# Lecture 4:
# Neural Networks and Backpropagation

So far: backprop with scalars

What about vector-valued functions?

# Recap: Vector derivatives

Scalar to Scalar

$x \in \mathbb{R}, y \in \mathbb{R}$

Regular derivative:

$$\frac{\partial y}{\partial x} \in \mathbb{R}$$

If x changes by a small amount, how much will y change?

# Recap: Vector derivatives

## Scalar to Scalar

$x \in \mathbb{R}, y \in \mathbb{R}$

Regular derivative:

$$\frac{\partial y}{\partial x} \in \mathbb{R}$$

If x changes by a small amount, how much will y change?

## Vector to Scalar

$x \in \mathbb{R}^N, y \in \mathbb{R}$

Derivative is **Gradient**:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left(\frac{\partial y}{\partial x}\right)_n = \frac{\partial y}{\partial x_n}$$

For each element of x, if it changes by a small amount then how much will y change?

# Recap: Vector derivatives

| Scalar to Scalar | Vector to Scalar | Vector to Vector |
|---|---|---|

**Scalar to Scalar**

$x \in \mathbb{R}, y \in \mathbb{R}$

Regular derivative:

$$\frac{\partial y}{\partial x} \in \mathbb{R}$$

If x changes by a small amount, how much will y change?

**Vector to Scalar**

$x \in \mathbb{R}^N, y \in \mathbb{R}$

Derivative is **Gradient**:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left(\frac{\partial y}{\partial x}\right)_n = \frac{\partial y}{\partial x_n}$$

For each element of x, if it changes by a small amount then how much will y change?
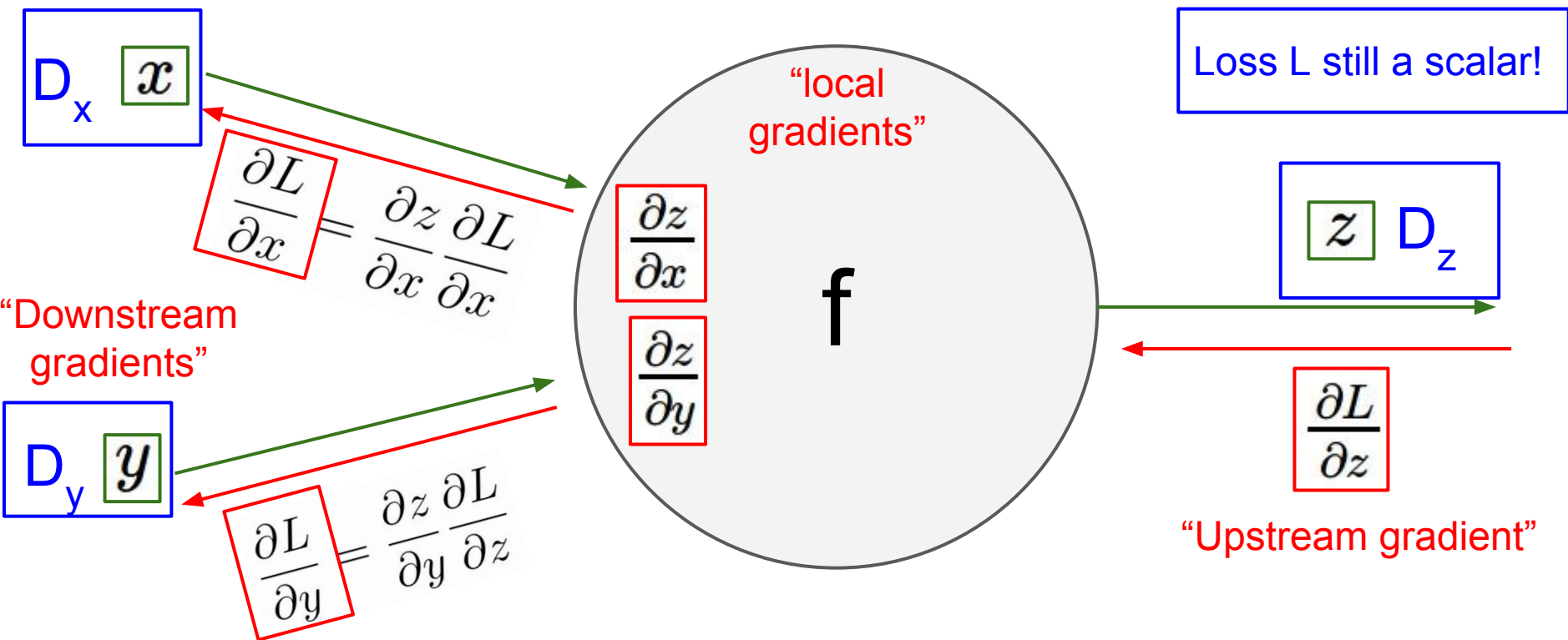
**Vector to Vector**

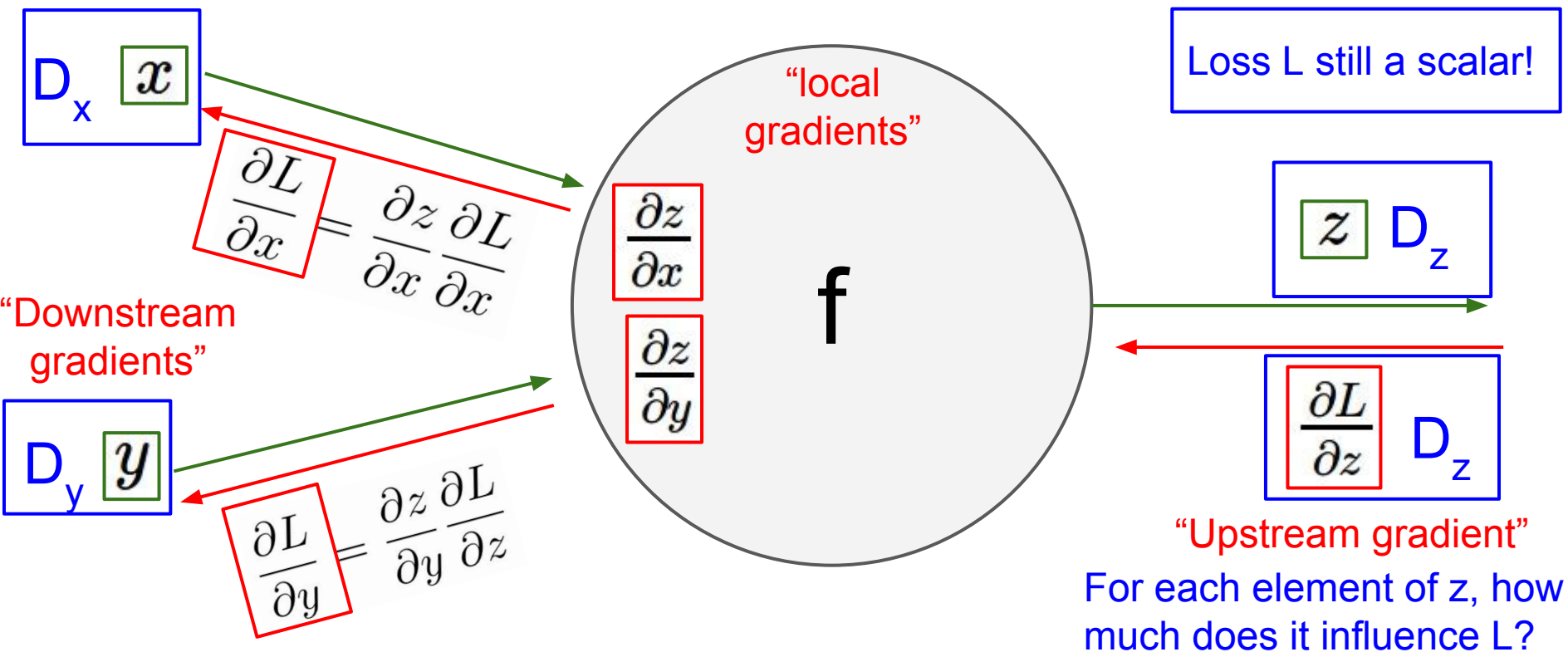$x \in \mathbb{R}^N, y \in \mathbb{R}^M$

Derivative is **Jacobian**:

$$\frac{\partial y}{\partial x} \in \mathbb{R}^{N \times M} \quad \left(\frac{\partial y}{\partial x}\right)_{n,m} = \frac{\partial y_m}{\partial x_n}$$

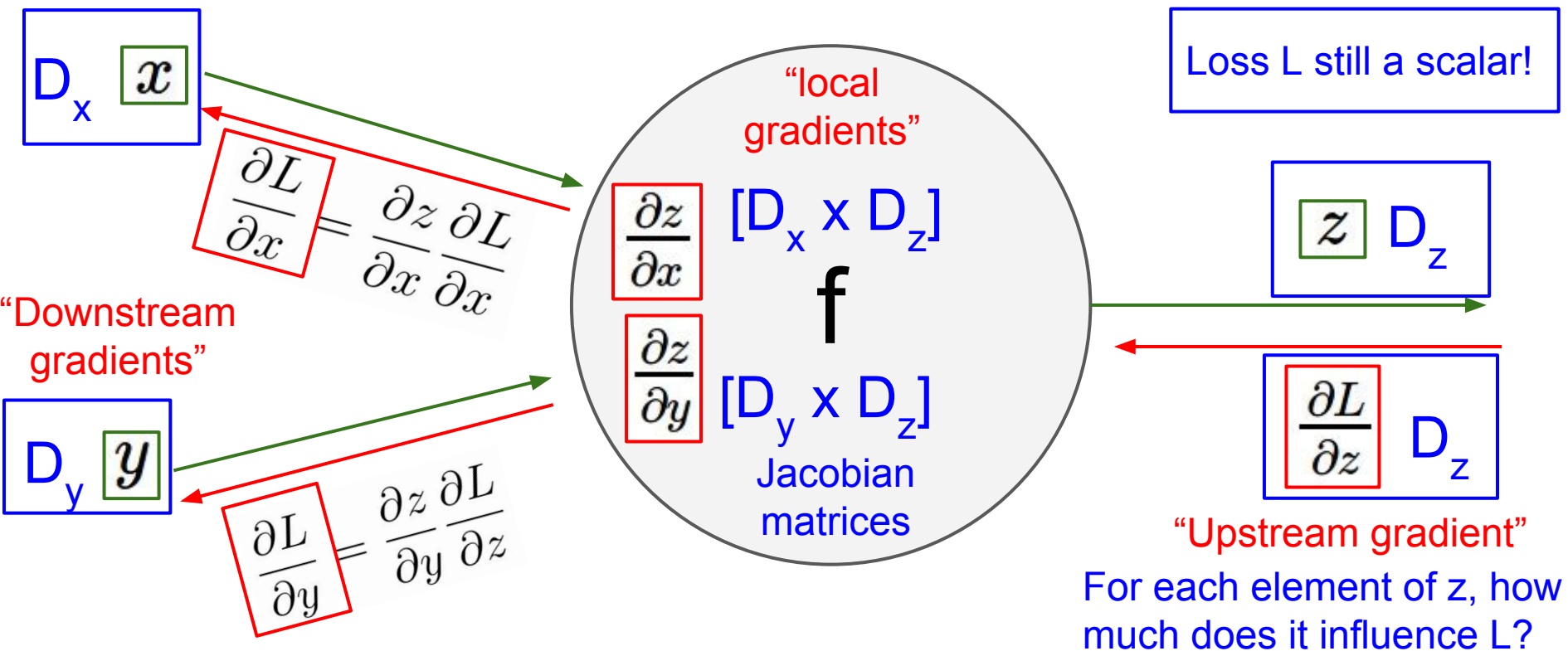For each element of x, if it changes by a small amount then how much will each element of y change?

# Backprop with Vectors



$D_x$ $x$

"local gradients"
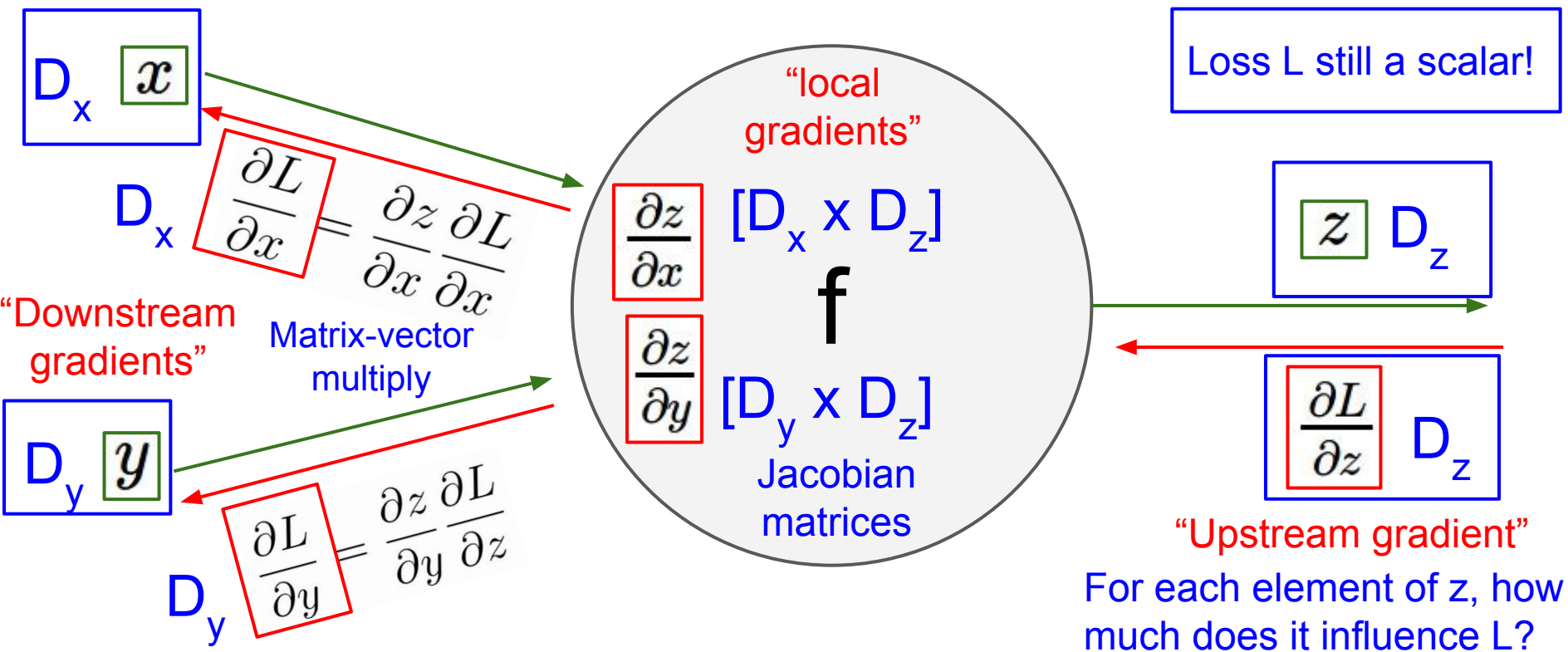
Loss L still a scalar!

$z$ $D_z$

$\dfrac{\partial L}{\partial x} = \dfrac{\partial z}{\partial x}\dfrac{\partial L}{\partial x}$

$\dfrac{\partial z}{\partial x}$

$\dfrac{\partial z}{\partial y}$

f

"Downstream gradients"

$D_y$ $y$

$\dfrac{\partial L}{\partial y} = \dfrac{\partial z}{\partial y}\dfrac{\partial L}{\partial z}$

$\dfrac{\partial L}{\partial z}$

"Upstream gradient"

# Backprop with Vectors



$D_x$ $x$

"local gradients"

Loss L still a scalar!

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

$z$ $D_z$

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

f

"Downstream gradients"

$\frac{\partial z}{\partial y}$

$D_y$ $y$

$$\frac{\partial L}{\partial y} = \frac{\partial z}{\partial y}\frac{\partial L}{\partial z}$$

$\frac{\partial L}{\partial z}$ $D_z$

"Upstream gradient"

For each element of z, how much does it influence L?

# Backprop with Vectors



Loss L still a scalar!

"local gradients"

$$\frac{\partial z}{\partial x} \quad [D_x \times D_z]$$

$$f$$

$$\frac{\partial z}{\partial y} \quad [D_y \times D_z]$$

Jacobian matrices

$D_x$ $x$

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

"Downstream gradients"

$D_y$ $y$

$$\frac{\partial L}{\partial y} = \frac{\partial z}{\partial y} \frac{\partial L}{\partial z}$$

$z$ $D_z$

$$\frac{\partial L}{\partial z} \quad D_z$$

"Upstream gradient"
For each element of z, how much does it influence L?

# Backprop with Vectors



"local gradients"

$$\frac{\partial z}{\partial x} \quad [D_x \times D_z]$$

$$\frac{\partial z}{\partial y} \quad [D_y \times D_z]$$

Jacobian matrices

$f$

$D_x$ $x$

$D_x$ $\frac{\partial L}{\partial x} = \frac{\partial z}{\partial x}\frac{\partial L}{\partial x}$

"Downstream gradients"

Matrix-vector multiply

$D_y$ $y$

$D_y$ $\frac{\partial L}{\partial y} = \frac{\partial z}{\partial y}\frac{\partial L}{\partial z}$

Loss L still a scalar!

$z$ $D_z$

$\frac{\partial L}{\partial z}$ $D_z$

"Upstream gradient"

For each element of z, how much does it influence L?

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

f(x) = max(0,x)
*(elementwise)*

4D output y:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

# Backprop with Vectors

4D input x:
[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

f(x) = max(0,x)
*(elementwise)*

4D output y:
[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

4D dL/dy:
[ 4 ]
[ -1 ]
[ 5 ]
[ 9 ]

Upstream gradient

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

f(x) = max(0,x)
*(elementwise)*

4D output y:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

Jacobian dy/dx

[ 1 0 0 0 ]
[ 0 0 0 0 ]
[ 0 0 1 0 ]
[ 0 0 0 0 ]

4D dL/dy:

[ 4 ]
[ -1 ]
[ 5 ]
[ 9 ]

Upstream gradient

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

f(x) = max(0,x)
*(elementwise)*

4D output y:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

[dy/dx] [dL/dy]
[ 1 0 0 0 ] [ 4 ]
[ 0 0 0 0 ] [ -1 ]
[ 0 0 1 0 ] [ 5 ]
[ 0 0 0 0 ] [ 9 ]

4D dL/dy:

[ 4 ]
[ -1 ]
[ 5 ]
[ 9 ]

Upstream gradient

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

f(x) = max(0,x)
*(elementwise)*

4D output y:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

4D dL/dx:

[ 4 ]
[ 0 ]
[ 5 ]
[ 0 ]

[dy/dx] [dL/dy]

[ 1 0 0 0 ] [ 4 ]
[ 0 0 0 0 ] [ -1 ]
[ 0 0 1 0 ] [ 5 ]
[ 0 0 0 0 ] [ 9 ]

4D dL/dy:

[ 4 ]
[ -1 ]
[ 5 ]
[ 9 ]

Upstream gradient

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

Jacobian is **sparse**: off-diagonal entries always zero! Never **explicitly** form Jacobian -- instead use **implicit** multiplication

$f(x) = \max(0,x)$
*(elementwise)*

4D output y:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

4D dL/dx:

[ 4 ]
[ 0 ]
[ 5 ]
[ 0 ]

[dy/dx] [dL/dy]

[ 1 0 0 0 ] [ 4 ]
[ 0 0 0 0 ] [ -1 ]
[ 0 0 1 0 ] [ 5 ]
[ 0 0 0 0 ] [ 9 ]

4D dL/dy:

[ 4 ]
[ -1 ]
[ 5 ]
[ 9 ]

Upstream gradient

# Backprop with Vectors

4D input x:

[ 1 ]
[ -2 ]
[ 3 ]
[ -1 ]

f(x) = max(0,x)
*(elementwise)*

4D output y:

[ 1 ]
[ 0 ]
[ 3 ]
[ 0 ]

Jacobian is **sparse**: off-diagonal entries always zero! Never **explicitly** form Jacobian -- instead use **implicit** multiplication

4D dL/dx:          [dy/dx] [dL/dy]          4D dL/dy:

[ 4 ]←
[ 0 ]←
[ 5 ]←
[ 0 ]←

$$\left(\frac{\partial L}{\partial x}\right)_i = \begin{cases} \left(\frac{\partial L}{\partial y}\right)_i & \text{if } x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$

← [ 4 ] ←
← [ -1 ] ←
← [ 5 ] ←
← [ 9 ] ←

Upstream gradient

# Backprop with Matrices (or Tensors)

Loss L still a scalar!

dL/dx always has the same shape as x!

$[D_x \times M_x]$ $x$

$[D_x \times M_x]$ $\dfrac{\partial L}{\partial x} = \dfrac{\partial z}{\partial x}\dfrac{\partial L}{\partial x}$

"local gradients"

$\dfrac{\partial z}{\partial x}$

$z$ $[D_z \times M_z]$

"Downstream gradients"

Matrix-vector multiply

$\dfrac{\partial z}{\partial y}$

$\dfrac{\partial L}{\partial z}$ $[D_z \times M_z]$

$[D_y \times M_y]$ $y$

$[D_y \times M_y]$ $\dfrac{\partial L}{\partial y} = \dfrac{\partial z}{\partial y}\dfrac{\partial L}{\partial z}$

Jacobian matrices

"Upstream gradient"

For each element of y, how much does it influence each element of z?

For each element of z, how much does it influence L?

# Backprop with Matrices (or Tensors)

Loss L still a scalar!

dL/dx always has the same shape as x!

$[D_x \times M_x]$ $x$

$[D_x \times M_x]$ $\dfrac{\partial L}{\partial x} = \dfrac{\partial z}{\partial x}\dfrac{\partial L}{\partial x}$

"local gradients"

$\dfrac{\partial z}{\partial x}$ $[(D_x \times M_x) \times (D_z \times M_z)]$

$z$ $[D_z \times M_z]$

"Downstream gradients"

Matrix-vector multiply

$\dfrac{\partial z}{\partial y}$ $[(D_y \times M_y) \times (D_z \times M_z)]$

$[D_y \times M_y]$ $y$

Jacobian matrices

$\dfrac{\partial L}{\partial z}$ $[D_z \times M_z]$

$\dfrac{\partial L}{\partial y} = \dfrac{\partial z}{\partial y}\dfrac{\partial L}{\partial z}$

$[D_y \times M_y]$

For each element of y, how much does it influence each element of z?

"Upstream gradient"
For each element of z, how much does it influence L?

# Backprop with Matrices

x: [N×D]

[ 2 **1** -3 ]
[ -3  4  2 ]

w: [D×M]

[ 3 2 1 -1]
[ 2 1 3  2]
[ 3 2 1 -2]

Matrix Multiply

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

y: [N×M]

[ **13  9  -2  -6** ]
[  5  2  17  1 ]

dL/dy: [N×M]

[  2  3 -3  9 ]
[ -8  1  4  6 ]

Also see derivation in the course notes:
http://cs231n.stanford.edu/handouts/linear-backprop.pdf

# Backprop with Matrices

x: [N×D]
[ 2  **1** -3 ]
[ -3  4  2 ]

w: [D×M]
[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

**Matrix Multiply**

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

y: [N×M]
[**13  9  -2  -6**]
[ 5  2  17  1 ]

dL/dy: [N×M]
[ 2  3 -3  9 ]
[ -8  1  4  6 ]

**Jacobians**:
dy/dx: [(N×D)×(N×M)]
dy/dw: [(D×M)×(N×M)]

For a neural net we may have
N=64, D=M=4096
Each Jacobian takes 256 GB of memory!
Must work with them implicitly!

# Backprop with Matrices

x: [N×D]

[ 2  **1**  -3 ]
[ -3  4   2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

**Matrix Multiply**

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

**Q**: What parts of y are affected by one element of x?

y: [N×M]

[ **13  9  -2  -6** ]
[  5  2  17  1 ]

dL/dy: [N×M]

[  2  3 -3  9 ]
[ -8  1  4  6 ]

# Backprop with Matrices

x: [N×D]

[ 2 **1** -3 ]
[ -3 4 2 ]

w: [D×M]

[ 3 2 1 -1]
[ 2 1 3 2]
[ 3 2 1 -2]

**Matrix Multiply**

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

y: [N×M]

[ **13 9 -2 -6** ]
[ 5 2 17 1 ]

dL/dy: [N×M]

[ 2 3 -3 9 ]
[ -8 1 4 6 ]

**Q**: What parts of y are affected by one element of x?

**A**: $x_{n,d}$ affects the whole row $y_{n,\cdot}$

$$\frac{\partial L}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}}$$

# Backprop with Matrices

x: [N×D]

[ 2  1 -3 ]
[ -3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

**Matrix Multiply**

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

y: [N×M]

[ **13  9  -2  -6** ]
[  5  2  17  1 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

**Q**: What parts of y are affected by one element of x?
**A**: $x_{n,d}$ affects the whole row $y_{n,\cdot}$

**Q**: How much does $x_{n,d}$ affect $y_{n,m}$?

$$\frac{\partial L}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}}$$

# Backprop with Matrices

x: [N×D]
[ 2  **1** -3 ]
[ -3  4  2 ]

w: [D×M]
[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

Matrix Multiply

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

y: [N×M]
[ **13  9** **-2** **-6** ]
[ 5  2  17  1 ]

dL/dy: [N×M]
[ 2  3 -3  9 ]
[ -8  1  4  6 ]

**Q**: What parts of y are affected by one element of x?
**A**: $x_{n,d}$ affects the whole row $y_{n,\cdot}$

**Q**: How much does $x_{n,d}$ affect $y_{n,m}$?
**A**: $w_{d,m}$

$$\frac{\partial L}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} w_{d,m}$$

# Backprop with Matrices

x: [N×D]

[ 2  **1** -3 ]
[ -3  4  2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

**Matrix Multiply**

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

y: [N×M]

[ **13  9  -2  -6** ]
[  5  2  17  1 ]

dL/dy: [N×M]

[  2  3 -3  9 ]
[ -8  1  4  6 ]

**Q**: What parts of y are affected by one element of x?
**A**: $x_{n,d}$ affects the whole row $y_{n,\cdot}$

**Q**: How much does $x_{n,d}$ affect $y_{n,m}$?
**A**: $w_{d,m}$

[N×D]  [N×M] [M×D]

$$\frac{\partial L}{\partial x} = \left(\frac{\partial L}{\partial y}\right) w^T$$

$$\frac{\partial L}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} \frac{\partial y_{n,m}}{\partial x_{n,d}} = \sum_m \frac{\partial L}{\partial y_{n,m}} w_{d,m}$$

# Backprop with Matrices

x: [N×D]

[ 2  **1** -3 ]
[ -3  4   2 ]

w: [D×M]

[ 3  2  1 -1]
[ 2  1  3  2]
[ 3  2  1 -2]

**Matrix Multiply**

$$y_{n,m} = \sum_d x_{n,d} w_{d,m}$$

y: [N×M]

[ **13  9** -2 **-6** ]
[  5  2  17  1 ]

dL/dy: [N×M]

[ 2  3 -3  9 ]
[ -8  1  4  6 ]

By similar logic:

[N×D]  [N×M] [M×D]

$$\frac{\partial L}{\partial x} = \left(\frac{\partial L}{\partial y}\right) w^T$$

[D×M]  [D×N] [N×M]

$$\frac{\partial L}{\partial w} = x^T \left(\frac{\partial L}{\partial y}\right)$$

These formulas are easy to remember: they are the only way to make shapes match up!