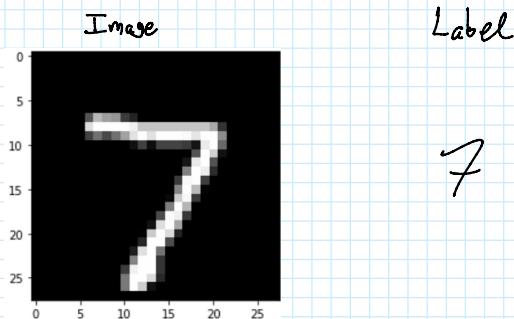


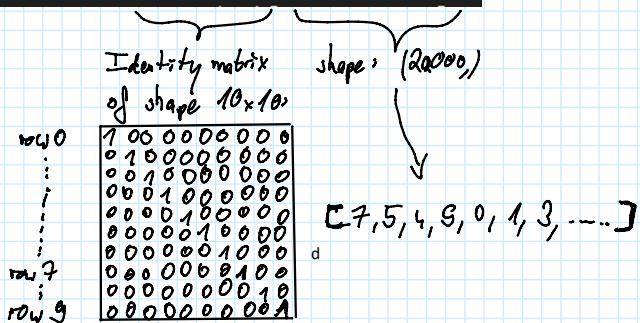
The dataset: MNIST

- Handwritten digit recognition dataset
- Consists of 70.000 labeled grayscale images
 - ↳ Resolution: 28×28 pixels, each pixel is value in $[0, 255]$
 - ↳ 60.000 training images, 10.000 test images
- Example:

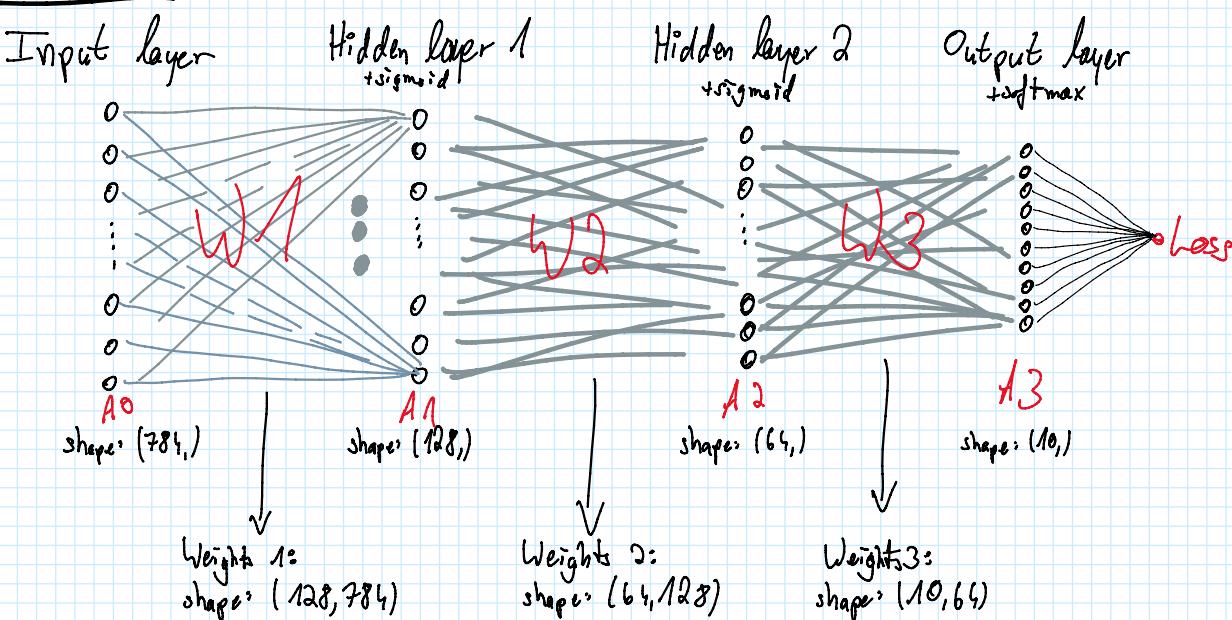


- Converting labels to one-hot representation:
- What is one-hot? $7 \rightarrow [0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 0 \ 1 \ 0 \ 0]^T$ shape: $(10,)$
- $(\text{shape}^2, 10)$

```
data_train_y_one_hot = np.eye(10)[data_train_y]
```



The model: MLP



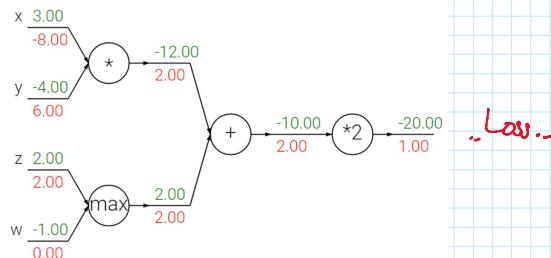
Softmax:

Softmax:

- Example: Input $\vec{z} = [1, 2, 3, 4, 1, 2, 3]$
Output $\vec{a} = [0.024, 0.064, 0.175, 0.475, 0.024, 0.064, 0.175]$
- Formula: $a_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$
- Properties:
 - $a_i \in [0, 1]$
 - $\sum_{i=1}^K a_i = 1$
 - Additive change to all \vec{z}_i 's doesn't change result.
(Not true for multiplication!)
 - Preserves ordering

Backpropagation:

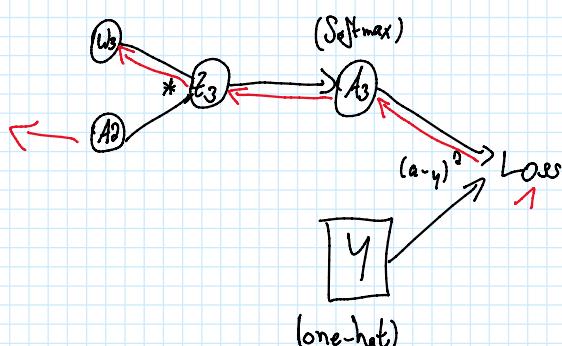
- Intuition:



Source: <https://cs231n.github.io/optimization-2/#grad>

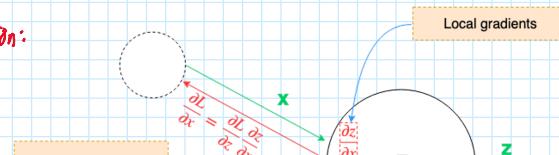
- Applying the intuition to our network:

Let's look at the last layer: $\frac{\partial L}{\partial w_3}$

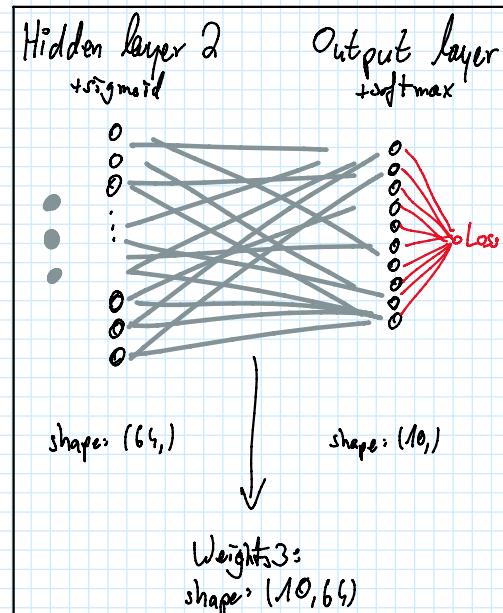


$$\hookrightarrow \text{We want: } \frac{\partial L}{\partial w_3} \Rightarrow \text{Chain rule: } \frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z_3} \cdot \frac{\partial z_3}{\partial w_3}$$

Illustration:

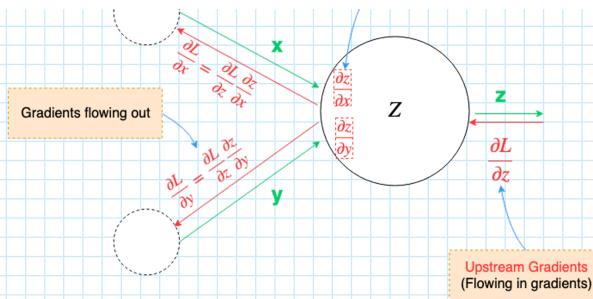


Reminder



Source: Stanford's CS231n course

Source: Stanford's CS231n course

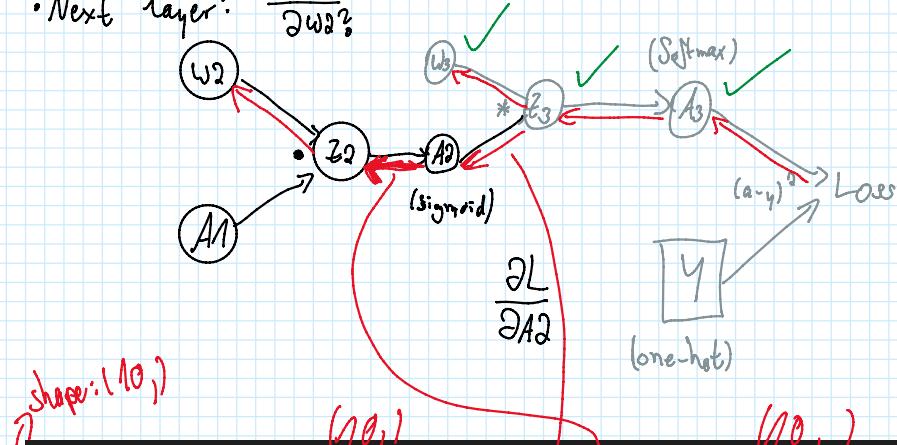


Calculating $\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z_3} \cdot \frac{\partial z_3}{\partial w_3}$:

```
# Calculate w3 update
error = [y * (predictions - data_y)] * [softmax_derivative(params['Z3'])]
change_w['W3'] = np.outer(error, params['A2'])
```

$$\frac{\partial L}{\partial w_3} = \begin{bmatrix} -x & -y & -z \\ x & y & z \end{bmatrix} \cdot \begin{bmatrix} A2 \\ w_3 \\ Z3 \end{bmatrix} = \begin{bmatrix} 1x + 2y + 3z \\ 4x + 5y + 6z \end{bmatrix}$$

• Next layer: $\frac{\partial L}{\partial w_2}$



```
# Calculate w2 update
error = [np.dot(params['W3'].T, error)] * [sigmoid_derivative(params['Z2'])]
change_w['W2'] = np.outer(error, params['A1'])
```

Extra Notes:

$$\bullet \text{Sigmoid derivative: } \sigma(x) = \frac{1}{1+e^{-x}} \Rightarrow \frac{d}{dx} \sigma(x) = \frac{d}{dx} (1+e^{-x})^{-1} = -\frac{d/dx (1+e^{-x})}{(1+e^{-x})^2} = -\frac{-e^{-x}}{(1+e^{-x})^2}$$

Python's
zip
function:

$x = [1, 2, 3, 4]$
 $y = ['a', 'b', 'c', 'd']$
 $\hookrightarrow \text{list(zip}(x, y)) = [(1, 'a'), (2, 'b'), (3, 'c'), (4, 'd')]$
 $\hookrightarrow \text{List of tuples}$

$\text{zip} \left(\begin{array}{|c|c|} \hline x & y \\ \hline \boxed{1} & \boxed{a} \\ \boxed{2} & \boxed{b} \\ \boxed{3} & \boxed{c} \\ \boxed{4} & \boxed{d} \\ \hline \end{array}, \quad \right)$

$= [(1, 'a'), (2, 'b'), (3, 'c'), (4, 'd')]$