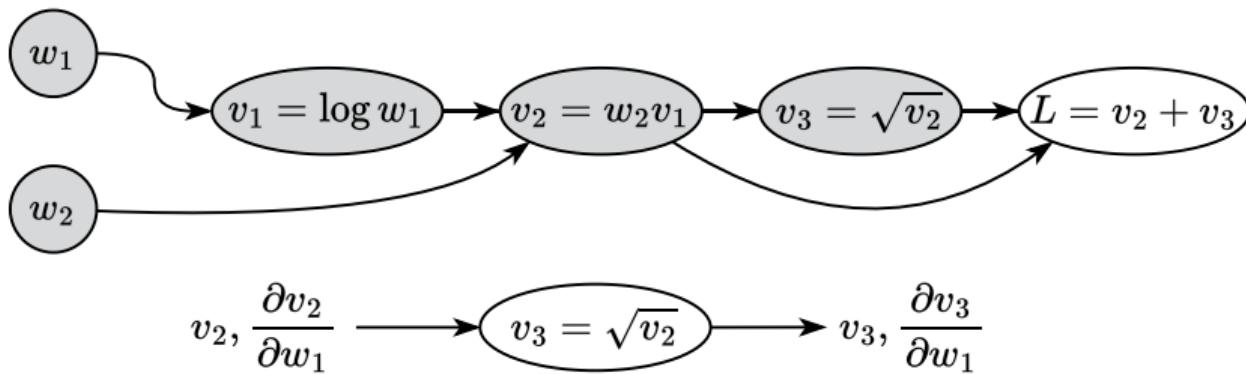


Figure 9: Illustration of forward mode automatic differentiation

### Function

$$v_3 = \sqrt{v_2}$$

## Forward mode automatic differentiation

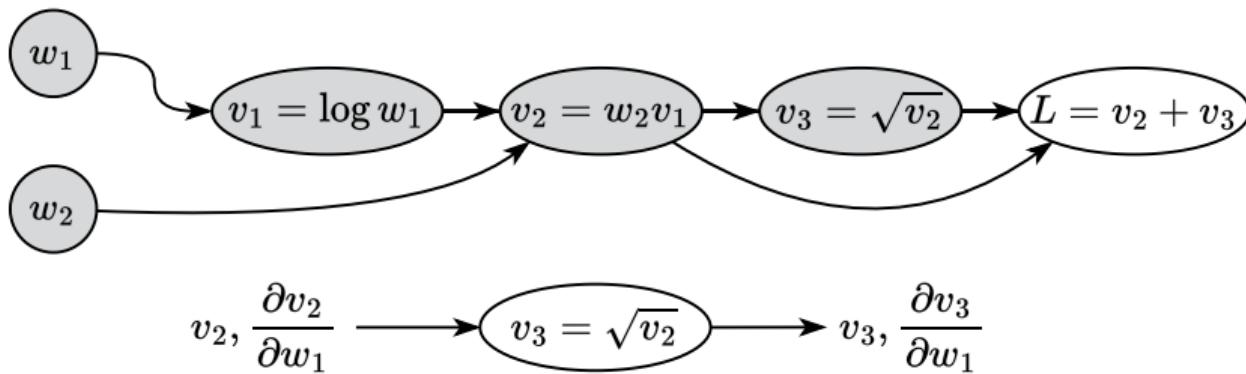


Figure 9: Illustration of forward mode automatic differentiation

### Function

$$v_3 = \sqrt{v_2}$$

### Derivative

$$\frac{\partial v_3}{\partial w_1} = \frac{\partial v_3}{\partial v_2} \frac{\partial v_2}{\partial w_1} = \frac{1}{2\sqrt{v_2}} \frac{\partial v_2}{\partial w_1}$$

## Forward mode automatic differentiation

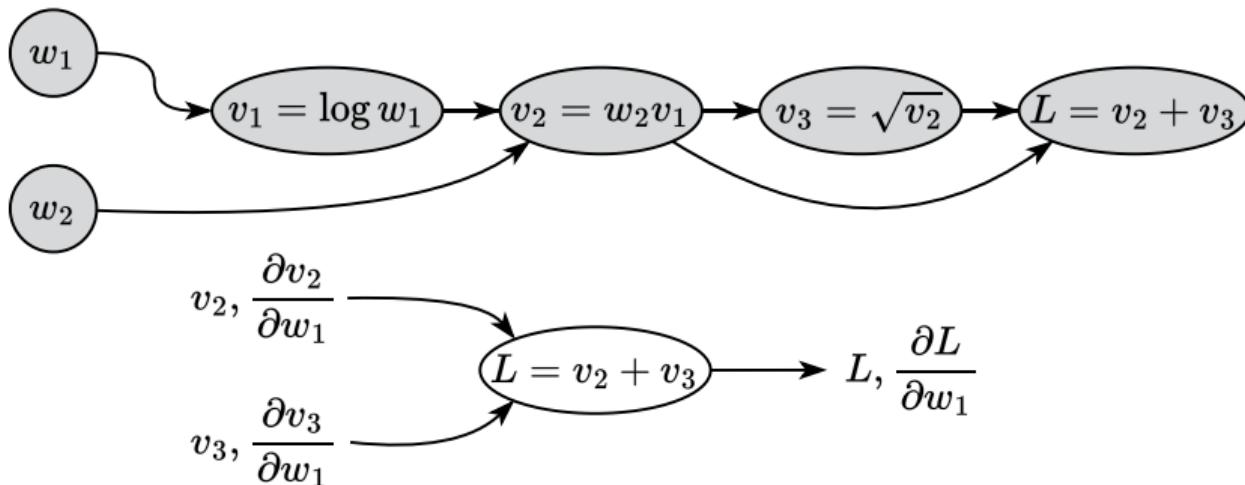


Figure 10: Illustration of forward mode automatic differentiation

## Forward mode automatic differentiation

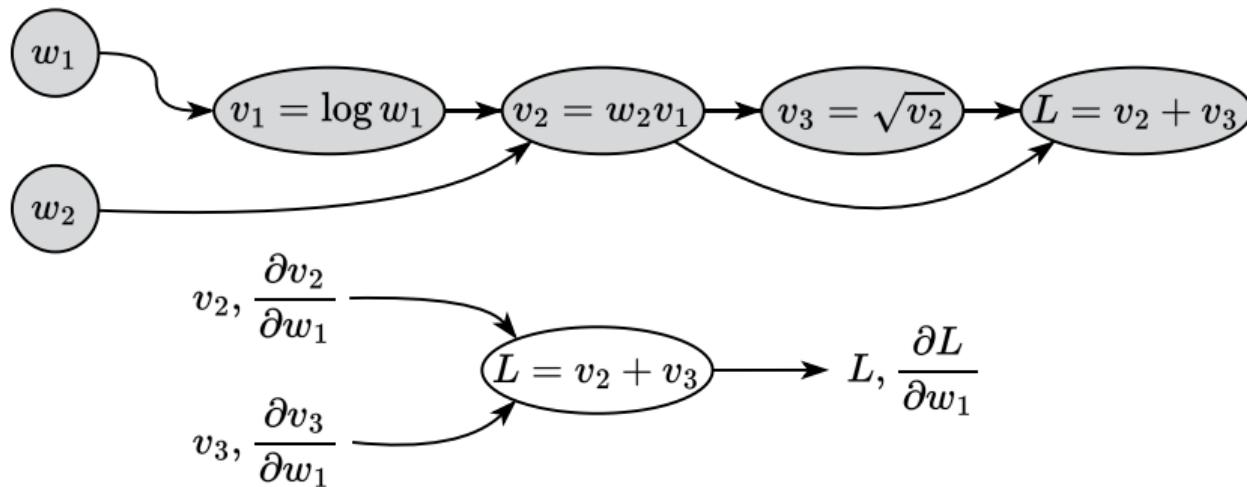


Figure 10: Illustration of forward mode automatic differentiation

### Function

$$L = v_2 + v_3$$

## Forward mode automatic differentiation

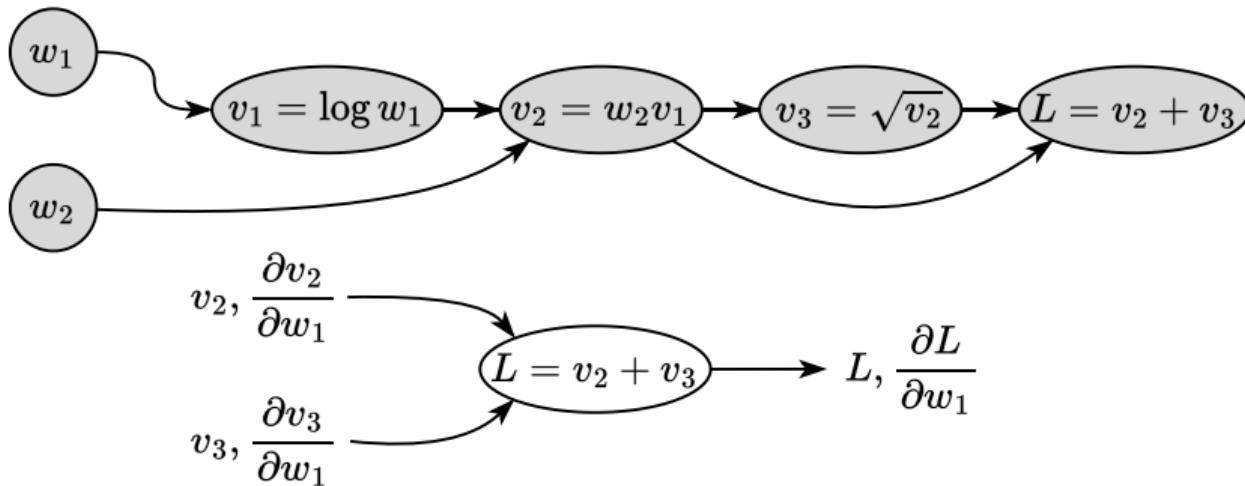


Figure 10: Illustration of forward mode automatic differentiation

### Function

$$L = v_2 + v_3$$

### Derivative

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial v_2} \frac{\partial v_2}{\partial w_1} + \frac{\partial L}{\partial v_3} \frac{\partial v_3}{\partial w_1} = 1 \frac{\partial v_2}{\partial w_1} + 1 \frac{\partial v_3}{\partial w_1}$$

Make the similar computations for  $\frac{\partial L}{\partial w_2}$

$$L(w_1, w_2) = w_2 \log w_1 + \sqrt{w_2 \log w_1}$$

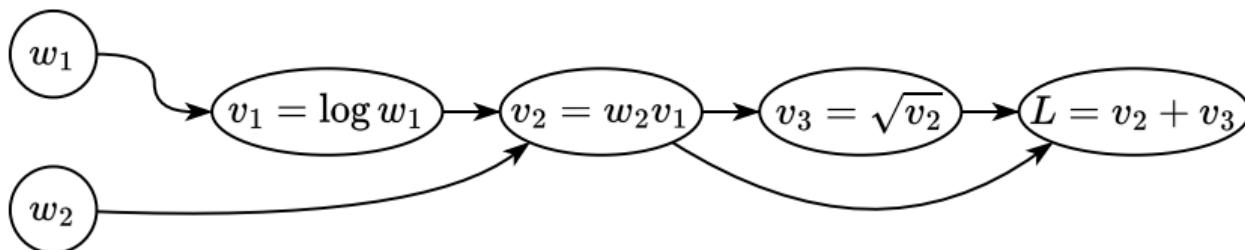


Figure 11: Illustration of computation graph of primitive arithmetic operations for the function  $L(w_1, w_2)$

## Forward mode automatic differentiation example

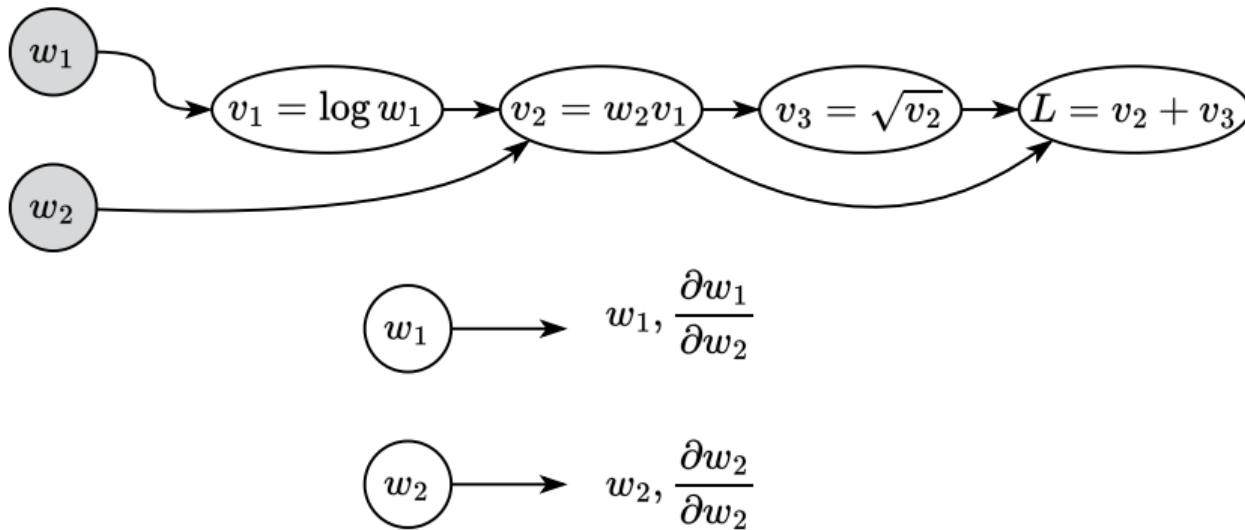


Figure 12: Illustration of forward mode automatic differentiation

### Function

$$w_1 = w_1, w_2 = w_2$$

### Derivative

$$\frac{\partial w_1}{\partial w_2} = 0, \frac{\partial w_2}{\partial w_2} = 1$$

## Forward mode automatic differentiation example

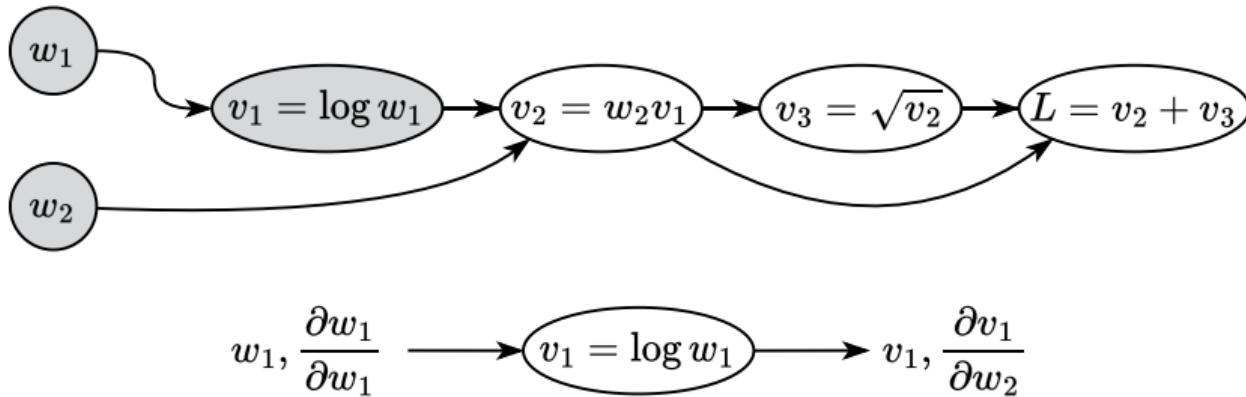


Figure 13: Illustration of forward mode automatic differentiation

### Function

$$v_1 = \log w_1$$

### Derivative

$$\frac{\partial v_1}{\partial w_2} = \frac{\partial v_1}{\partial w_2} \frac{\partial w_2}{\partial w_2} = 0 \cdot 1$$

## Forward mode automatic differentiation example

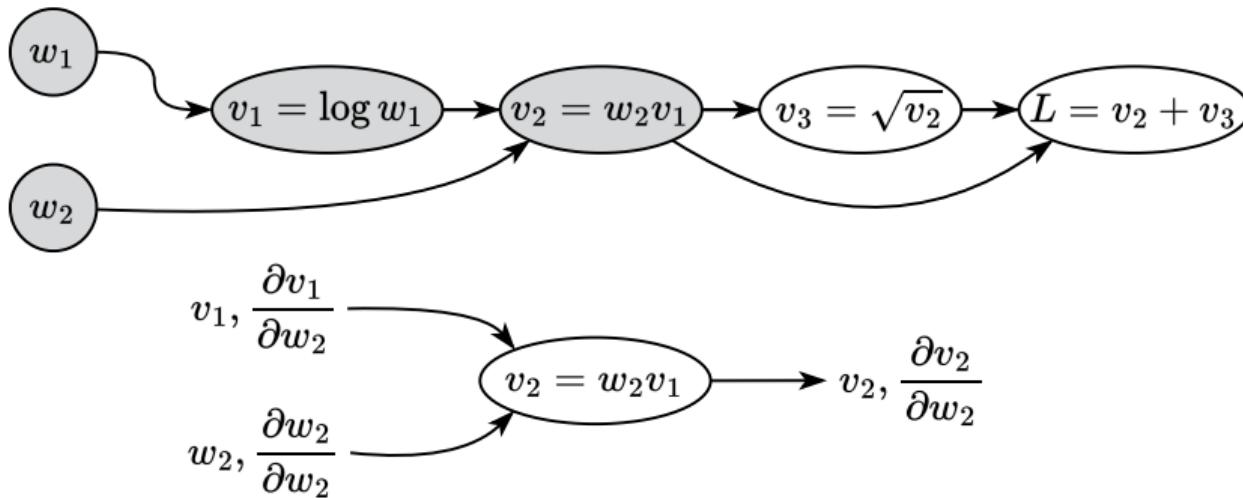


Figure 14: Illustration of forward mode automatic differentiation

### Function

$$v_2 = w_2 v_1$$

### Derivative

$$\frac{\partial v_2}{\partial w_2} = \frac{\partial v_2}{\partial v_1} \frac{\partial v_1}{\partial w_2} + \frac{\partial v_2}{\partial w_2} \frac{\partial w_2}{\partial w_2} = w_2 \frac{\partial v_1}{\partial w_2} + v_1 \frac{\partial w_2}{\partial w_2}$$

## Forward mode automatic differentiation example

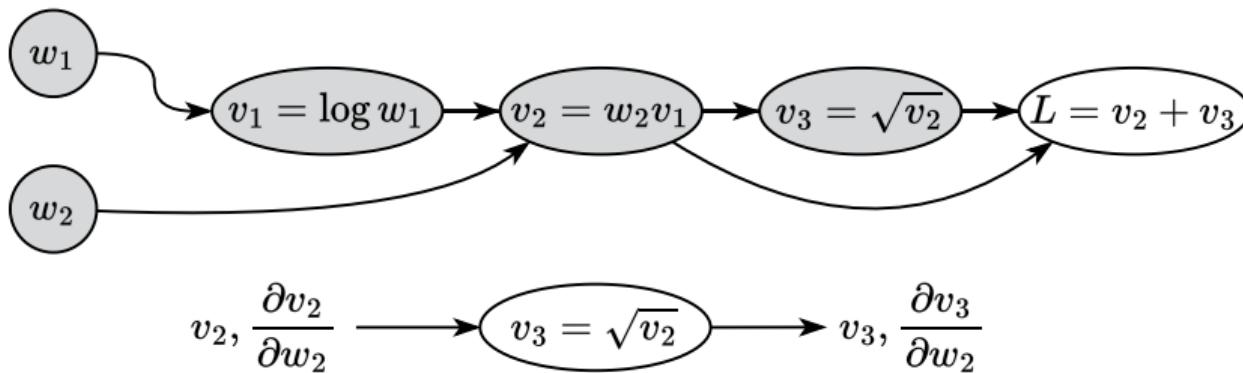


Figure 15: Illustration of forward mode automatic differentiation

### Function

$$v_3 = \sqrt{v_2}$$

### Derivative

$$\frac{\partial v_3}{\partial w_2} = \frac{\partial v_3}{\partial v_2} \frac{\partial v_2}{\partial w_2} = \frac{1}{2\sqrt{v_2}} \frac{\partial v_2}{\partial w_2}$$

## Forward mode automatic differentiation example

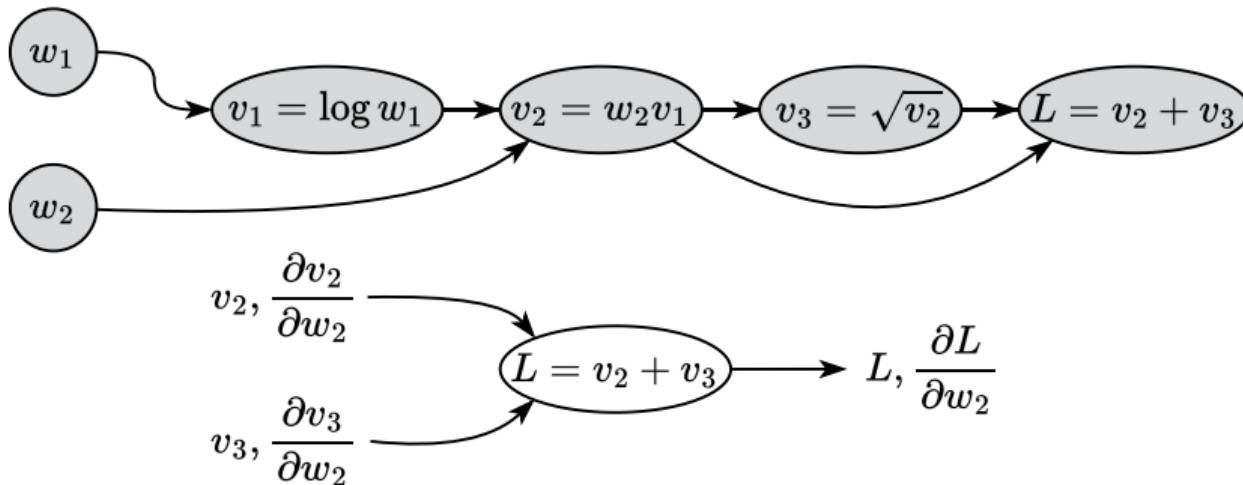


Figure 16: Illustration of forward mode automatic differentiation

### Function

$$L = v_2 + v_3$$

### Derivative

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial v_2} \frac{\partial v_2}{\partial w_2} + \frac{\partial L}{\partial v_3} \frac{\partial v_3}{\partial w_2} = 1 \frac{\partial v_2}{\partial w_2} + 1 \frac{\partial v_3}{\partial w_2}$$

## Forward mode automatic differentiation algorithm

Suppose, we have a computational graph  $v_i, i \in [1; N]$ .

Our goal is to calculate the derivative of the output of this graph with respect to some input variable  $w_k$ ,

i.e.  $\frac{\partial v_N}{\partial w_k}$ . This idea implies propagation of the gradient

with respect to the input variable from start to end, that is why we can introduce the notation:

## Forward mode automatic differentiation algorithm

Suppose, we have a computational graph  $v_i, i \in [1; N]$ .

Our goal is to calculate the derivative of the output of this graph with respect to some input variable  $w_k$ ,

i.e.  $\frac{\partial v_N}{\partial w_k}$ . This idea implies propagation of the gradient with respect to the input variable from start to end, that

is why we can introduce the notation:

$$\overline{v_i} = \frac{\partial v_i}{\partial w_k}$$

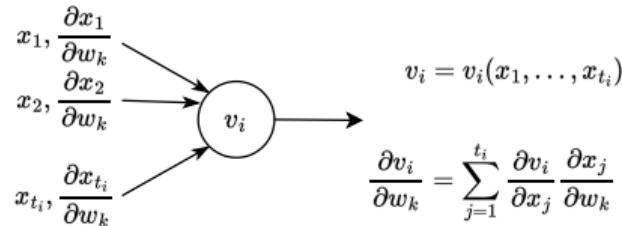


Figure 17: Illustration of forward chain rule to calculate the derivative of the function  $L$  with respect to  $w_k$ .

## Forward mode automatic differentiation algorithm

Suppose, we have a computational graph  $v_i, i \in [1; N]$ .

- For  $i = 1, \dots, N$ :

Our goal is to calculate the derivative of the output of this graph with respect to some input variable  $w_k$ ,

i.e.  $\frac{\partial v_N}{\partial w_k}$ . This idea implies propagation of the gradient with respect to the input variable from start to end, that

is why we can introduce the notation:

$$\overline{v_i} = \frac{\partial v_i}{\partial w_k}$$

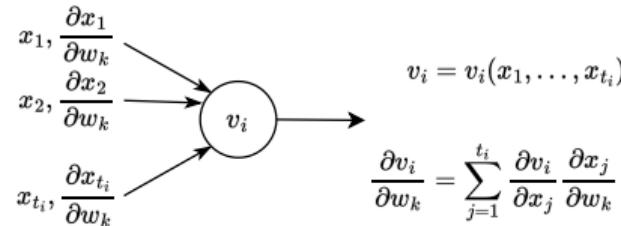


Figure 17: Illustration of forward chain rule to calculate the derivative of the function  $L$  with respect to  $w_k$ .

## Forward mode automatic differentiation algorithm

Suppose, we have a computational graph  $v_i, i \in [1; N]$ .

Our goal is to calculate the derivative of the output of this graph with respect to some input variable  $w_k$ ,

i.e.  $\frac{\partial v_N}{\partial w_k}$ . This idea implies propagation of the gradient

with respect to the input variable from start to end, that is why we can introduce the notation:

$$\overline{v_i} = \frac{\partial v_i}{\partial w_k}$$

- For  $i = 1, \dots, N$ :
  - Compute  $v_i$  as a function of its parents (inputs)  $x_1, \dots, x_{t_i}$ :

$$v_i = v_i(x_1, \dots, x_{t_i})$$

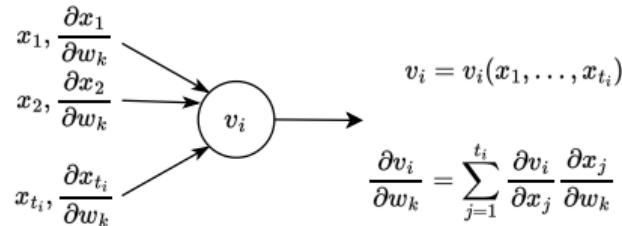


Figure 17: Illustration of forward chain rule to calculate the derivative of the function  $L$  with respect to  $w_k$ .

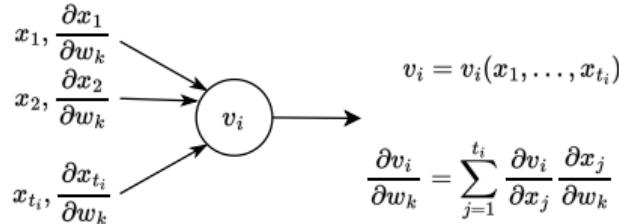
## Forward mode automatic differentiation algorithm

Suppose, we have a computational graph  $v_i, i \in [1; N]$ .

Our goal is to calculate the derivative of the output of this graph with respect to some input variable  $w_k$ ,

i.e.  $\frac{\partial v_N}{\partial w_k}$ . This idea implies propagation of the gradient with respect to the input variable from start to end, that is why we can introduce the notation:

$$\bar{v}_i = \frac{\partial v_i}{\partial w_k}$$



- For  $i = 1, \dots, N$ :
  - Compute  $v_i$  as a function of its parents (inputs)  $x_1, \dots, x_{t_i}$ :
- Compute the derivative  $\bar{v}_i$  using the forward chain rule:

$$\bar{v}_i = \sum_{j=1}^{t_i} \frac{\partial v_i}{\partial x_j} \frac{\partial x_j}{\partial w_k}$$

Figure 17: Illustration of forward chain rule to calculate the derivative of the function  $L$  with respect to  $w_k$ .

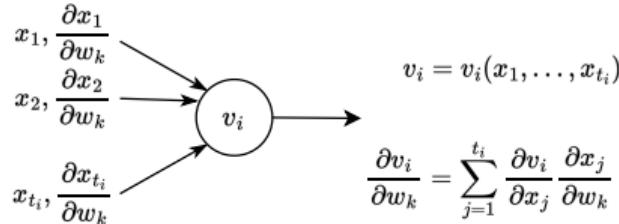
## Forward mode automatic differentiation algorithm

Suppose, we have a computational graph  $v_i, i \in [1; N]$ .

Our goal is to calculate the derivative of the output of this graph with respect to some input variable  $w_k$ ,

i.e.  $\frac{\partial v_N}{\partial w_k}$ . This idea implies propagation of the gradient with respect to the input variable from start to end, that is why we can introduce the notation:

$$\bar{v}_i = \frac{\partial v_i}{\partial w_k}$$



- For  $i = 1, \dots, N$ :
  - Compute  $v_i$  as a function of its parents (inputs)  $x_1, \dots, x_{t_i}$ :
- Compute the derivative  $\bar{v}_i$  using the forward chain rule:

$$\bar{v}_i = \sum_{j=1}^{t_i} \frac{\partial v_i}{\partial x_j} \frac{\partial x_j}{\partial w_k}$$

Figure 17: Illustration of forward chain rule to calculate the derivative of the function  $L$  with respect to  $w_k$ .

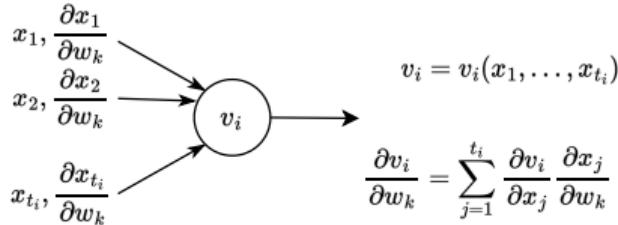
## Forward mode automatic differentiation algorithm

Suppose, we have a computational graph  $v_i, i \in [1; N]$ .

Our goal is to calculate the derivative of the output of this graph with respect to some input variable  $w_k$ ,

i.e.  $\frac{\partial v_N}{\partial w_k}$ . This idea implies propagation of the gradient with respect to the input variable from start to end, that is why we can introduce the notation:

$$\bar{v}_i = \frac{\partial v_i}{\partial w_k}$$



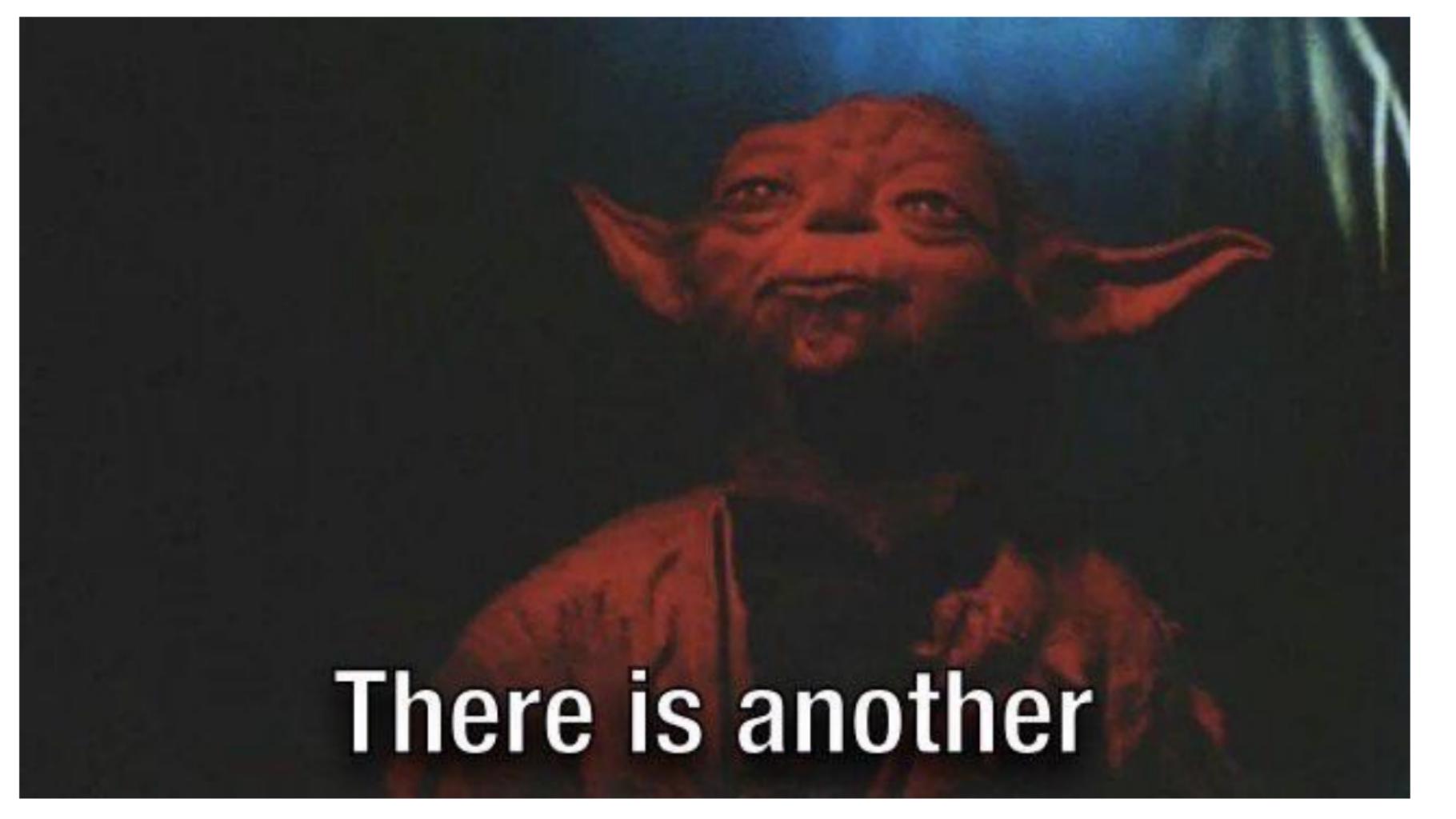
- For  $i = 1, \dots, N$ :
- Compute  $v_i$  as a function of its parents (inputs)  $x_1, \dots, x_{t_i}$ :
- Compute the derivative  $\bar{v}_i$  using the forward chain rule:

$$v_i = v_i(x_1, \dots, x_{t_i})$$

$$\bar{v}_i = \sum_{j=1}^{t_i} \frac{\partial v_i}{\partial x_j} \frac{\partial x_j}{\partial w_k}$$

Note, that this approach does not require storing all intermediate computations, but one can see, that for calculating the derivative  $\frac{\partial L}{\partial w_k}$  we need  $\mathcal{O}(T)$  operations. This means, that for the whole gradient, we need  $d\mathcal{O}(T)$  operations, which is the same as for finite differences, but we do not have stability issues, or inaccuracies now (the formulas above are exact).

Figure 17: Illustration of forward chain rule to calculate the derivative of the function  $L$  with respect to  $w_k$ .

A portrait of Yoda from Star Wars, looking slightly upwards with his arms outstretched.

There is another

## Backward mode automatic differentiation

We will consider the same function with a computational graph:

$$L(w_1, w_2) = w_2 \log w_1 + \sqrt{w_2 \log w_1}$$

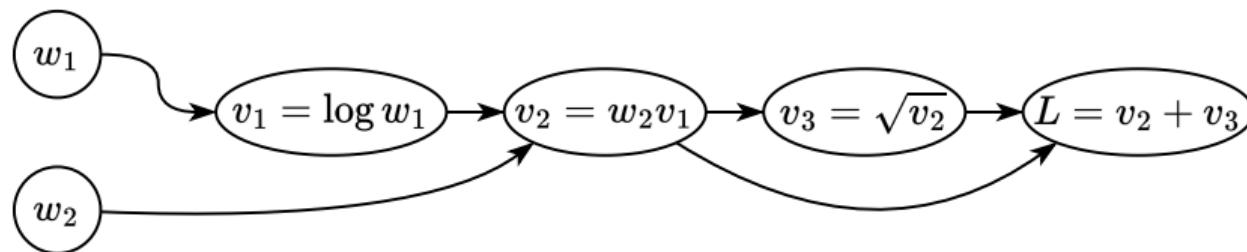


Figure 18: Illustration of computation graph of primitive arithmetic operations for the function  $L(w_1, w_2)$

## Backward mode automatic differentiation

We will consider the same function with a computational graph:

$$L(w_1, w_2) = w_2 \log w_1 + \sqrt{w_2 \log w_1}$$

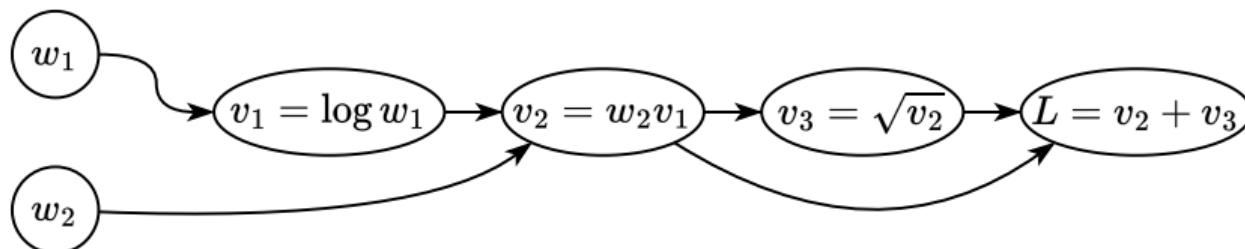


Figure 18: Illustration of computation graph of primitive arithmetic operations for the function  $L(w_1, w_2)$

Assume, that we have some values of the parameters  $w_1, w_2$  and we have already performed a forward pass (i.e. single propagation through the computational graph from left to right). Suppose, also, that we somehow saved all intermediate values of  $v_i$ . Let's go from the end of the graph to the beginning and calculate the derivatives  $\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}$ :

## Backward mode automatic differentiation example

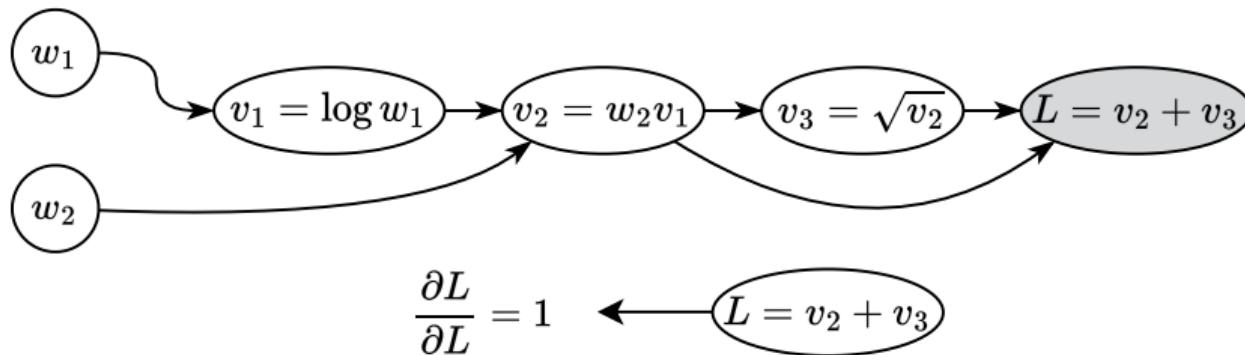


Figure 19: Illustration of backward mode automatic differentiation

## Backward mode automatic differentiation example

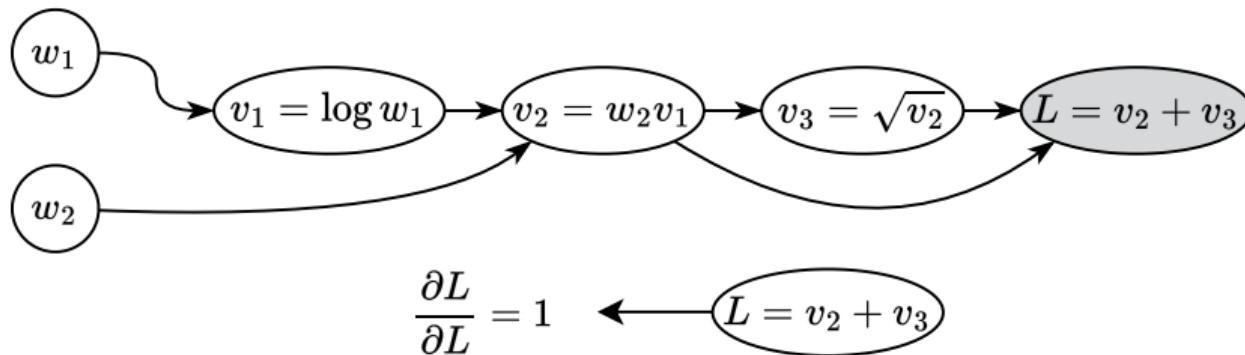


Figure 19: Illustration of backward mode automatic differentiation

Derivatives

## Backward mode automatic differentiation example

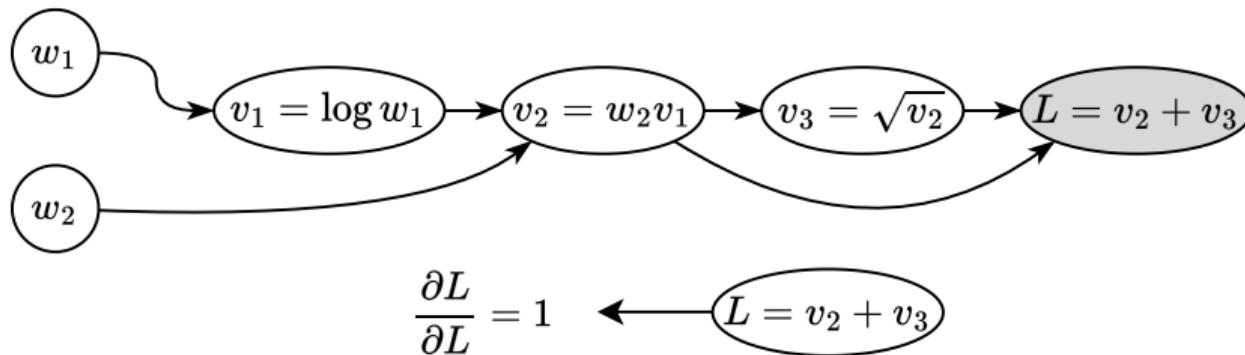


Figure 19: Illustration of backward mode automatic differentiation

## Derivatives

$$\frac{\partial L}{\partial L} = 1$$

## Backward mode automatic differentiation example

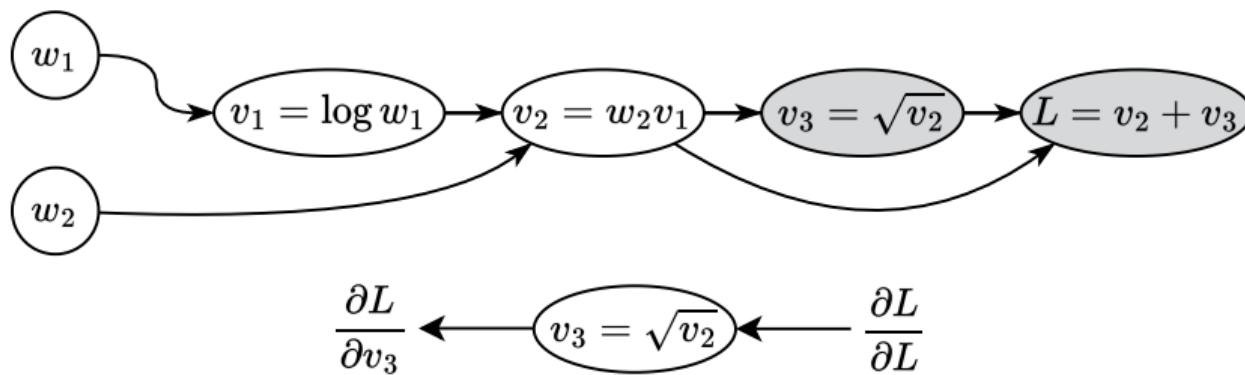


Figure 20: Illustration of backward mode automatic differentiation

## Backward mode automatic differentiation example

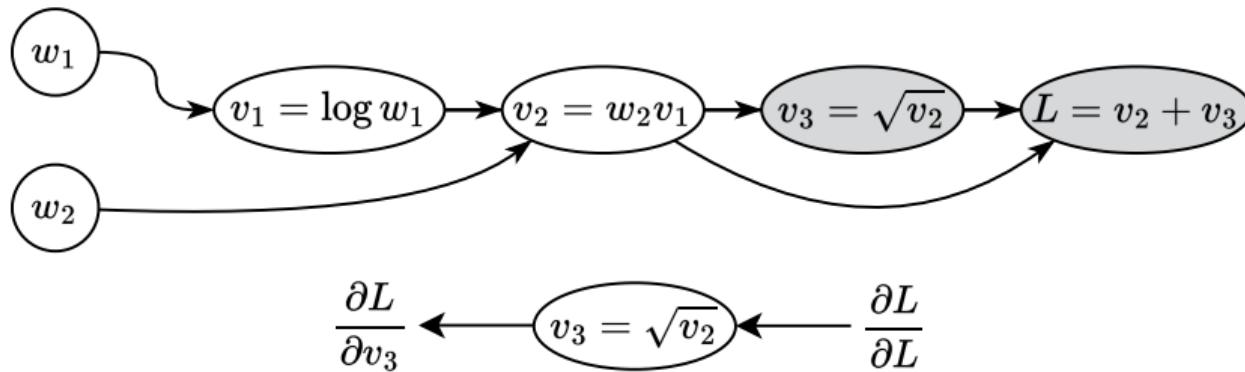


Figure 20: Illustration of backward mode automatic differentiation

Derivatives

## Backward mode automatic differentiation example

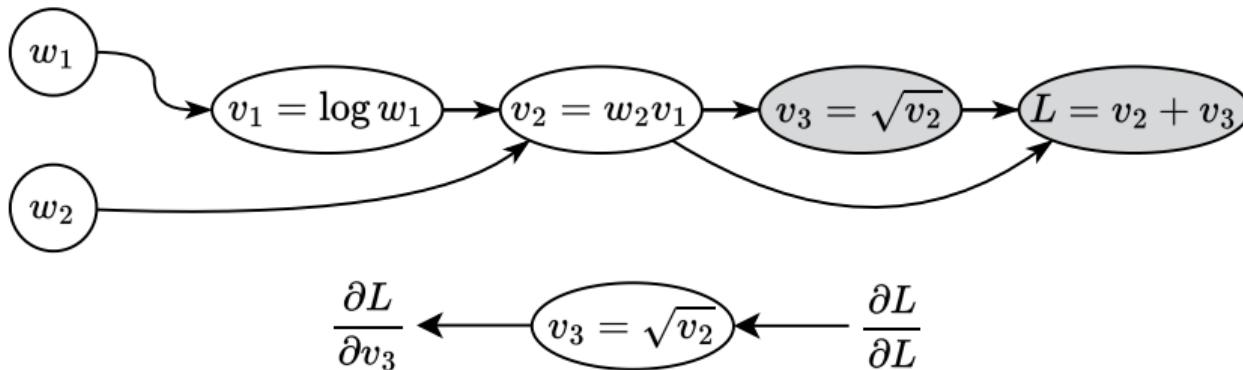


Figure 20: Illustration of backward mode automatic differentiation

## Derivatives

$$\begin{aligned}\frac{\partial L}{\partial v_3} &= \frac{\partial L}{\partial L} \frac{\partial L}{\partial v_3} \\ &= \frac{\partial L}{\partial L} 1\end{aligned}$$

## Backward mode automatic differentiation example

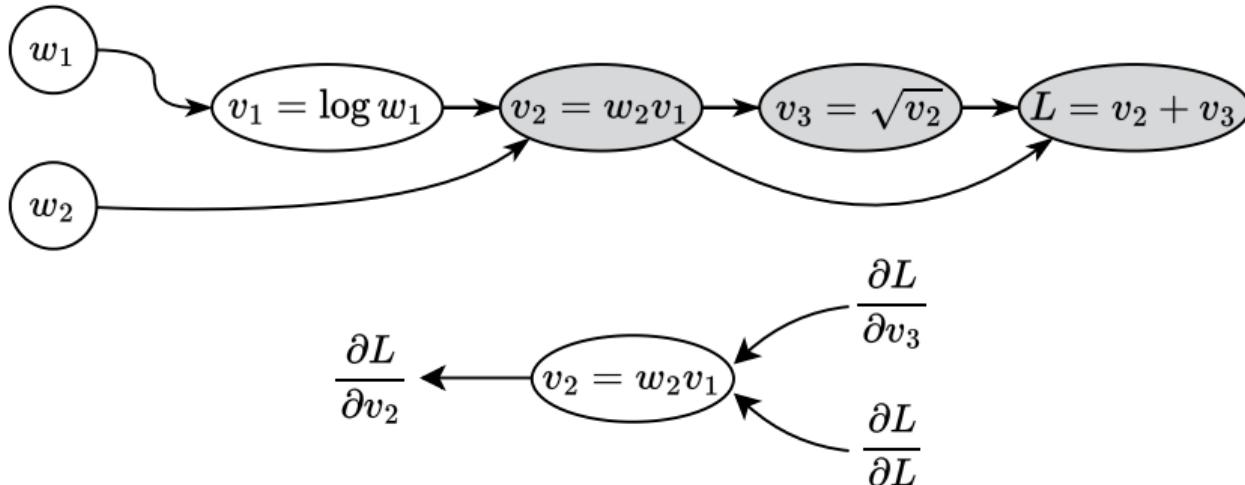


Figure 21: Illustration of backward mode automatic differentiation

## Backward mode automatic differentiation example

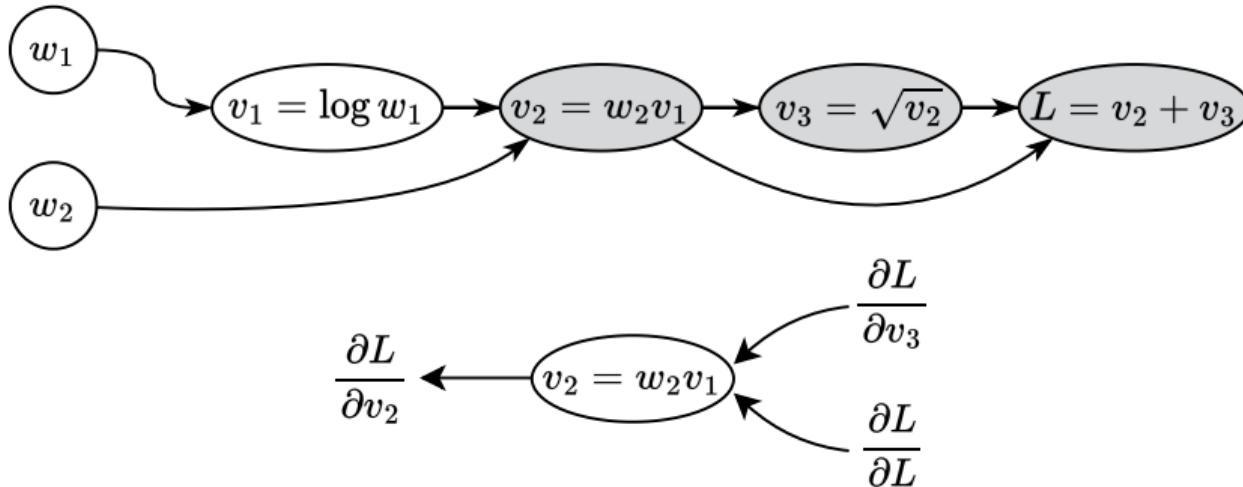


Figure 21: Illustration of backward mode automatic differentiation

Derivatives

## Backward mode automatic differentiation example

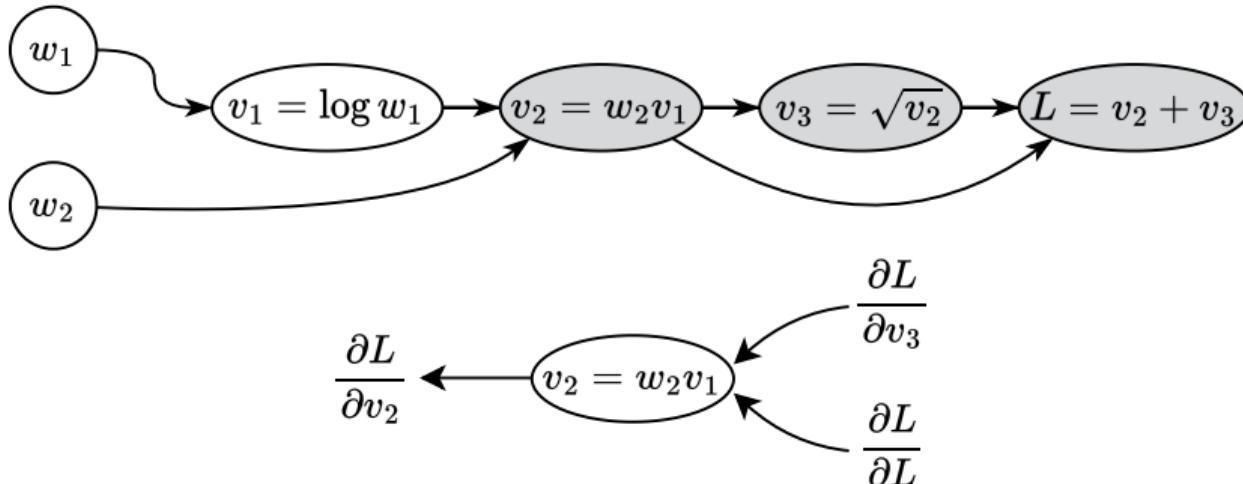


Figure 21: Illustration of backward mode automatic differentiation

## Derivatives

$$\begin{aligned}\frac{\partial L}{\partial v_2} &= \frac{\partial L}{\partial v_3} \frac{\partial v_3}{\partial v_2} + \frac{\partial L}{\partial L} \frac{\partial L}{\partial v_2} \\ &= \frac{\partial L}{\partial v_3} \frac{1}{2\sqrt{v_2}} + \frac{\partial L}{\partial L} 1\end{aligned}$$

## Backward mode automatic differentiation example

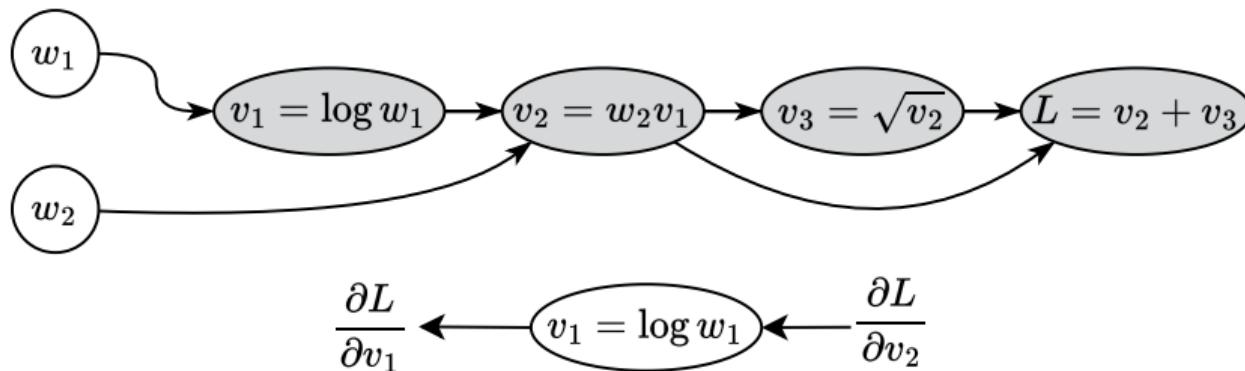


Figure 22: Illustration of backward mode automatic differentiation

## Backward mode automatic differentiation example

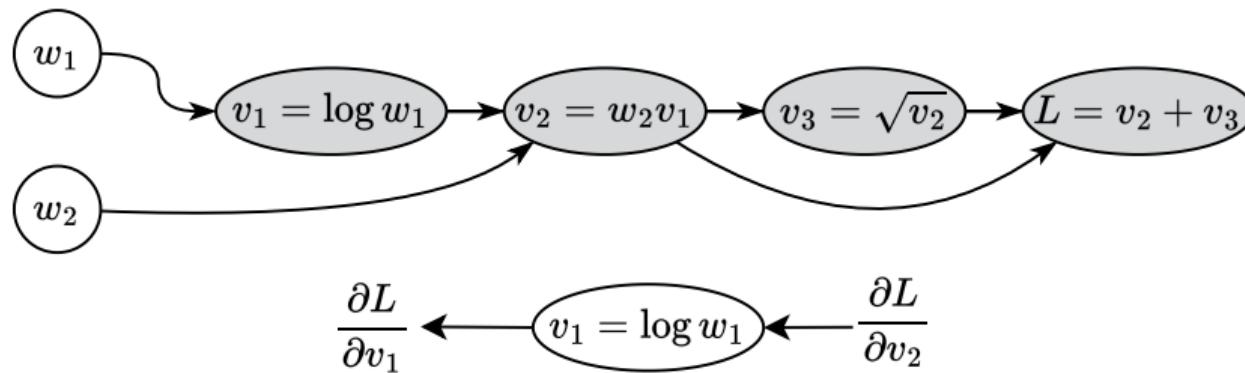


Figure 22: Illustration of backward mode automatic differentiation

Derivatives

## Backward mode automatic differentiation example

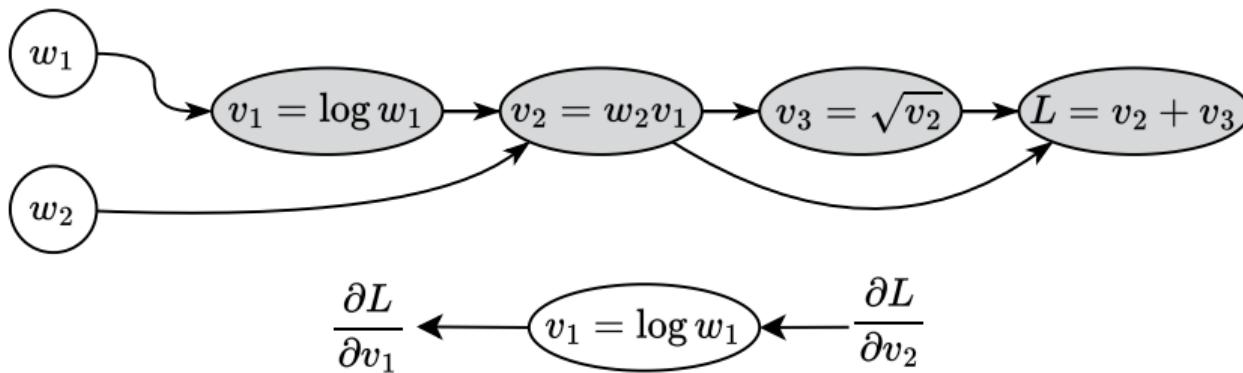


Figure 22: Illustration of backward mode automatic differentiation

## Derivatives

$$\begin{aligned}\frac{\partial L}{\partial v_1} &= \frac{\partial L}{\partial v_2} \frac{\partial v_2}{\partial v_1} \\ &= \frac{\partial L}{\partial v_2} w_2\end{aligned}$$

## Backward mode automatic differentiation example

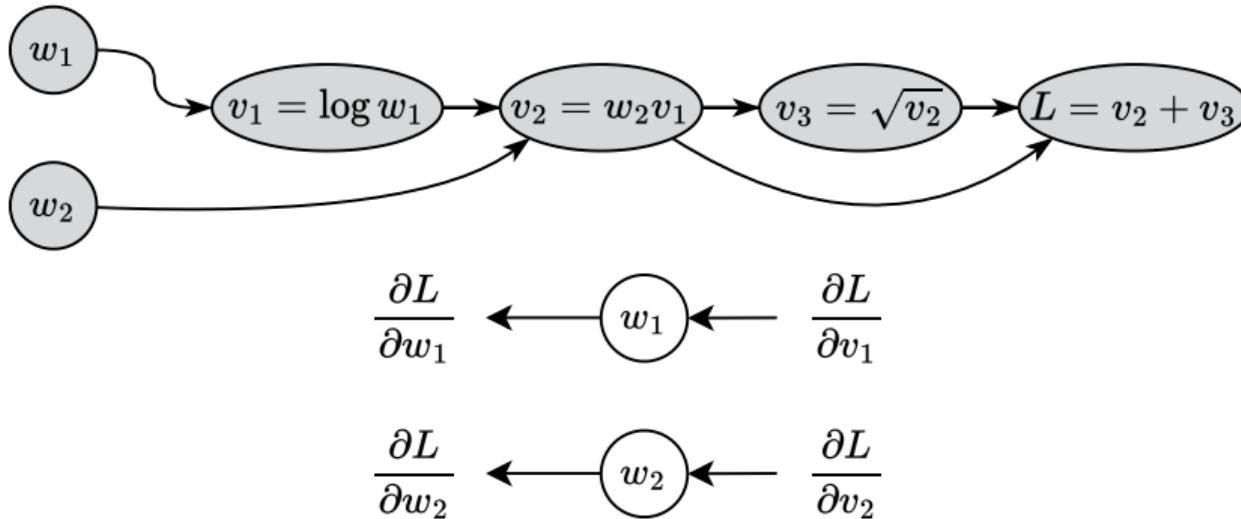


Figure 23: Illustration of backward mode automatic differentiation

## Backward mode automatic differentiation example

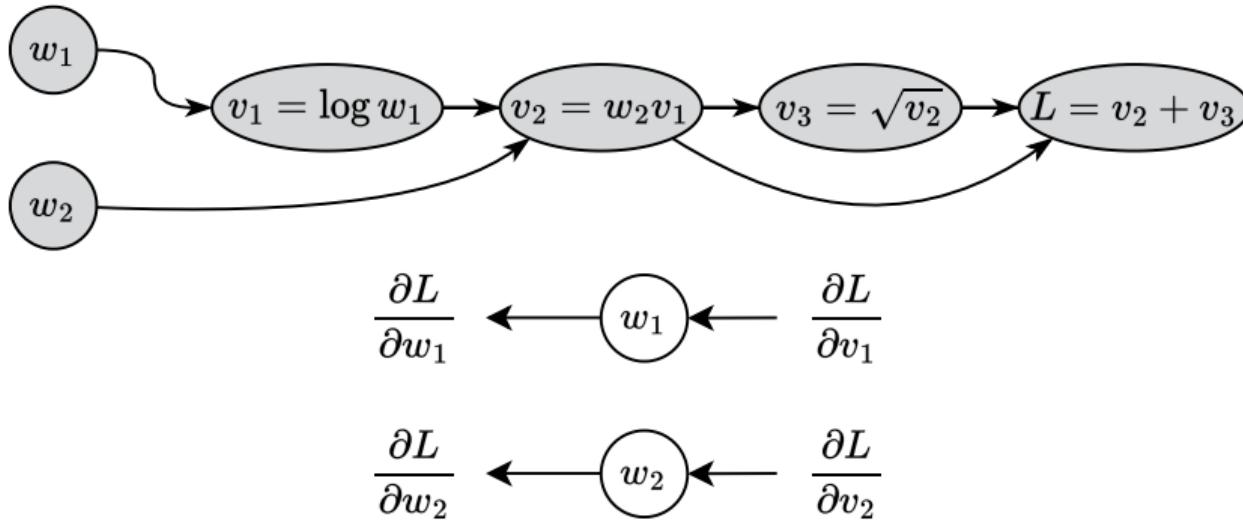


Figure 23: Illustration of backward mode automatic differentiation

Derivatives

## Backward mode automatic differentiation example

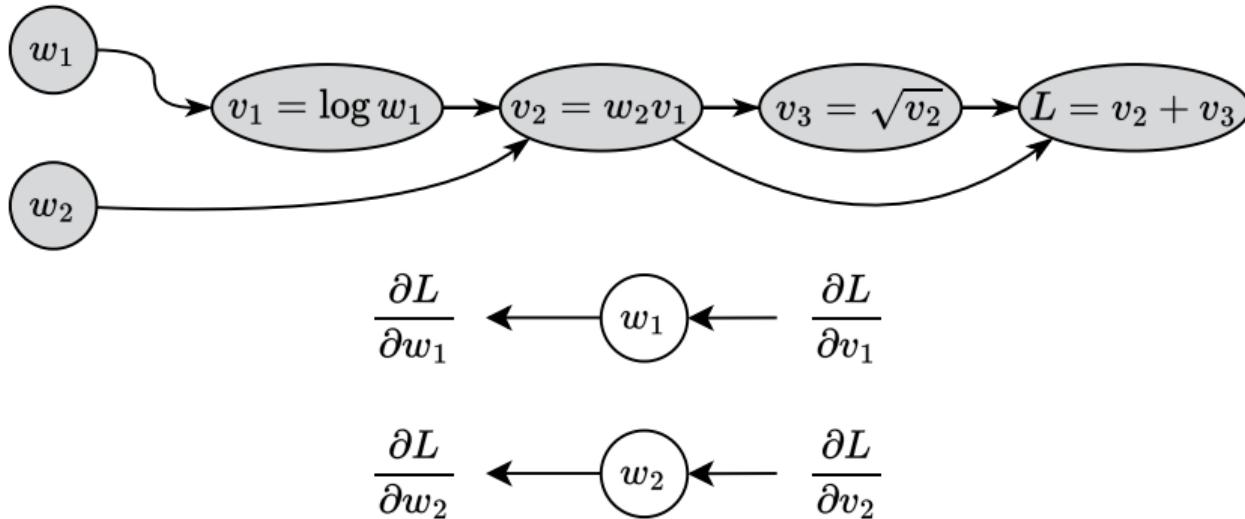


Figure 23: Illustration of backward mode automatic differentiation

## Derivatives

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial v_1} \frac{\partial v_1}{\partial w_1} = \frac{\partial L}{\partial v_1} \frac{1}{w_1} \quad \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial v_2} \frac{\partial v_2}{\partial w_2} = \frac{\partial L}{\partial v_1} v_1$$

## Backward (reverse) mode automatic differentiation

### Question

Note, that for the same price of computations as it was in the forward mode we have the full vector of gradient  $\nabla_w L$ . Is it a free lunch? What is the cost of acceleration?

## Backward (reverse) mode automatic differentiation

### Question

Note, that for the same price of computations as it was in the forward mode we have the full vector of gradient  $\nabla_w L$ . Is it a free lunch? What is the cost of acceleration?

**Answer** Note, that for using the reverse mode AD you need to store all intermediate computations from the forward pass. This problem could be somehow mitigated with the gradient checkpointing approach, which involves necessary recomputations of some intermediate values. This could significantly reduce the memory footprint of the large machine-learning model.

## Reverse mode automatic differentiation algorithm

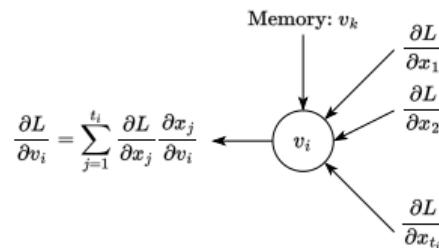
Suppose, we have a computational graph  $v_i, i \in [1; N]$ .

Our goal is to calculate the derivative of the output of this graph with respect to all inputs variable  $w$ ,

i.e.  $\nabla_w v_N = \left( \frac{\partial v_N}{\partial w_1}, \dots, \frac{\partial v_N}{\partial w_d} \right)^T$ . This idea implies

propagation of the gradient of the function with respect to the intermediate variables from the end to the origin, that is why we can introduce the notation:

$$\overline{v_i} = \frac{\partial L}{\partial v_i} = \frac{\partial v_N}{\partial v_i}$$



- FORWARD PASS

For  $i = 1, \dots, N$ :

Figure 24: Illustration of reverse chain rule to calculate the derivative of the function  $L$  with respect to the node  $v_i$ .

## Reverse mode automatic differentiation algorithm

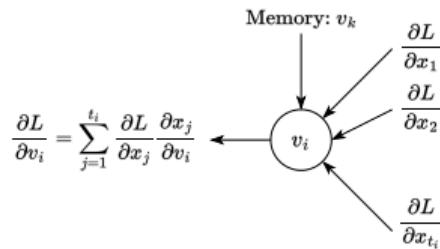
Suppose, we have a computational graph  $v_i, i \in [1; N]$ .

Our goal is to calculate the derivative of the output of this graph with respect to all inputs variable  $w$ ,

i.e.  $\nabla_w v_N = \left( \frac{\partial v_N}{\partial w_1}, \dots, \frac{\partial v_N}{\partial w_d} \right)^T$ . This idea implies

propagation of the gradient of the function with respect to the intermediate variables from the end to the origin, that is why we can introduce the notation:

$$\overline{v_i} = \frac{\partial L}{\partial v_i} = \frac{\partial v_N}{\partial v_i}$$



- **FORWARD PASS**

For  $i = 1, \dots, N$ :

- Compute and store the values of  $v_i$  as a function of its parents (inputs)

Figure 24: Illustration of reverse chain rule to calculate the derivative of the function  $L$  with respect to the node  $v_i$ .

## Reverse mode automatic differentiation algorithm

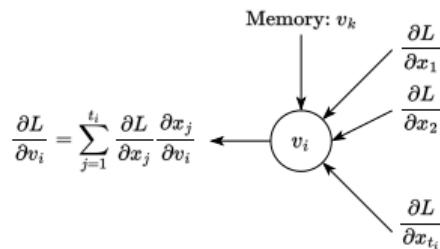
Suppose, we have a computational graph  $v_i, i \in [1; N]$ .

Our goal is to calculate the derivative of the output of this graph with respect to all inputs variable  $w$ ,

i.e.  $\nabla_w v_N = \left( \frac{\partial v_N}{\partial w_1}, \dots, \frac{\partial v_N}{\partial w_d} \right)^T$ . This idea implies

propagation of the gradient of the function with respect to the intermediate variables from the end to the origin, that is why we can introduce the notation:

$$\overline{v_i} = \frac{\partial L}{\partial v_i} = \frac{\partial v_N}{\partial v_i}$$



- **FORWARD PASS**

For  $i = 1, \dots, N$ :

- Compute and store the values of  $v_i$  as a function of its parents (inputs)

- **BACKWARD PASS**

For  $i = N, \dots, 1$ :

Figure 24: Illustration of reverse chain rule to calculate the derivative of the function  $L$  with respect to the node  $v_i$ .

## Reverse mode automatic differentiation algorithm

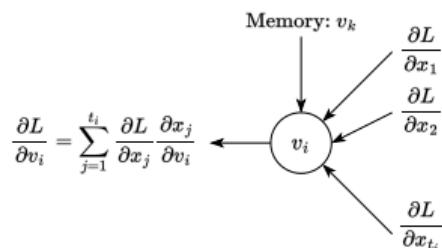
Suppose, we have a computational graph  $v_i, i \in [1; N]$ .

Our goal is to calculate the derivative of the output of this graph with respect to all inputs variable  $w$ ,

i.e.  $\nabla_w v_N = \left( \frac{\partial v_N}{\partial w_1}, \dots, \frac{\partial v_N}{\partial w_d} \right)^T$ . This idea implies

propagation of the gradient of the function with respect to the intermediate variables from the end to the origin, that is why we can introduce the notation:

$$\overline{v}_i = \frac{\partial L}{\partial v_i} = \frac{\partial v_N}{\partial v_i}$$



- **FORWARD PASS**

For  $i = 1, \dots, N$ :

- Compute and store the values of  $v_i$  as a function of its parents (inputs)

- **BACKWARD PASS**

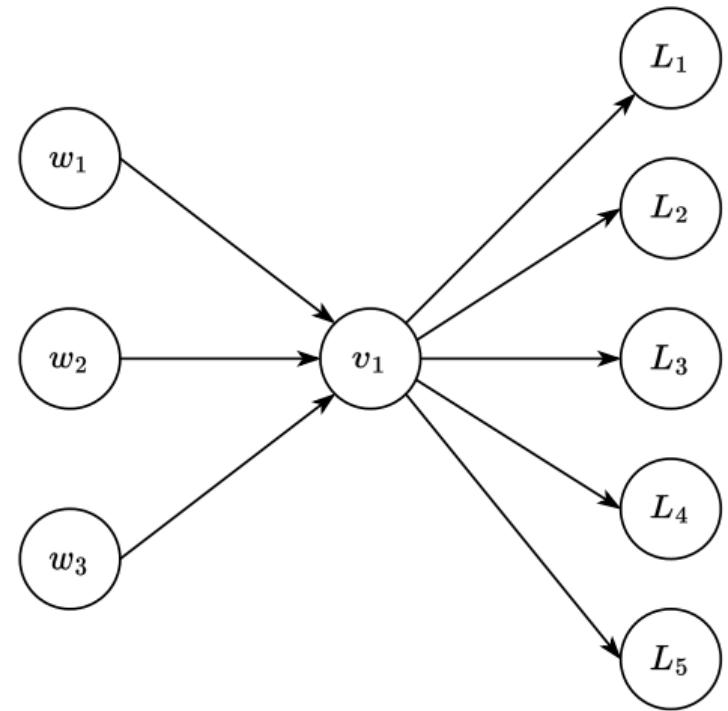
For  $i = N, \dots, 1$ :

- Compute the derivative  $\overline{v}_i$  using the backward chain rule and information from all of its children (outputs)  $(x_1, \dots, x_{t_i})$ :

$$\overline{v}_i = \frac{\partial L}{\partial v_i} = \sum_{j=1}^{t_i} \frac{\partial L}{\partial x_j} \frac{\partial x_j}{\partial v_i}$$

Figure 24: Illustration of reverse chain rule to calculate the derivative of the function  $L$  with respect to the node  $v_i$ .

## Choose your fighter



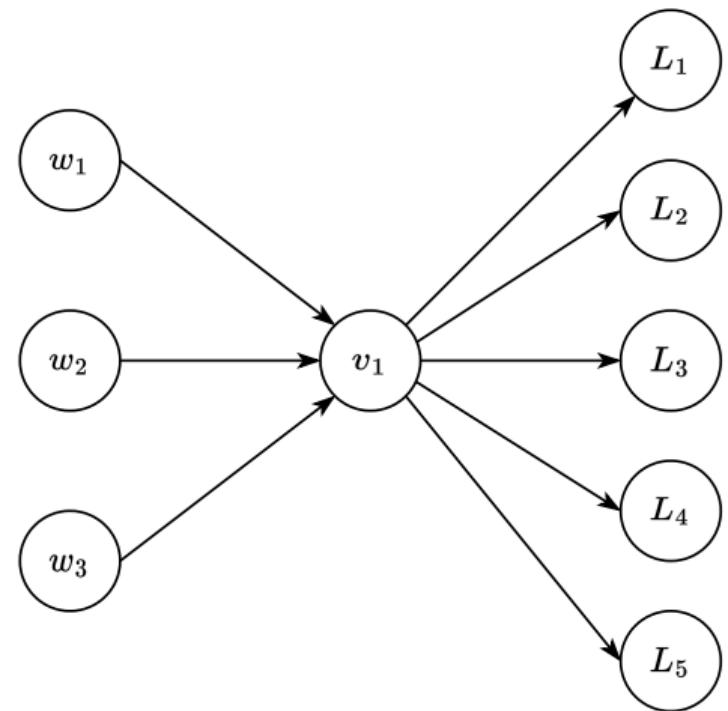
### Question

Which of the AD modes would you choose (forward/ reverse) for the following computational graph of primitive arithmetic operations? Suppose, you are needed to compute the jacobian

$$J = \left\{ \frac{\partial L_i}{\partial w_j} \right\}_{i,j}$$

Figure 25: Which mode would you choose for calculating gradients there?

## Choose your fighter



### Question

Which of the AD modes would you choose (forward/ reverse) for the following computational graph of primitive arithmetic operations? Suppose, you are needed to compute the jacobian

$$J = \left\{ \frac{\partial L_i}{\partial w_j} \right\}_{i,j}$$

**Answer** Note, that the reverse mode computational time is proportional to the number of outputs here, while the forward mode works proportionally to the number of inputs there. This is why it would be a good idea to consider the forward mode AD.

Figure 25: Which mode would you choose for calculating gradients there?

# Choose your fighter

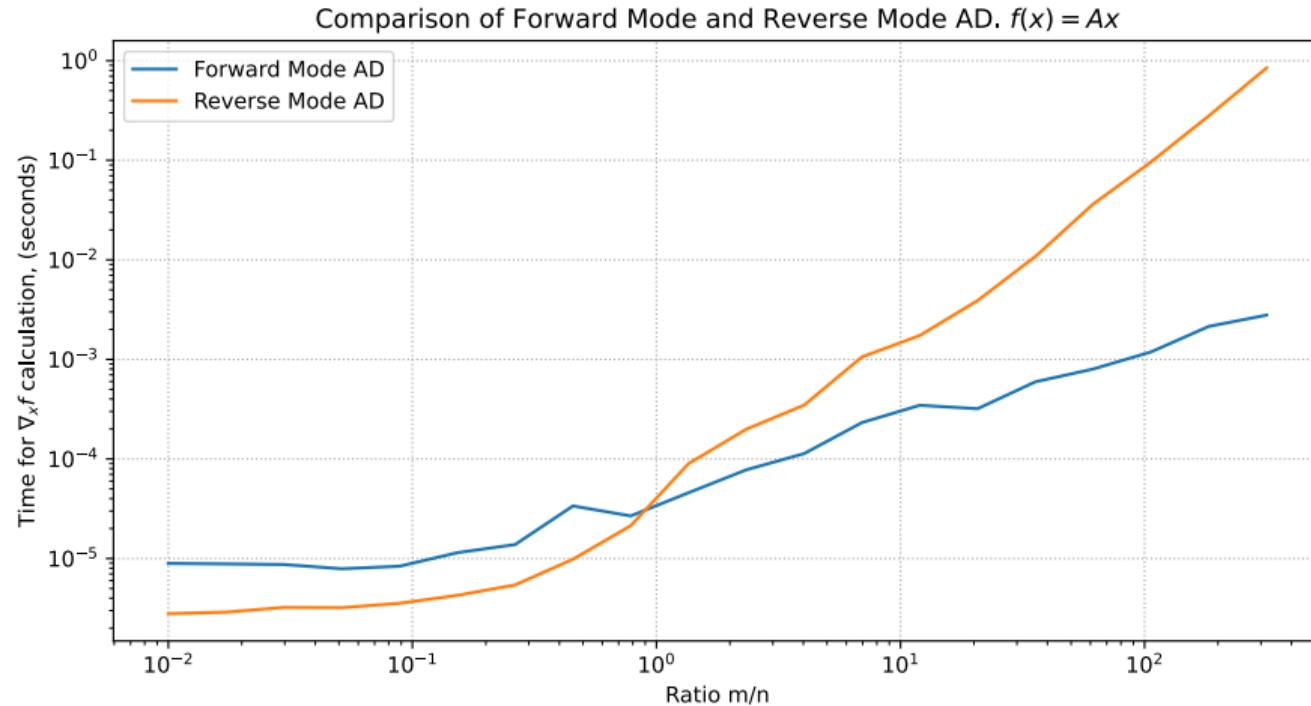
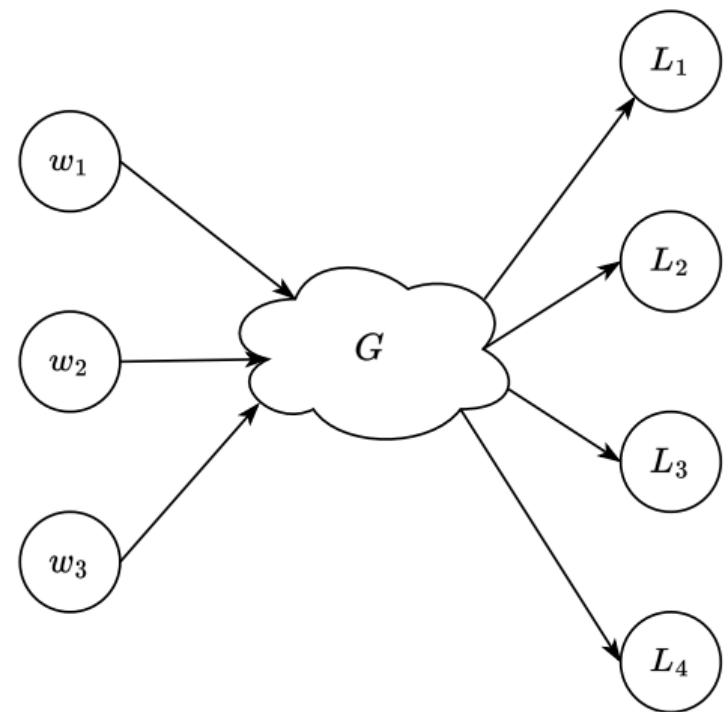


Figure 26: ♣ This graph nicely illustrates the idea of choice between the modes. The  $n = 100$  dimension is fixed and the graph presents the time needed for Jacobian calculation w.r.t.  $x$  for  $f(x) = Ax$

## Choose your fighter

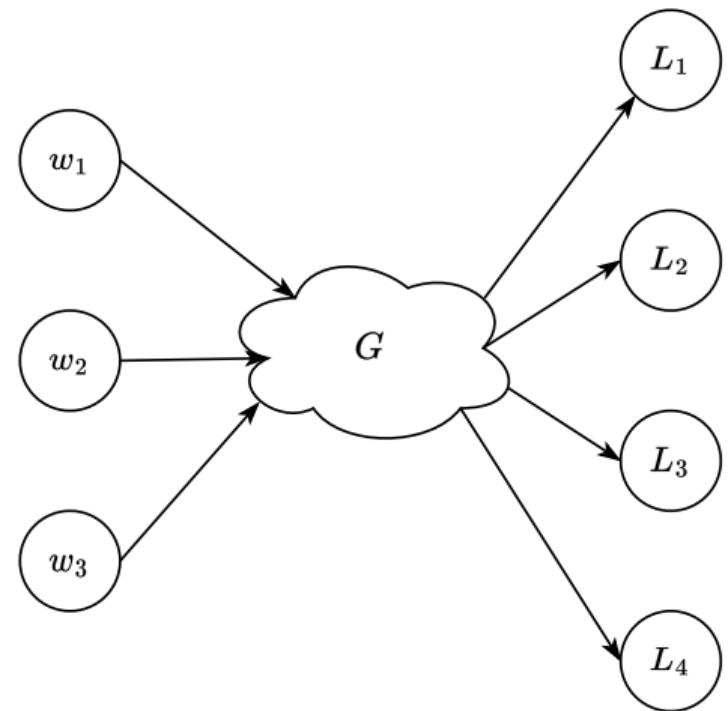


### Question

Which of the AD modes would you choose (forward/ reverse) for the following computational graph of primitive arithmetic operations? Suppose, you are needed to compute the jacobian  $J = \left\{ \frac{\partial L_i}{\partial w_j} \right\}_{i,j}$ . Note, that  $G$  is an arbitrary computational graph

Figure 27: Which mode would you choose for calculating gradients there?

## Choose your fighter



### Question

Which of the AD modes would you choose (forward/ reverse) for the following computational graph of primitive arithmetic operations? Suppose, you are needed to compute the jacobian  $J = \left\{ \frac{\partial L_i}{\partial w_j} \right\}_{i,j}$ . Note, that  $G$  is an arbitrary computational graph

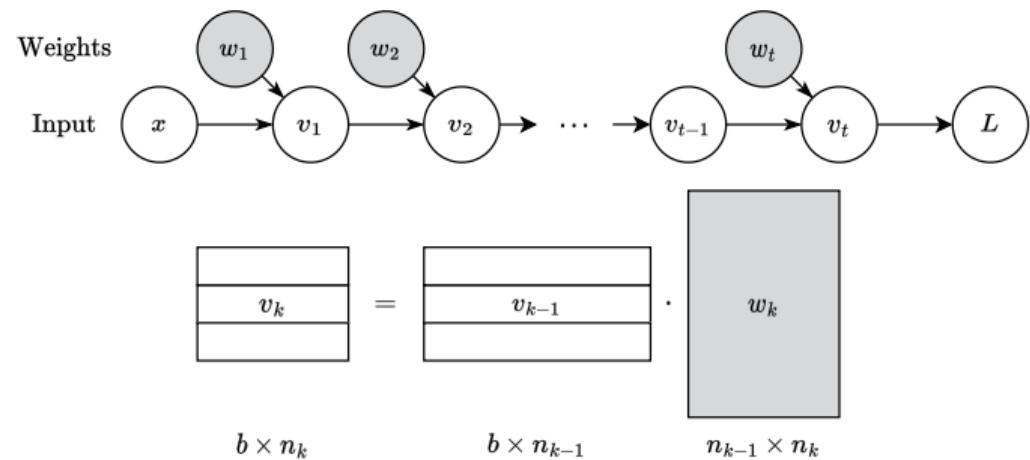
**Answer** It is generally impossible to say it without some knowledge about the specific structure of the graph  $G$ . Note, that there are also plenty of advanced approaches to mix forward and reverse mode AD, based on the specific  $G$  structure.

Figure 27: Which mode would you choose for calculating gradients there?

# Feedforward Architecture

## FORWARD

- $v_0 = x$  typically we have a batch of data  $x$  here as an input.



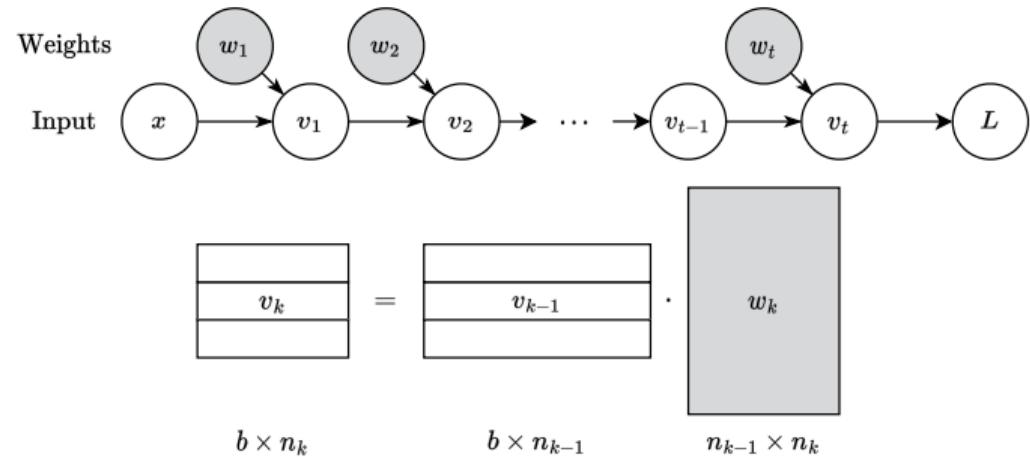
## BACKWARD

Figure 28: Feedforward neural network architecture

# Feedforward Architecture

## FORWARD

- $v_0 = x$  typically we have a batch of data  $x$  here as an input.
- For  $k = 1, \dots, t-1, t$ :



## BACKWARD

Figure 28: Feedforward neural network architecture

# Feedforward Architecture

## FORWARD

- $v_0 = x$  typically we have a batch of data  $x$  here as an input.
- For  $k = 1, \dots, t-1, t$ :
  - $v_k = \sigma(v_{k-1} w_k)$ . Note, that practically speaking the data has dimension  $x \in \mathbb{R}^{b \times d}$ , where  $b$  is the batch size (for the single data point  $b = 1$ ). While the weight matrix  $w_k$  of a  $k$  layer has a shape  $n_{k-1} \times n_k$ , where  $n_k$  is the dimension of an inner representation of the data.

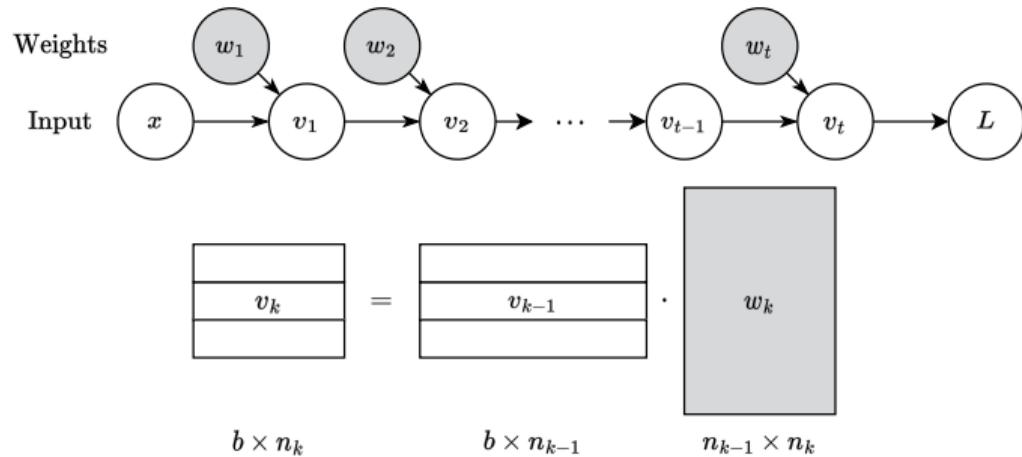


Figure 28: Feedforward neural network architecture

## BACKWARD

# Feedforward Architecture

## FORWARD

- $v_0 = x$  typically we have a batch of data  $x$  here as an input.
- For  $k = 1, \dots, t-1, t$ :
  - $v_k = \sigma(v_{k-1} w_k)$ . Note, that practically speaking the data has dimension  $x \in \mathbb{R}^{b \times d}$ , where  $b$  is the batch size (for the single data point  $b = 1$ ). While the weight matrix  $w_k$  of a  $k$  layer has a shape  $n_{k-1} \times n_k$ , where  $n_k$  is the dimension of an inner representation of the data.
  - $L = L(v_t)$  - calculate the loss function.

## BACKWARD

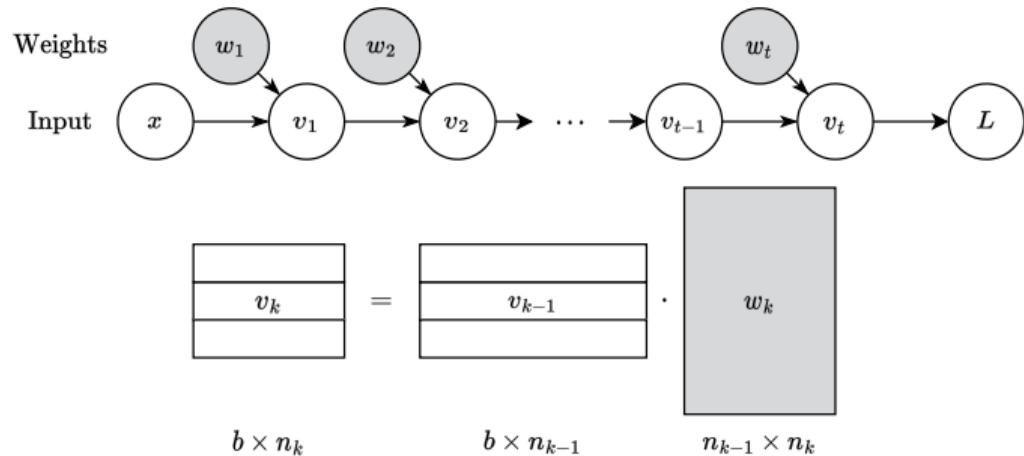


Figure 28: Feedforward neural network architecture

# Feedforward Architecture

## FORWARD

- $v_0 = x$  typically we have a batch of data  $x$  here as an input.
- For  $k = 1, \dots, t-1, t$ :
  - $v_k = \sigma(v_{k-1} w_k)$ . Note, that practically speaking the data has dimension  $x \in \mathbb{R}^{b \times d}$ , where  $b$  is the batch size (for the single data point  $b = 1$ ). While the weight matrix  $w_k$  of a  $k$  layer has a shape  $n_{k-1} \times n_k$ , where  $n_k$  is the dimension of an inner representation of the data.
  - $L = L(v_t)$  - calculate the loss function.

## BACKWARD

- $v_{t+1} = L, \frac{\partial L}{\partial L} = 1$

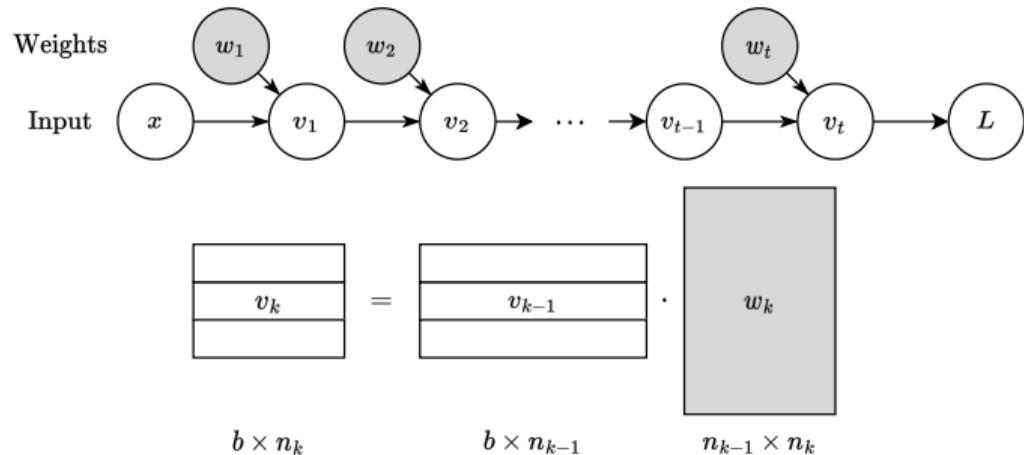


Figure 28: Feedforward neural network architecture

# Feedforward Architecture

## FORWARD

- $v_0 = x$  typically we have a batch of data  $x$  here as an input.
- For  $k = 1, \dots, t-1, t$ :
  - $v_k = \sigma(v_{k-1} w_k)$ . Note, that practically speaking the data has dimension  $x \in \mathbb{R}^{b \times d}$ , where  $b$  is the batch size (for the single data point  $b = 1$ ). While the weight matrix  $w_k$  of a  $k$  layer has a shape  $n_{k-1} \times n_k$ , where  $n_k$  is the dimension of an inner representation of the data.
- $L = L(v_t)$  - calculate the loss function.

## BACKWARD

- $v_{t+1} = L, \frac{\partial L}{\partial L} = 1$
- For  $k = t, t-1, \dots, 1$ :

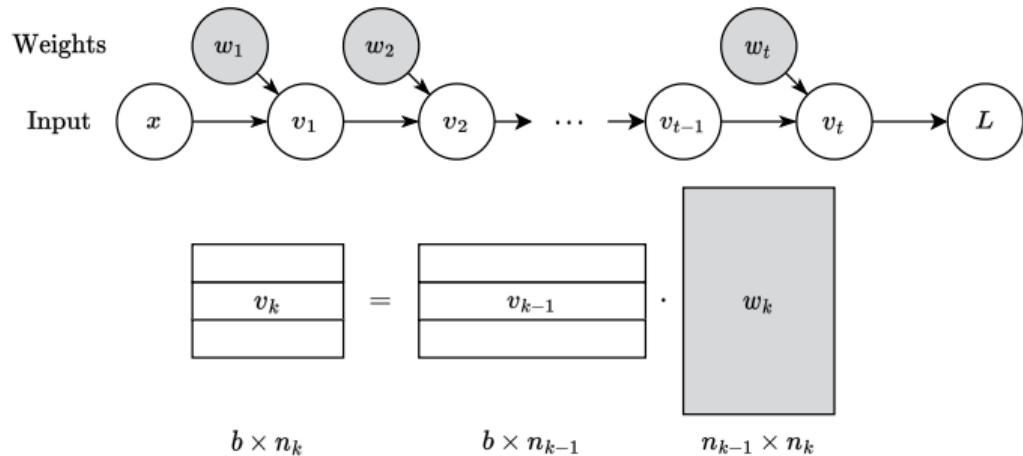


Figure 28: Feedforward neural network architecture

# Feedforward Architecture

## FORWARD

- $v_0 = x$  typically we have a batch of data  $x$  here as an input.
- For  $k = 1, \dots, t-1, t$ :
  - $v_k = \sigma(v_{k-1} w_k)$ . Note, that practically speaking the data has dimension  $x \in \mathbb{R}^{b \times d}$ , where  $b$  is the batch size (for the single data point  $b = 1$ ). While the weight matrix  $w_k$  of a  $k$  layer has a shape  $n_{k-1} \times n_k$ , where  $n_k$  is the dimension of an inner representation of the data.
- $L = L(v_t)$  - calculate the loss function.

## BACKWARD

- $v_{t+1} = L, \frac{\partial L}{\partial L} = 1$
- For  $k = t, t-1, \dots, 1$ :
  - $\frac{\partial L}{\partial v_k} = \frac{\partial L}{\partial v_{k+1}} \frac{\partial v_{k+1}}{\partial v_k}$

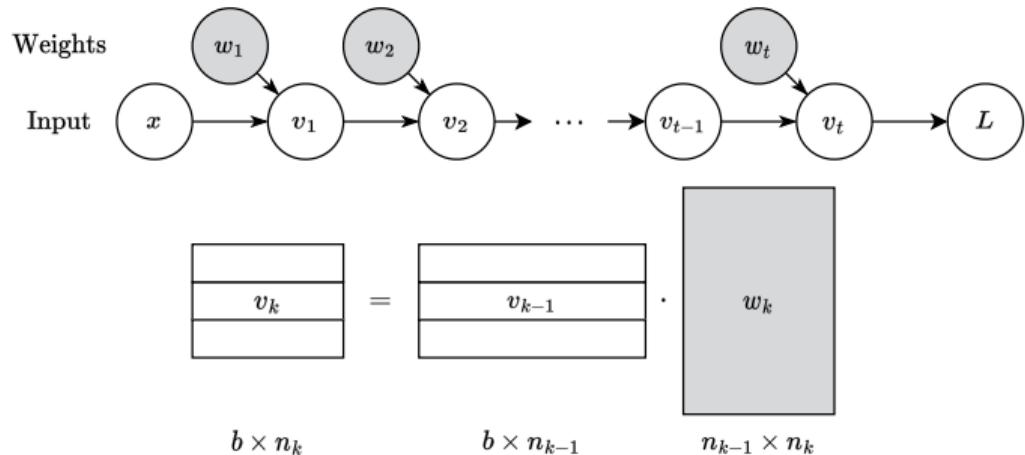


Figure 28: Feedforward neural network architecture

# Feedforward Architecture

## FORWARD

- $v_0 = x$  typically we have a batch of data  $x$  here as an input.
- For  $k = 1, \dots, t-1, t$ :
  - $v_k = \sigma(v_{k-1} w_k)$ . Note, that practically speaking the data has dimension  $x \in \mathbb{R}^{b \times d}$ , where  $b$  is the batch size (for the single data point  $b = 1$ ). While the weight matrix  $w_k$  of a  $k$  layer has a shape  $n_{k-1} \times n_k$ , where  $n_k$  is the dimension of an inner representation of the data.
- $L = L(v_t)$  - calculate the loss function.

## BACKWARD

- $v_{t+1} = L, \frac{\partial L}{\partial L} = 1$
- For  $k = t, t-1, \dots, 1$ :
  - $\frac{\partial L}{\partial v_k} = \frac{\partial L}{\partial v_{k+1}} \frac{\partial v_{k+1}}{\partial v_k}$   
 $b \times n_k \quad b \times n_{k+1} \quad n_{k+1} \times n_k$
  - $\frac{\partial L}{\partial w_k} = \frac{\partial L}{\partial v_{k+1}} \cdot \frac{\partial v_{k+1}}{\partial w_k}$   
 $b \times n_{k-1} \cdot n_k \quad b \times n_{k+1} \quad n_{k+1} \times n_{k-1} \cdot n_k$

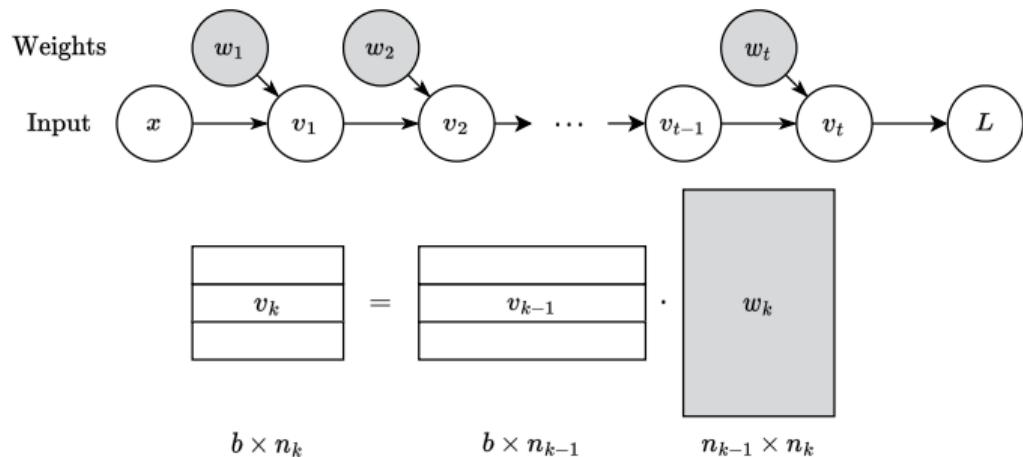
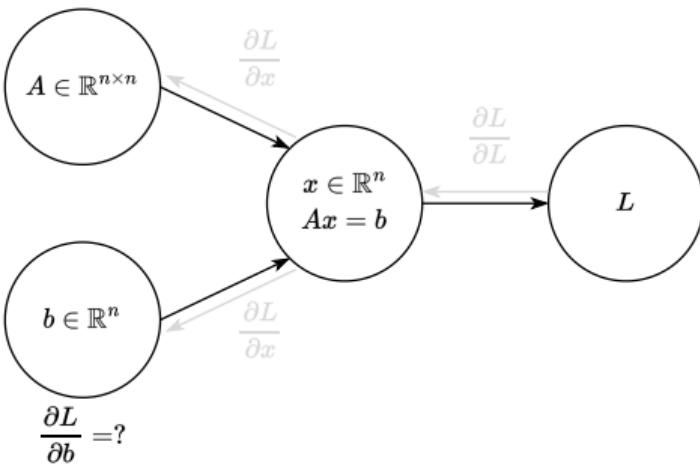


Figure 28: Feedforward neural network architecture

## Gradient propagation through the linear least squares

$$\frac{\partial L}{\partial A} = ?$$

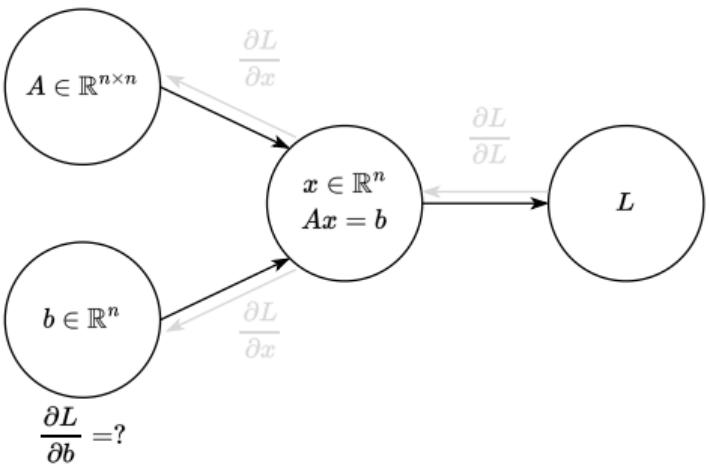


Suppose, we have an invertible matrix  $A$  and a vector  $b$ , the vector  $x$  is the solution of the linear system  $Ax = b$ , namely one can write down an analytical solution  $x = A^{-1}b$ , in this example we will show, that computing all derivatives  $\frac{\partial L}{\partial A}, \frac{\partial L}{\partial b}, \frac{\partial L}{\partial x}$ , i.e. the backward pass, costs approximately the same as the forward pass.

Figure 29:  $x$  could be found as a solution of linear system

## Gradient propagation through the linear least squares

$$\frac{\partial L}{\partial A} = ?$$



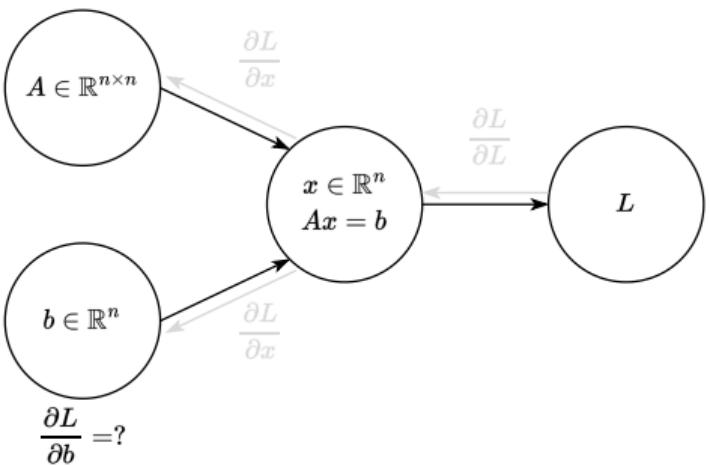
Suppose, we have an invertible matrix  $A$  and a vector  $b$ , the vector  $x$  is the solution of the linear system  $Ax = b$ , namely one can write down an analytical solution  $x = A^{-1}b$ , in this example we will show, that computing all derivatives  $\frac{\partial L}{\partial A}, \frac{\partial L}{\partial b}, \frac{\partial L}{\partial x}$ , i.e. the backward pass, costs approximately the same as the forward pass.  
It is known, that the differential of the function does not depend on the parametrization:

$$dL = \left\langle \frac{\partial L}{\partial x}, dx \right\rangle = \left\langle \frac{\partial L}{\partial A}, dA \right\rangle + \left\langle \frac{\partial L}{\partial b}, db \right\rangle$$

Figure 29:  $x$  could be found as a solution of linear system

## Gradient propagation through the linear least squares

$$\frac{\partial L}{\partial A} = ?$$



Suppose, we have an invertible matrix  $A$  and a vector  $b$ , the vector  $x$  is the solution of the linear system  $Ax = b$ , namely one can write down an analytical solution  $x = A^{-1}b$ , in this example we will show, that computing all derivatives  $\frac{\partial L}{\partial A}, \frac{\partial L}{\partial b}, \frac{\partial L}{\partial x}$ , i.e. the backward pass, costs approximately the same as the forward pass.  
It is known, that the differential of the function does not depend on the parametrization:

$$dL = \left\langle \frac{\partial L}{\partial x}, dx \right\rangle = \left\langle \frac{\partial L}{\partial A}, dA \right\rangle + \left\langle \frac{\partial L}{\partial b}, db \right\rangle$$

Given the linear system, we have:

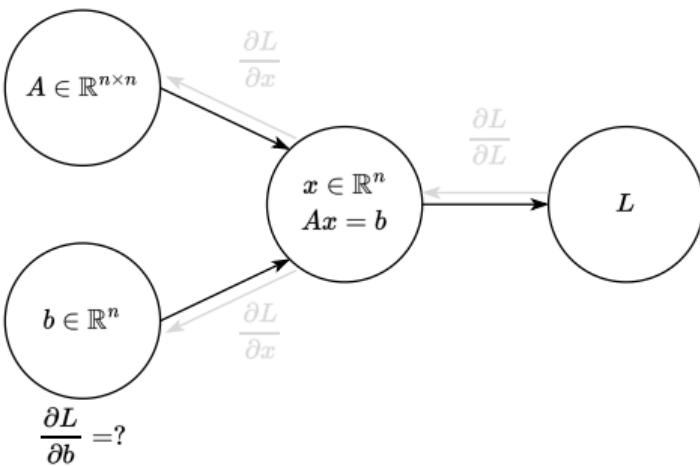
$$Ax = b$$

$$dAx + Adx = db \rightarrow dx = A^{-1}(db - dAx)$$

Figure 29:  $x$  could be found as a solution of linear system

## Gradient propagation through the linear least squares

$$\frac{\partial L}{\partial A} = ?$$



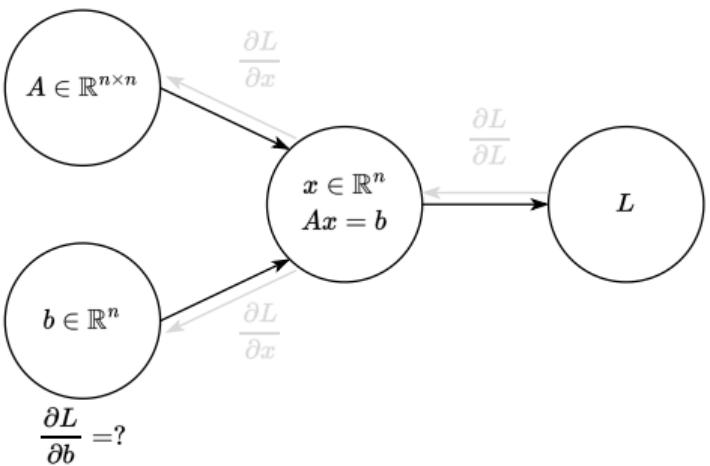
The straightforward substitution gives us:

$$\left\langle \frac{\partial L}{\partial x}, A^{-1}(db - dAx) \right\rangle = \left\langle \frac{\partial L}{\partial A}, dA \right\rangle + \left\langle \frac{\partial L}{\partial b}, db \right\rangle$$

Figure 30:  $x$  could be found as a solution of linear system

## Gradient propagation through the linear least squares

$$\frac{\partial L}{\partial A} = ?$$



The straightforward substitution gives us:

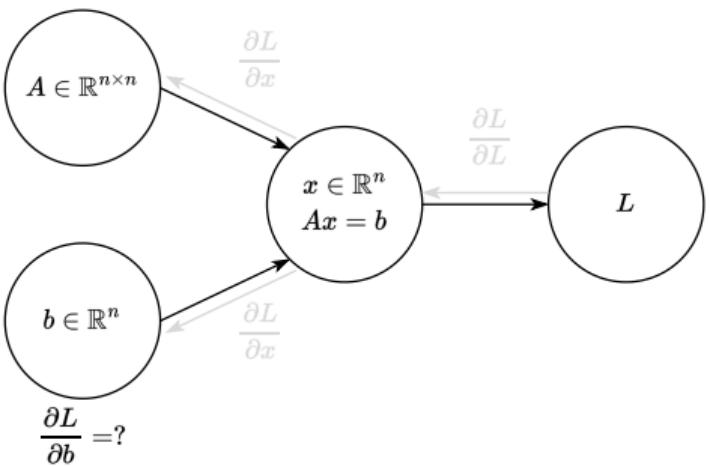
$$\left\langle \frac{\partial L}{\partial x}, A^{-1}(db - dAx) \right\rangle = \left\langle \frac{\partial L}{\partial A}, dA \right\rangle + \left\langle \frac{\partial L}{\partial b}, db \right\rangle$$

$$\left\langle -A^{-T} \frac{\partial L}{\partial x} x^T, dA \right\rangle + \left\langle A^{-T} \frac{\partial L}{\partial x}, db \right\rangle = \left\langle \frac{\partial L}{\partial A}, dA \right\rangle + \left\langle \frac{\partial L}{\partial b}, db \right\rangle$$

Figure 30:  $x$  could be found as a solution of linear system

## Gradient propagation through the linear least squares

$$\frac{\partial L}{\partial A} = ?$$



The straightforward substitution gives us:

$$\left\langle \frac{\partial L}{\partial x}, A^{-1}(db - dAx) \right\rangle = \left\langle \frac{\partial L}{\partial A}, dA \right\rangle + \left\langle \frac{\partial L}{\partial b}, db \right\rangle$$

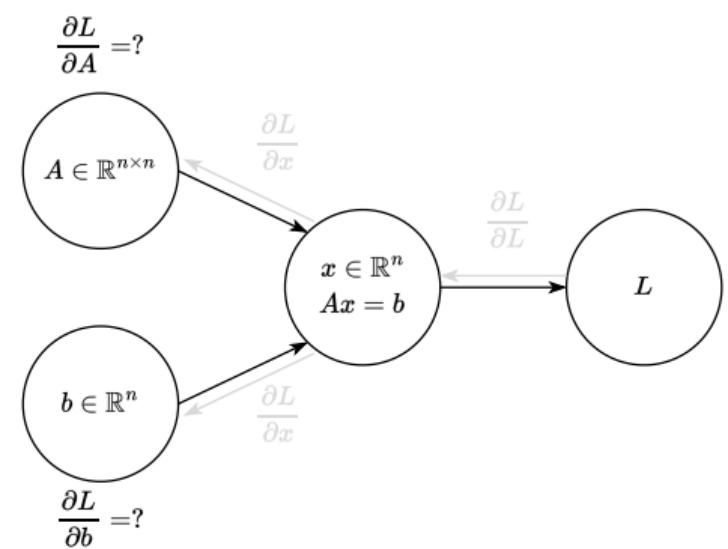
$$\left\langle -A^{-T} \frac{\partial L}{\partial x} x^T, dA \right\rangle + \left\langle A^{-T} \frac{\partial L}{\partial x}, db \right\rangle = \left\langle \frac{\partial L}{\partial A}, dA \right\rangle + \left\langle \frac{\partial L}{\partial b}, db \right\rangle$$

Therefore:

$$\frac{\partial L}{\partial A} = -A^{-T} \frac{\partial L}{\partial x} x^T \quad \frac{\partial L}{\partial b} = A^{-T} \frac{\partial L}{\partial x}$$

Figure 30:  $x$  could be found as a solution of linear system

## Gradient propagation through the linear least squares



The straightforward substitution gives us:

$$\left\langle \frac{\partial L}{\partial x}, A^{-1}(db - dAx) \right\rangle = \left\langle \frac{\partial L}{\partial A}, dA \right\rangle + \left\langle \frac{\partial L}{\partial b}, db \right\rangle$$

$$\left\langle -A^{-T} \frac{\partial L}{\partial x} x^T, dA \right\rangle + \left\langle A^{-T} \frac{\partial L}{\partial x}, db \right\rangle = \left\langle \frac{\partial L}{\partial A}, dA \right\rangle + \left\langle \frac{\partial L}{\partial b}, db \right\rangle$$

Therefore:

$$\frac{\partial L}{\partial A} = -A^{-T} \frac{\partial L}{\partial x} x^T \quad \frac{\partial L}{\partial b} = A^{-T} \frac{\partial L}{\partial x}$$

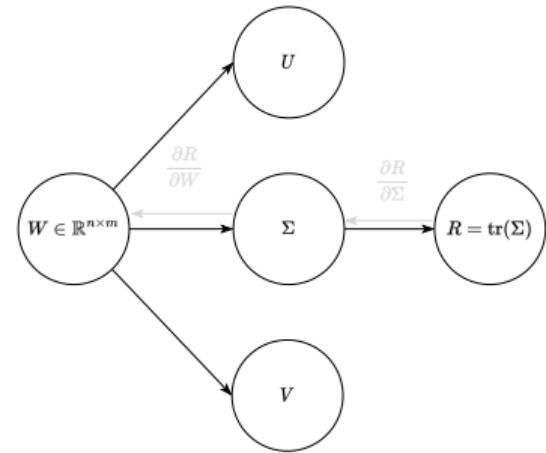
It is interesting, that the most computationally intensive part here is the matrix inverse, which is the same as for the forward pass.

Sometimes it is even possible to store the result itself, which makes the backward pass even cheaper.

Figure 30:  $x$  could be found as a solution of linear system

## Gradient propagation through the SVD

Suppose, we have the rectangular matrix  $W \in \mathbb{R}^{m \times n}$ , which has a singular value decomposition:



$$W = U\Sigma V^T, \quad U^T U = I, \quad V^T V = I, \quad \Sigma = \text{diag}(\sigma_1, \dots, \sigma_{\min(m,n)})$$

1. Similarly to the previous example:

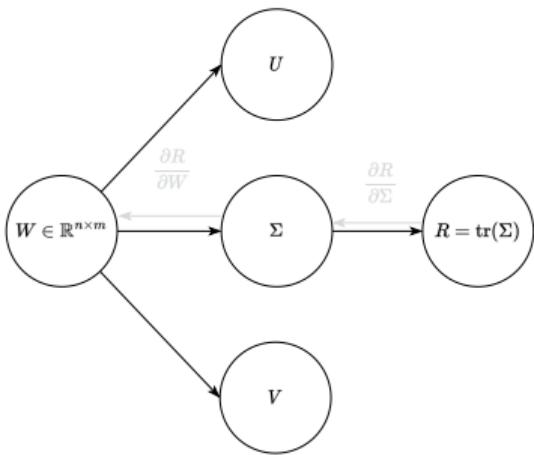
$$W = U\Sigma V^T$$

$$dW = dU\Sigma V^T + U d\Sigma V^T + U\Sigma dV^T$$

$$U^T dWV = U^T dU\Sigma V^T V + U^T U d\Sigma V^T V + U^T U\Sigma dV^T V$$

$$U^T dWV = U^T dU\Sigma + d\Sigma + \Sigma dV^T V$$

## Gradient propagation through the SVD



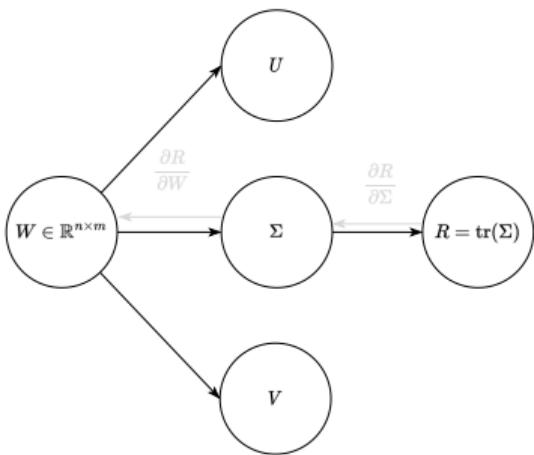
2. Note, that  $U^T U = I \rightarrow dU^T U + U^T dU = 0$ . But also  $dU^T U = (U^T dU)^T$ , which actually involves, that the matrix  $U^T dU$  is antisymmetric:

$$(U^T dU)^T + U^T dU = 0 \rightarrow \text{diag}(U^T dU) = (0, \dots, 0)$$

The same logic could be applied to the matrix  $V$  and

$$\text{diag}(dV^T V) = (0, \dots, 0)$$

# Gradient propagation through the SVD



2. Note, that  $U^T U = I \rightarrow dU^T U + U^T dU = 0$ . But also  $dU^T U = (U^T dU)^T$ , which actually involves, that the matrix  $U^T dU$  is antisymmetric:

$$(U^T dU)^T + U^T dU = 0 \rightarrow \text{diag}(U^T dU) = (0, \dots, 0)$$

The same logic could be applied to the matrix  $V$  and

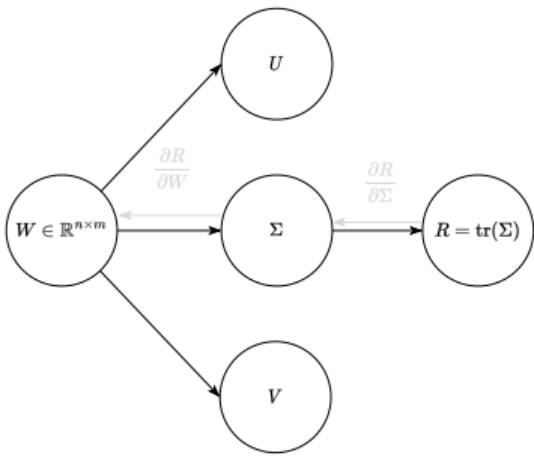
$$\text{diag}(dV^T V) = (0, \dots, 0)$$

3. At the same time, the matrix  $d\Sigma$  is diagonal, which means (look at the 1.) that

$$\text{diag}(U^T dWV) = d\Sigma$$

Here on both sides, we have diagonal matrices.

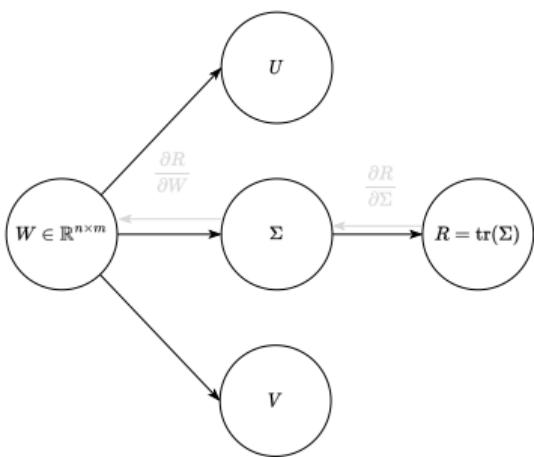
## Gradient propagation through the SVD



- Now, we can decompose the differential of the loss function as a function of  $\Sigma$  - such problems arise in ML problems, where we need to restrict the matrix rank:

$$\begin{aligned} dL &= \left\langle \frac{\partial L}{\partial \Sigma}, d\Sigma \right\rangle \\ &= \left\langle \frac{\partial L}{\partial \Sigma}, \text{diag}(U^T dWV) \right\rangle \\ &= \text{tr} \left( \frac{\partial L}{\partial \Sigma}^T \text{diag}(U^T dWV) \right) \end{aligned}$$

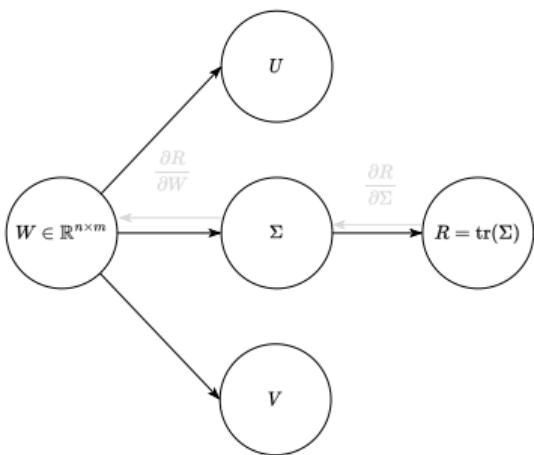
## Gradient propagation through the SVD



5. As soon as we have diagonal matrices inside the product, the trace of the diagonal part of the matrix will be equal to the trace of the whole matrix:

$$\begin{aligned} dL &= \text{tr} \left( \frac{\partial L^T}{\partial \Sigma} \text{diag}(U^T dWV) \right) \\ &= \text{tr} \left( \frac{\partial L^T}{\partial \Sigma} U^T dWV \right) \\ &= \left\langle \frac{\partial L}{\partial \Sigma}, U^T dWV \right\rangle \\ &= \left\langle U \frac{\partial L}{\partial \Sigma} V^T, dW \right\rangle \end{aligned}$$

# Gradient propagation through the SVD



6. Finally, using another parametrization of the differential

$$\left\langle U \frac{\partial L}{\partial \Sigma} V^T, dW \right\rangle = \left\langle \frac{\partial L}{\partial W}, dW \right\rangle$$

$$\frac{\partial L}{\partial W} = U \frac{\partial L}{\partial \Sigma} V^T,$$

This nice result allows us to connect the gradients  $\frac{\partial L}{\partial W}$  and  $\frac{\partial L}{\partial \Sigma}$ .

## Hessian vector product without the Hessian

When you need some information about the curvature of the function you usually need to work with the hessian. However, when the dimension of the problem is large it is challenging. For a scalar-valued function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , the Hessian at a point  $x \in \mathbb{R}^n$  is written as  $\nabla^2 f(x)$ . A Hessian-vector product function is then able to evaluate

$$v \mapsto \nabla^2 f(x) \cdot v$$

for any vector  $v \in \mathbb{R}^n$ . We have to use the identity

$$\nabla^2 f(x)v = \nabla[x \mapsto \nabla f(x) \cdot v] = \nabla g(x),$$

where  $g(x) = \nabla f(x)^T \cdot v$  is a new vector-valued function that dots the gradient of  $f$  at  $x$  with the vector  $v$ .

```
import jax.numpy as jnp

def hvp(f, x, v):
    return grad(lambda x: jnp.vdot(grad(f)(x), v))(x)
```

## Hutchinson Trace Estimation <sup>2</sup>

This example illustrates the estimation the Hessian trace of a neural network using Hutchinson's method, which is an algorithm to obtain such an estimate from matrix-vector products:

Let  $X \in \mathbb{R}^{d \times d}$  and  $v \in \mathbb{R}^d$  be a random vector such that  $\mathbb{E}[vv^T] = I$ . Then,

$$\text{Tr}(X) = \mathbb{E}[v^T X v] = \frac{1}{V} \sum_{i=1}^V v_i^T X v_i$$

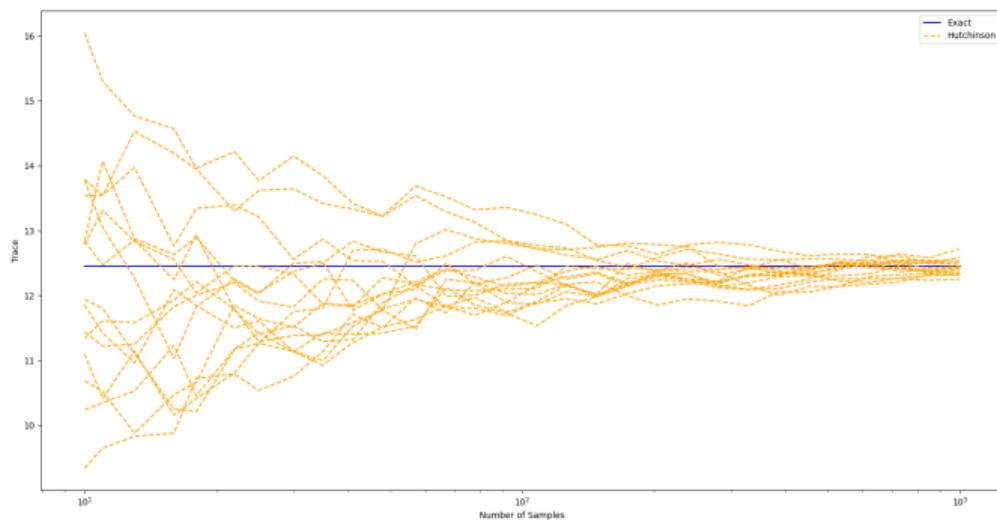


Figure 31: Source

<sup>2</sup>A stochastic estimator of the trace of the influence matrix for Laplacian smoothing splines - M.F. Hutchinson, 1990  
Automatic differentiation

## Activation checkpointing

The animated visualization of the above approaches 

An example of using a gradient checkpointing 

## What automatic differentiation (AD) is NOT:

- AD is not a finite differences

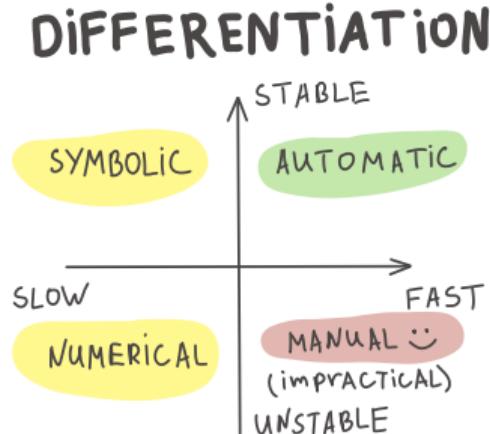


Figure 32: Different approaches for taking derivatives

## What automatic differentiation (AD) is NOT:

- AD is not a finite differences
- AD is not a symbolic derivative

### DIFFERENTIATION

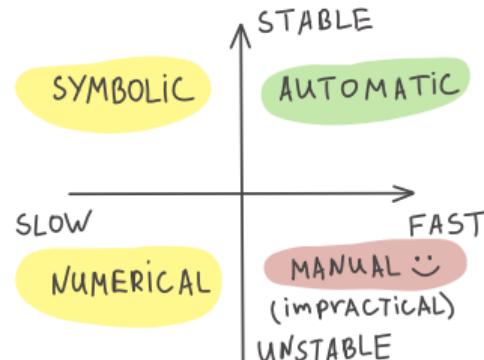


Figure 32: Different approaches for taking derivatives

## What automatic differentiation (AD) is NOT:

- AD is not a finite differences
- AD is not a symbolic derivative
- AD is not just the chain rule

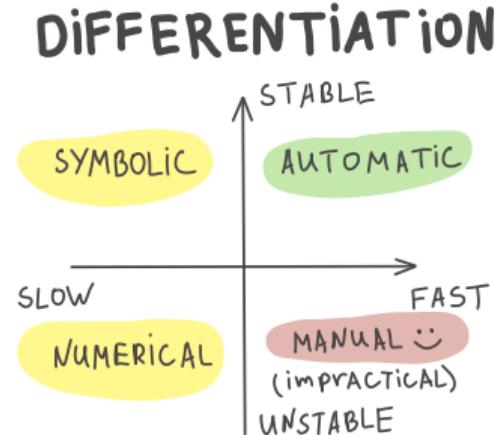


Figure 32: Different approaches for taking derivatives

## What automatic differentiation (AD) is NOT:

- AD is not a finite differences
- AD is not a symbolic derivative
- AD is not just the chain rule
- AD is not just backpropagation

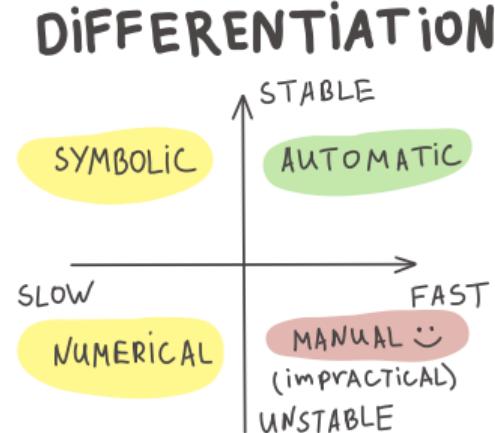


Figure 32: Different approaches for taking derivatives

## What automatic differentiation (AD) is NOT:

- AD is not a finite differences
- AD is not a symbolic derivative
- AD is not just the chain rule
- AD is not just backpropagation
- AD (reverse mode) is time-efficient and numerically stable

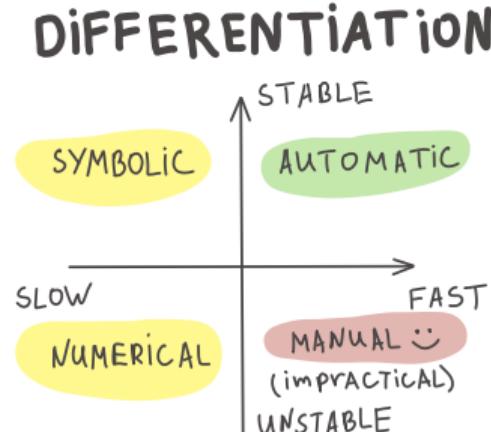


Figure 32: Different approaches for taking derivatives

## What automatic differentiation (AD) is NOT:

- AD is not a finite differences
- AD is not a symbolic derivative
- AD is not just the chain rule
- AD is not just backpropagation
- AD (reverse mode) is time-efficient and numerically stable
- AD (reverse mode) is memory inefficient (you need to store all intermediate computations from the forward pass).

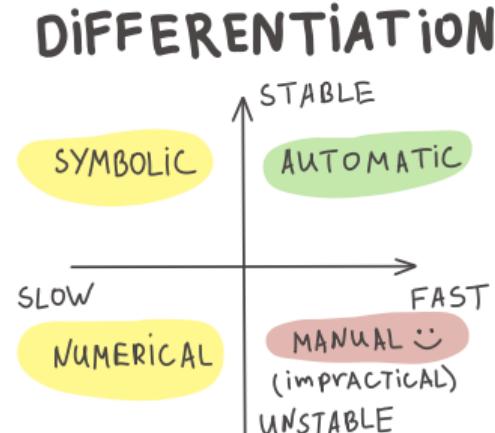


Figure 32: Different approaches for taking derivatives

# Code

Open In Colab 

## Gradient checkpointing

## Feedforward Architecture

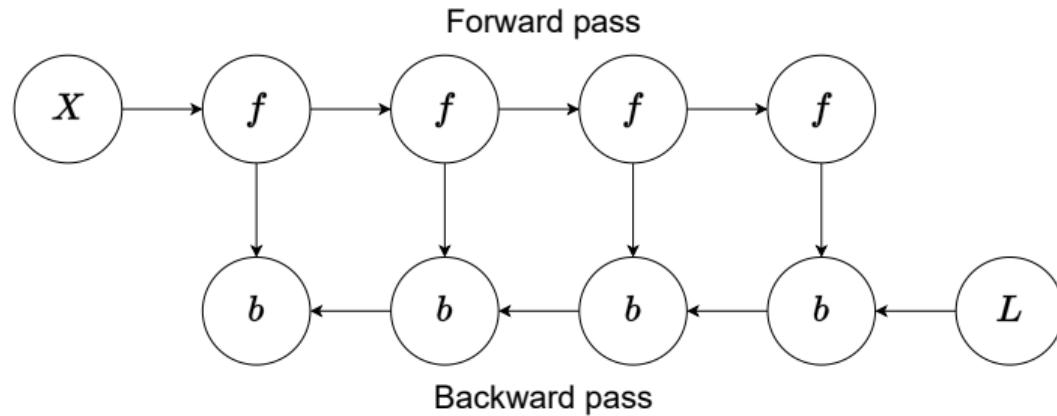


Figure 33: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The activations marked with an  $f$ . The gradient of the loss with respect to the activations and parameters marked with  $b$ .

# Feedforward Architecture

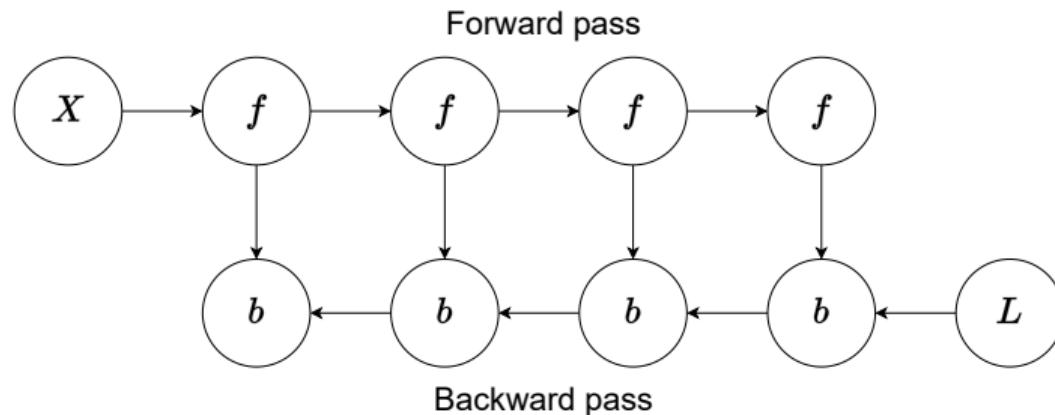


Figure 33: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The activations marked with an  $f$ . The gradient of the loss with respect to the activations and parameters marked with  $b$ .

## ! Important

The results obtained for the  $f$  nodes are needed to compute the  $b$  nodes.

## Vanilla backpropagation

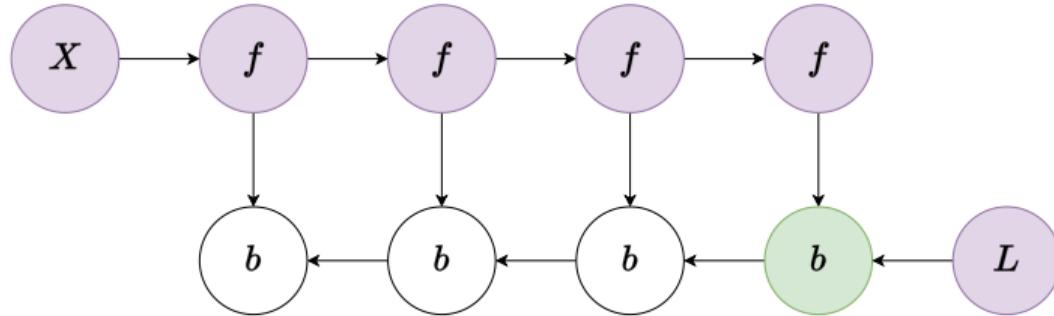


Figure 34: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

## Vanilla backpropagation

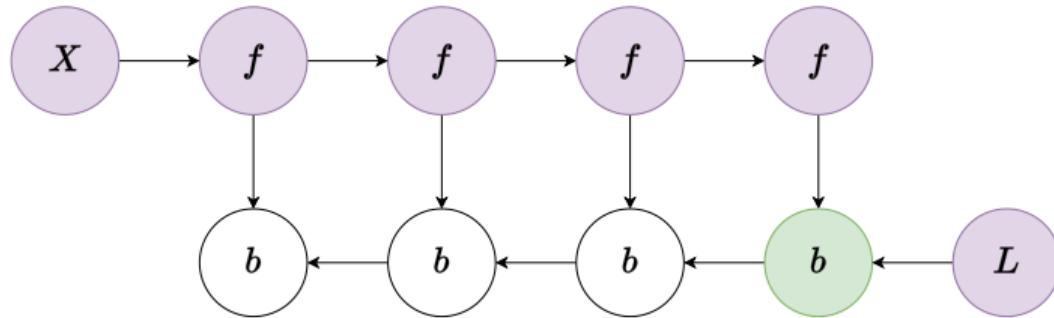


Figure 34: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- All activations  $f$  are kept in memory after the forward pass.

## Vanilla backpropagation

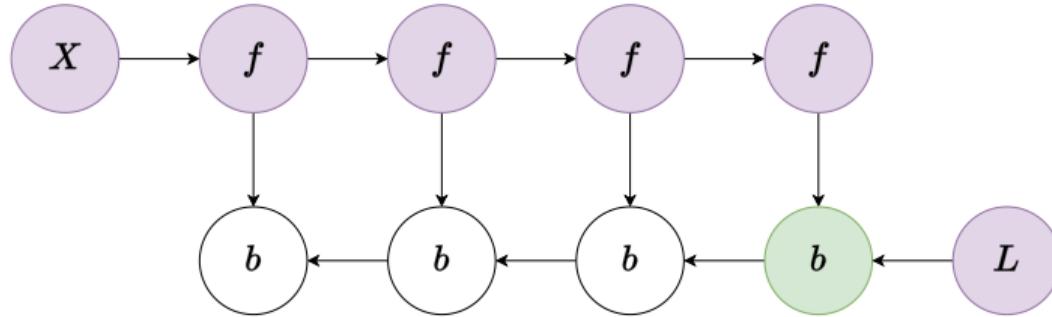


Figure 34: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- All activations  $f$  are kept in memory after the forward pass.

## Vanilla backpropagation

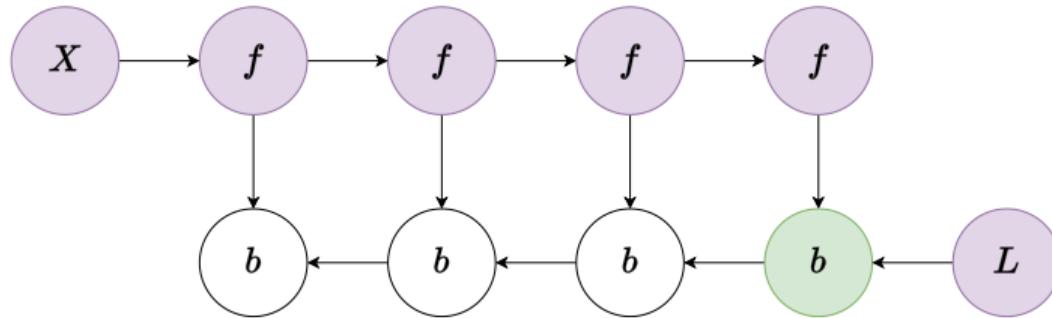


Figure 34: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- All activations  $f$  are kept in memory after the forward pass.
- Optimal in terms of computation: it only computes each node once.

## Vanilla backpropagation

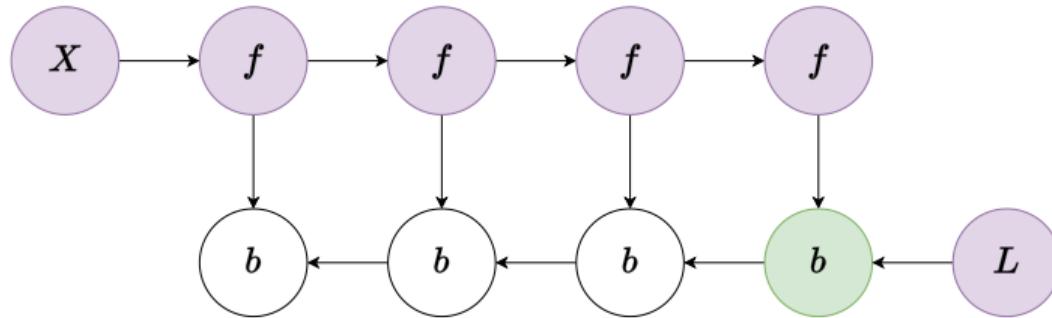


Figure 34: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- All activations  $f$  are kept in memory after the forward pass.
- Optimal in terms of computation: it only computes each node once.

## Vanilla backpropagation

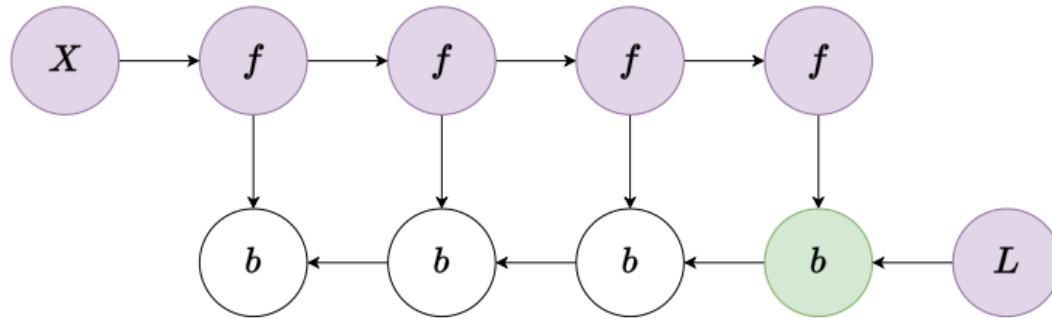


Figure 34: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- All activations  $f$  are kept in memory after the forward pass.
  - Optimal in terms of computation: it only computes each node once.
  - High memory usage. The memory usage grows linearly with the number of layers in the neural network.

## Memory poor backpropagation

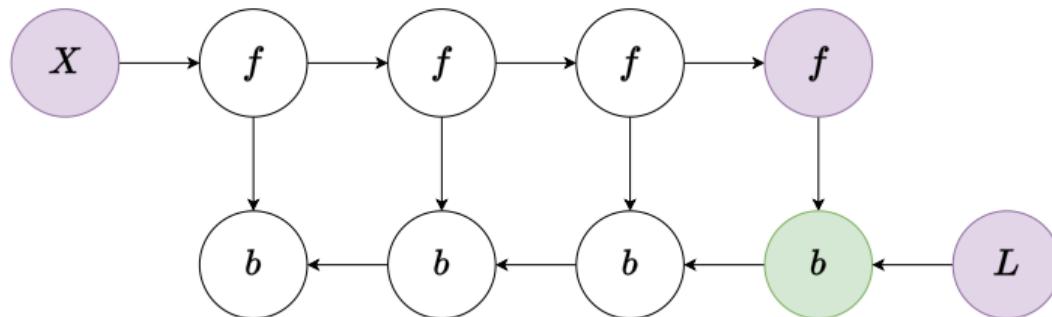


Figure 35: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

## Memory poor backpropagation

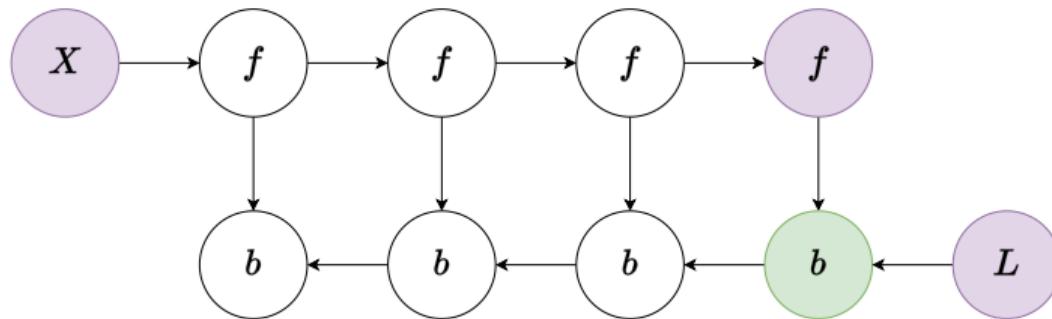


Figure 35: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- Each activation  $f$  is recalculated as needed.

## Memory poor backpropagation

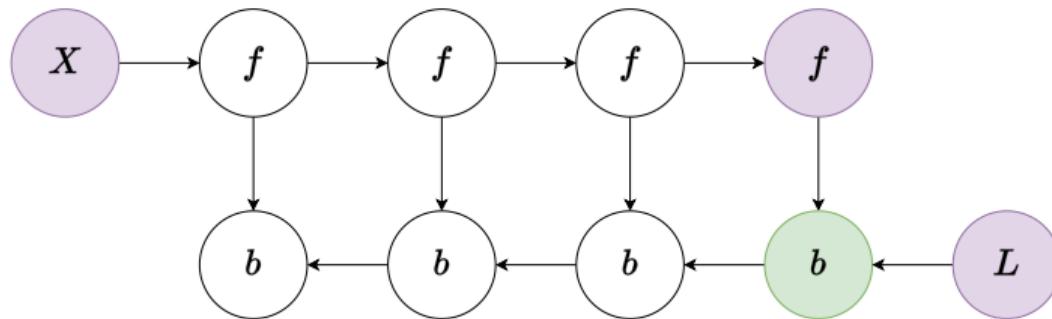


Figure 35: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- Each activation  $f$  is recalculated as needed.

## Memory poor backpropagation

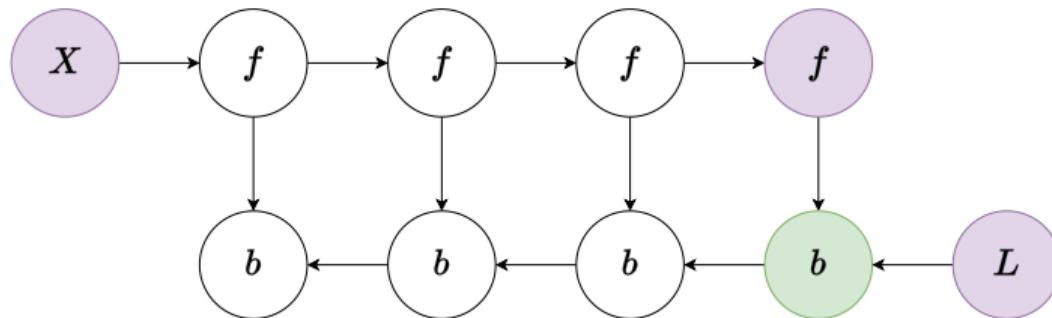


Figure 35: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- Each activation  $f$  is recalculated as needed.
  - Optimal in terms of memory: there is no need to store all activations in memory.

## Memory poor backpropagation

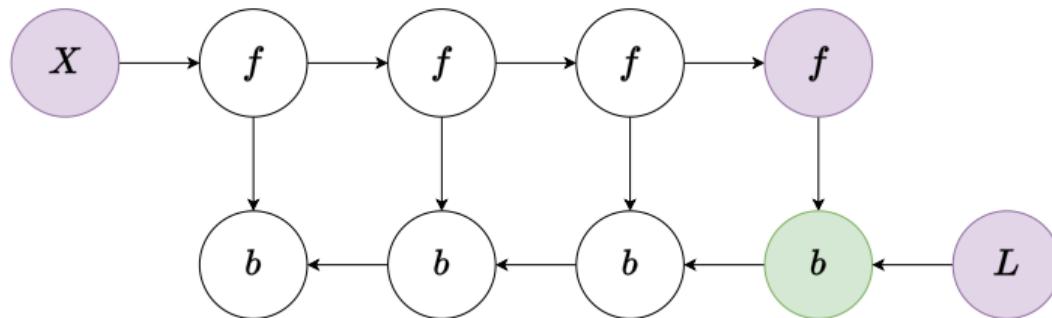


Figure 35: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- Each activation  $f$  is recalculated as needed.
  - Optimal in terms of memory: there is no need to store all activations in memory.

## Memory poor backpropagation

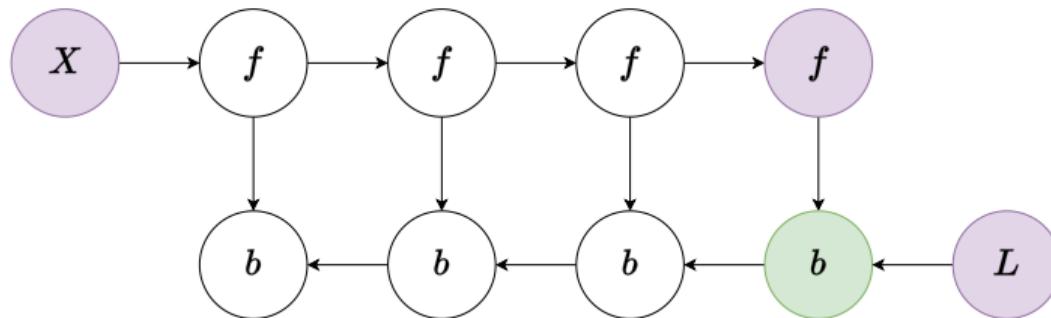


Figure 35: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- Each activation  $f$  is recalculated as needed.
  - Optimal in terms of memory: there is no need to store all activations in memory.
- Computationally inefficient. The number of node evaluations scales with  $n^2$ , whereas vanilla backprop scaled as  $n$ : each of the  $n$  nodes is recomputed on the order of  $n$  times.

## Checkpointed backpropagation

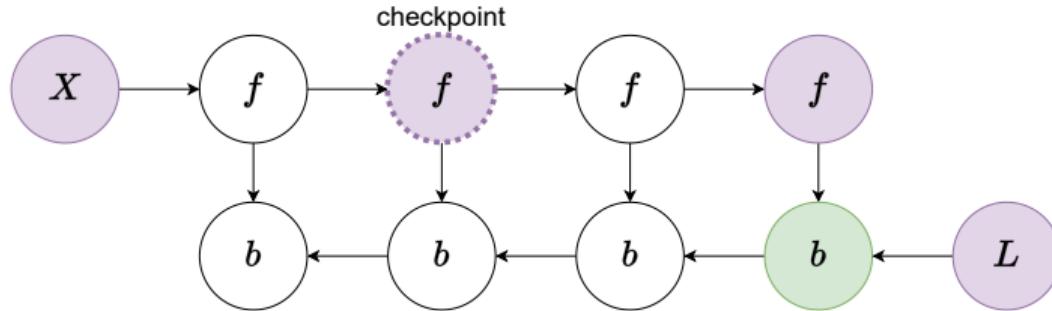


Figure 36: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

## Checkpointed backpropagation

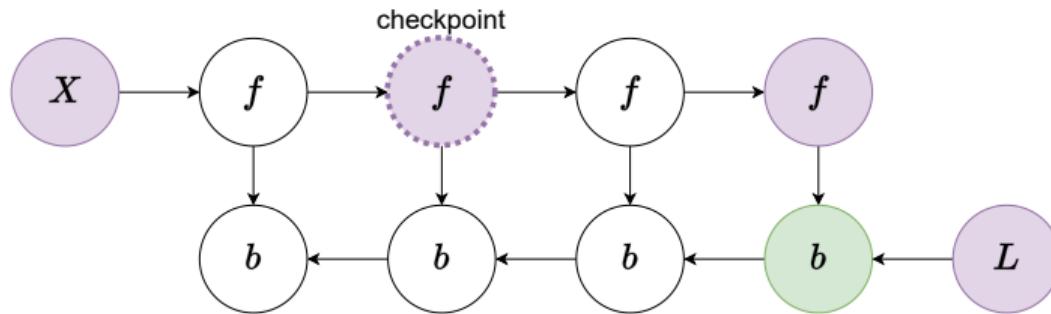


Figure 36: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- Trade-off between the **vanilla** and **memory poor** approaches. The strategy is to mark a subset of the neural net activations as checkpoint nodes, that will be stored in memory.

## Checkpointed backpropagation

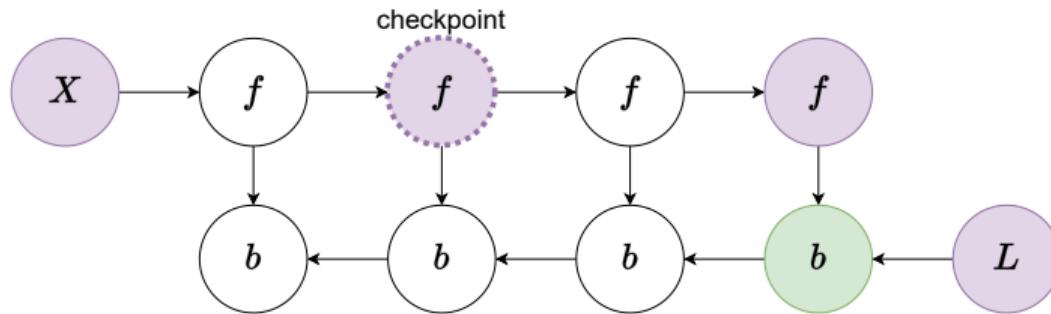


Figure 36: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- Trade-off between the **vanilla** and **memory poor** approaches. The strategy is to mark a subset of the neural net activations as checkpoint nodes, that will be stored in memory.

## Checkpointed backpropagation

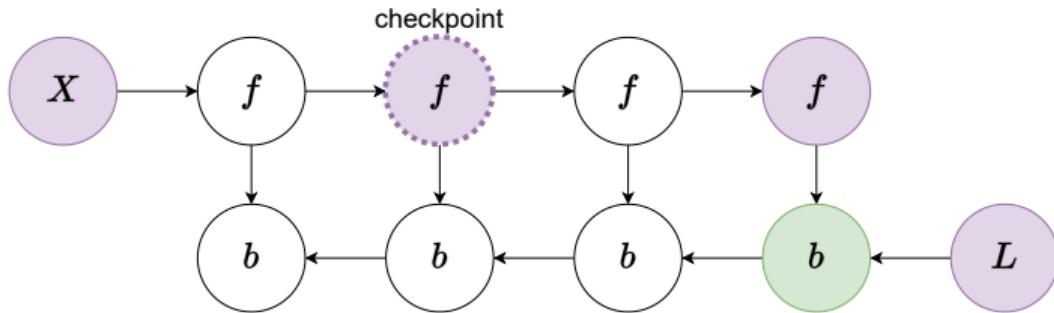


Figure 36: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- Trade-off between the **vanilla** and **memory poor** approaches. The strategy is to mark a subset of the neural net activations as checkpoint nodes, that will be stored in memory.
- Faster recalculation of activations  $f$ . We only need to recompute the nodes between a  $b$  node and the last checkpoint preceding it when computing that  $b$  node during backprop.

## Checkpointed backpropagation

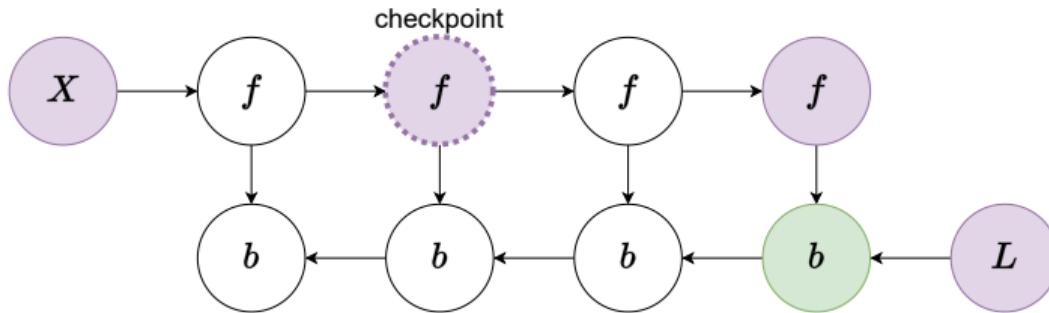


Figure 36: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- Trade-off between the **vanilla** and **memory poor** approaches. The strategy is to mark a subset of the neural net activations as checkpoint nodes, that will be stored in memory.
- Faster recalculation of activations  $f$ . We only need to recompute the nodes between a  $b$  node and the last checkpoint preceding it when computing that  $b$  node during backprop.

## Checkpointed backpropagation

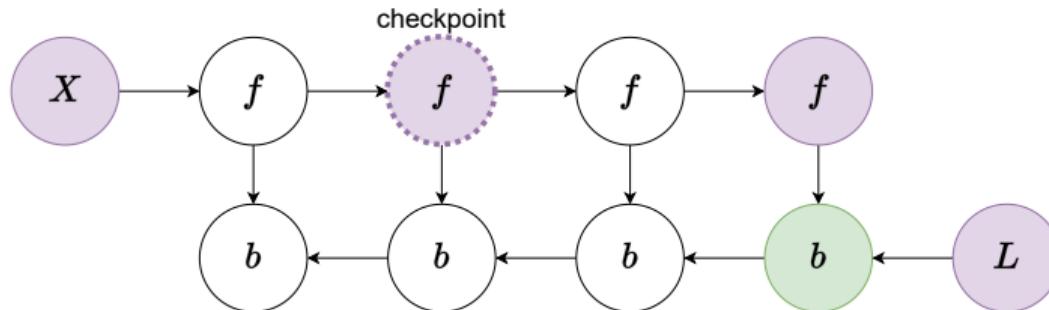


Figure 36: Computation graph for obtaining gradients for a simple feed-forward neural network with  $n$  layers. The purple color indicates nodes that are stored in memory.

- Trade-off between the **vanilla** and **memory poor** approaches. The strategy is to mark a subset of the neural net activations as checkpoint nodes, that will be stored in memory.
  - Faster recalculation of activations  $f$ . We only need to recompute the nodes between a  $b$  node and the last checkpoint preceding it when computing that  $b$  node during backprop.
  - Memory consumption depends on the number of checkpoints. More effective than **vanilla** approach.

## Gradient checkpointing visualization

The animated visualization of the above approaches 

An example of using a gradient checkpointing 