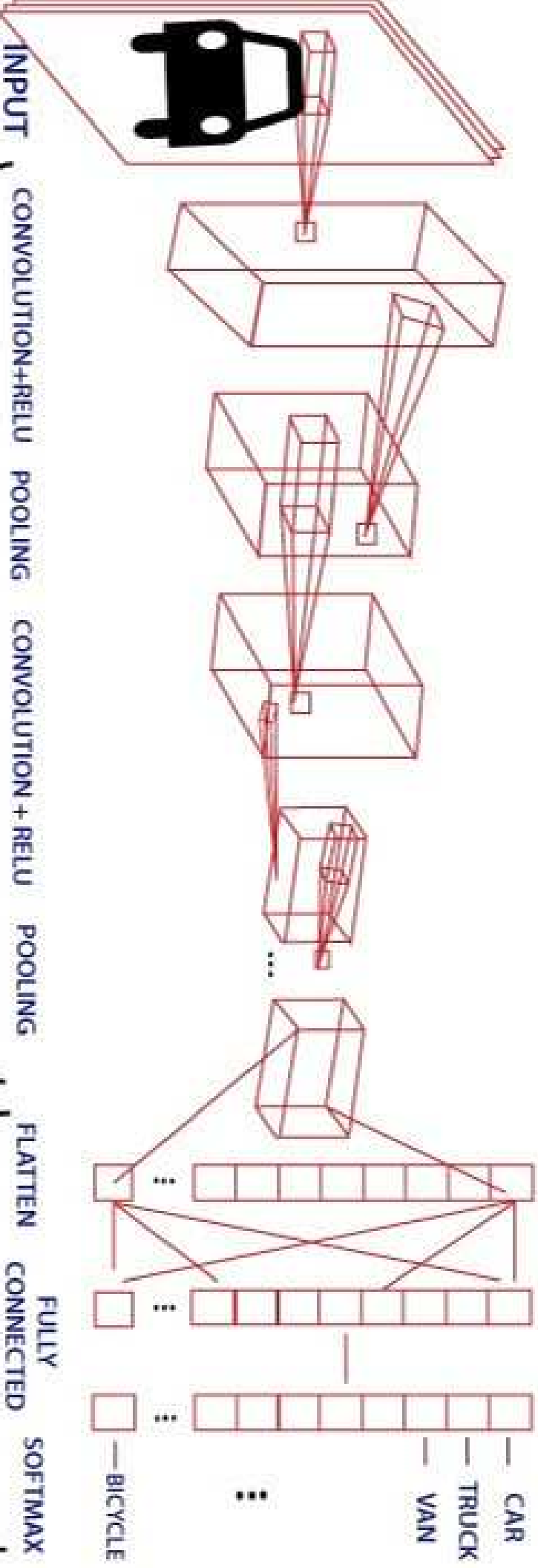


convolutional neural network

پروژه محاسبات علمی

نام و نام خانوادگی : صالح زارع زاده

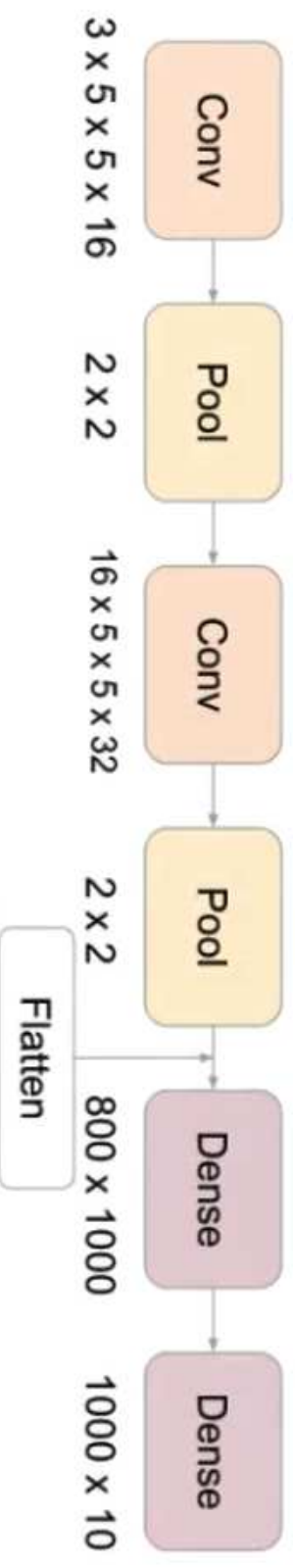
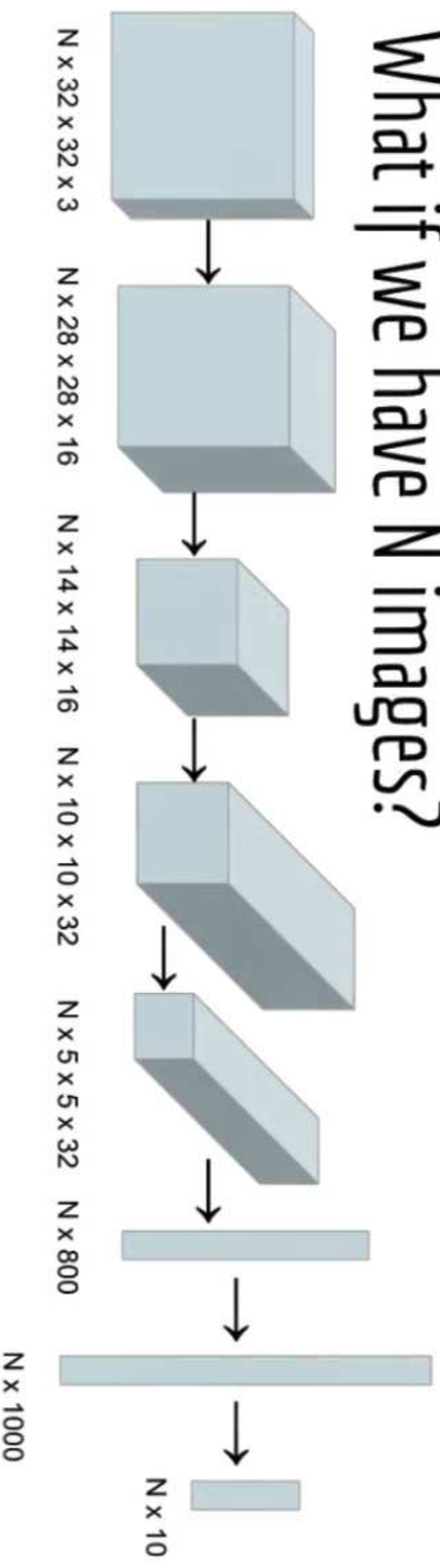
شماره دانشجویی : ۶۱۰۳۹۶۱۰۹



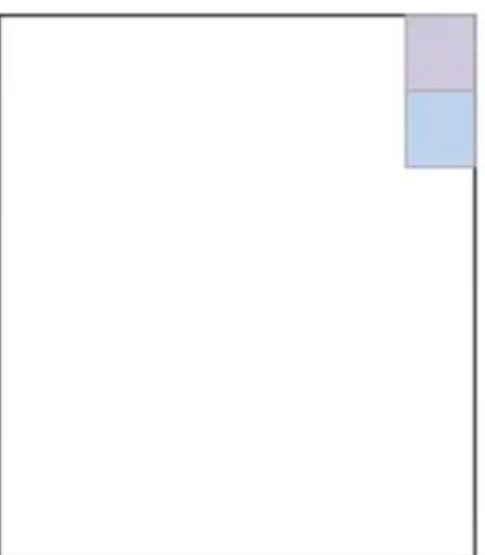
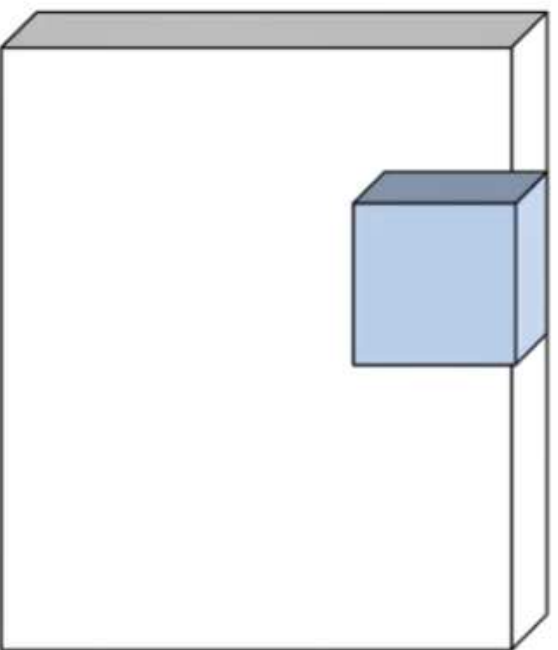
FEATURE LEARNING

CLASSIFICATION

What if we have N images?



Convolution with 3-D input



$$out[i,j] = \sum_{k_1=1}^K \sum_{k_2=1}^K \sum_{c=1}^C in[i-k_1, j-k_2, c] filter[c, k_1, k_2]$$

Max Pooling

29	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

2 x 2
pool size



100	184
12	45

Average Pooling

31	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

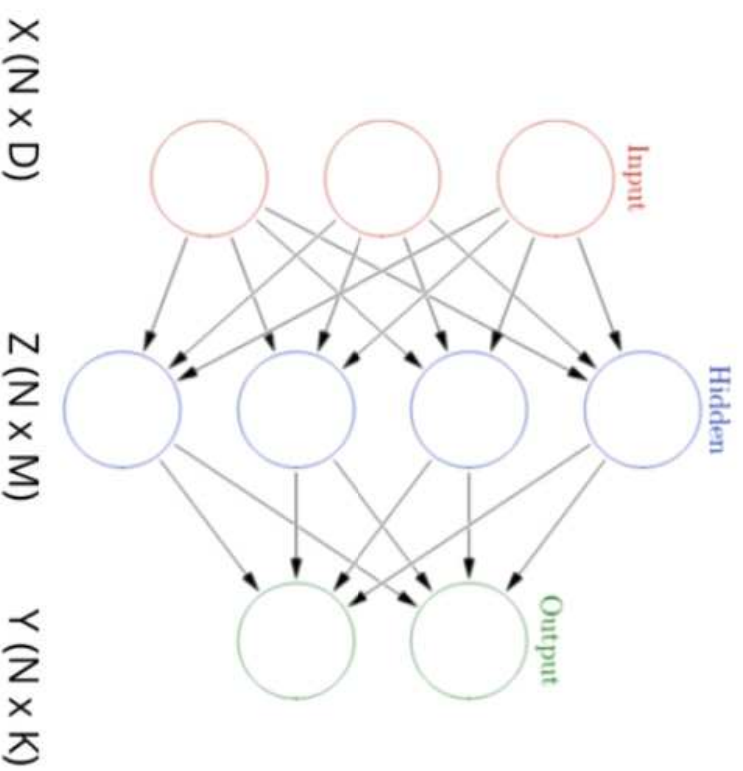
2 x 2
pool size



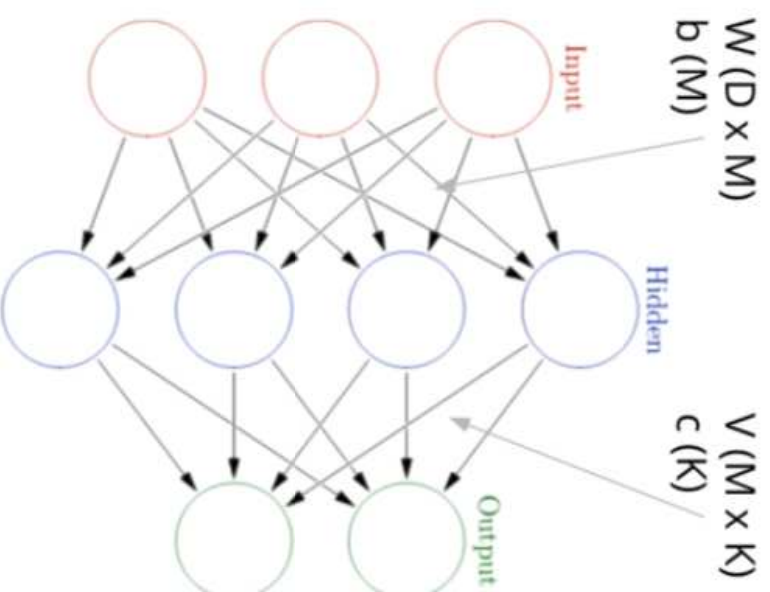
36	80
12	15

Backpropagation

- Let's continue our discussion on how to train a neural network



Backpropagation



Neural Network Equations

- Goal is still the same: build up our cost function, find gradients
- Gradient ascent (maximize) / Gradient descent (minimize)

Input \rightarrow Hidden

$$z = \sigma(W^T x + b)$$

Hidden \rightarrow Output

$$y = \textit{softmax}(V^T z + c)$$

Output \rightarrow Loss

$$J = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log y_{nk}$$

Chain Rule

- Let's start with gradients of V and c
- As before, we split it into 3 separate derivatives using the chain rule

$$J = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log y_{nk}$$

$$y_n = \text{softmax}(a_n)$$

$$a_{nk} = V_{:k}^T z_n + c_k$$

$$\frac{\partial J}{\partial V}_{nk} = \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a_{nk}} \frac{\partial a_{nk}}{\partial V_{mk}}$$

$$\frac{\partial J}{\partial c_k} = \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a_{nk}} \frac{\partial a_{nk}}{\partial c_k}$$

Gradients of V and c

$$\partial J / \partial V_{mk} = \sum_{n=1}^N \sum_{k'=1}^K \boxed{\frac{\partial J_{nk'}}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a_{nk}}} \frac{\partial a_{nk}}{\partial V_{mk}}$$

$$\partial J / \partial c_k = \sum_{n=1}^N \sum_{k'=1}^K \boxed{\frac{\partial J_{nk'}}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a_{nk}}} \frac{\partial a_{nk}}{\partial c_k}$$

- These gradients are the ones we already found in the Logistic Regression lectures! Therefore, no need to derive them again

Gradients of V and c

- The third derivatives are easy to derive (just linear terms)

$$\frac{\partial J}{\partial V}_{mk} = \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J_{nk'}}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a_{nk}} \boxed{\frac{\partial a_{nk}}{\partial V_{mk}}}$$

$$\frac{\partial J}{\partial c_k} = \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J_{nk'}}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a_{nk}} \boxed{\frac{\partial a_{nk}}{\partial c_k}}$$

$$a_{nk} = V_{:k}^T z_n + c_k \quad \longrightarrow \quad \frac{\partial a_{nk}}{\partial V_{mk}} = z_{nm}, \quad \frac{\partial a_{nk}}{\partial c_k} = 1$$

Put it all together

$$\partial J / \partial V_{mk} = \sum_{n=1}^N (t_{nk} - y_{nk}) z_{nm}$$

$$\partial J / \partial c_k = \sum_{n=1}^N (t_{nk} - y_{nk})$$

Vectorize

$$\nabla_V J = Z^T (T - Y)$$

```
grad_c = np.sum(T - Y, axis=0)
```

Separate activations

$$\alpha = W^T x + b$$

$$z = \sigma(\alpha)$$

$$a = V^T z + c$$

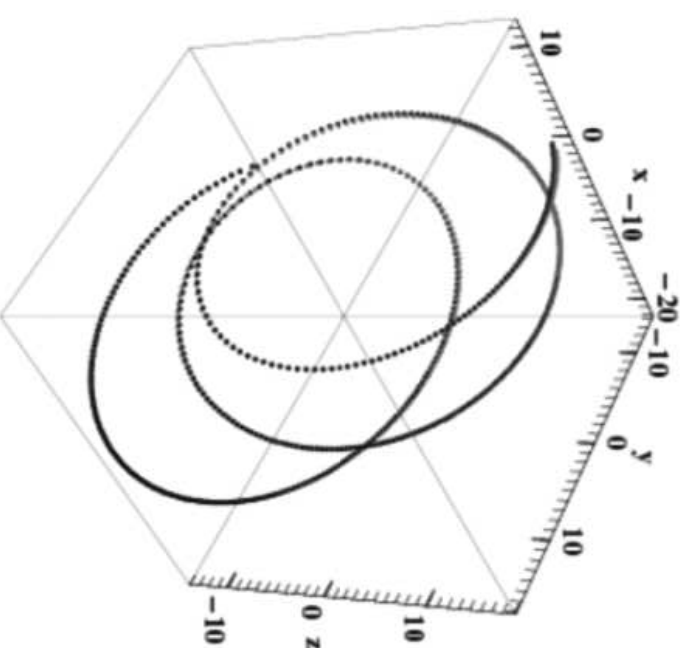
$$y = \textit{softmax}(a)$$

$$J = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log y_{nk}$$

Law of total derivatives

- Suppose we are tracking the position of some particle in 3-D space at regular time steps t
- We have position = $x(t), y(t), z(t)$
- Suppose we have some function of position:
 $f(x, y, z)$ (e.g. potential energy)

$$\frac{df}{dt} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt} + \frac{\partial f}{\partial z} \frac{dz}{dt}$$



Law of total derivatives

- We can use the same rule for any number of variables parameterized by t

$x_k(t)$ for $k = 1 \dots K$

$$\frac{df}{dt} = \sum_{k=1}^K \frac{\partial f}{\partial x_k} \frac{dx_k}{dt}$$

Applying the law of total derivatives

- Now it makes sense: the dummy index k goes away because we sum over it, and it doesn't appear on the LHS

$$\frac{\partial J}{\partial W_{dm}} = \sum_{k=1}^K \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a_{nk}} \frac{\partial a_{nk}}{\partial z_{nm}} \frac{\partial z_{nm}}{\partial \alpha_{nm}} \frac{\partial \alpha_{nm}}{\partial W_{dm}}$$

$$\frac{\partial J}{\partial b_m} = \sum_{k=1}^K \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a_{nk}} \frac{\partial a_{nk}}{\partial z_{nm}} \frac{\partial z_{nm}}{\partial \alpha_{nm}} \frac{\partial \alpha_{nm}}{\partial b_m}$$

We solved some of these already!

$$\frac{\partial J}{\partial W}_{dm} = \sum_{k=1}^K \sum_{n=1}^N \sum_{k'=1}^K \boxed{\frac{\partial J_{nk'}}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a_{nk}}} \frac{\partial a_{nk}}{\partial z_{nm}} \frac{\partial z_{nm}}{\partial \alpha_{nm}} \frac{\partial \alpha_{nm}}{\partial W}_{dm}$$

$$\frac{\partial J}{\partial W}_{dm} = \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) \frac{\partial a_{nk}}{\partial z_{nm}} \frac{\partial z_{nm}}{\partial \alpha_{nm}} \frac{\partial \alpha_{nm}}{\partial W}_{dm}$$

3 more derivatives

$$\frac{\partial J}{\partial W_{dm}} = \sum_{k=1}^K \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a_{nk}}$$

$$\frac{\partial a_{nk}}{\partial z_{nm}} \frac{\partial z_{nm}}{\partial \alpha_{nm}} \frac{\partial \alpha_{nm}}{\partial W_{dm}}$$

$$a_{nk} = V_{:k}^T z_n + c_k$$

$$z_{nm} = \sigma(\alpha_{nm})$$

$$\alpha_{nm} = W_{:m}^T x_n + b_m$$

$$\frac{\partial a_{nk}}{\partial z_{nm}} = V_{mk}, \quad \frac{\partial \alpha_{nm}}{\partial W_{dm}} = x_{nd}$$

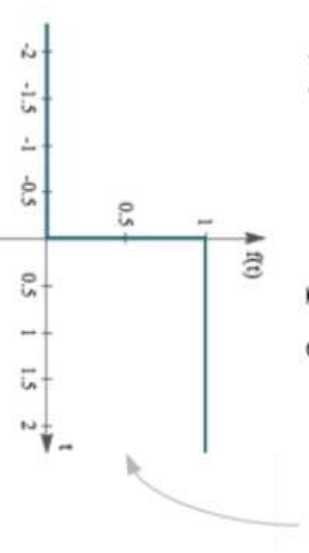
Activation function derivative

- For these lectures we'll assume sigmoid (mostly)

$$\frac{\partial z_{nm}}{\partial \alpha_{nm}} = z_{nm}(1 - z_{nm}) \quad \text{if } \sigma(x) = \frac{1}{1 + \exp(-x)}$$

$$\frac{\partial z_{nm}}{\partial \alpha_{nm}} = 1 - z_{nm}^2 \quad \text{if } \sigma(x) = \tanh(x)$$

$$\frac{\partial z_{nm}}{\partial \alpha_{nm}} = u(z_{nm}) \quad \text{if } \sigma(x) = \text{relu}(x), \text{ where } u(\cdot) = \text{step function}$$



Put it all together

- As before, the derivative wrt bias term just replaces the input to the layer with 1

$$\partial J / \partial W_{dm} = \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) V_{mk} z_{nm} (1 - z_{nm}) x_{nd}$$

$$\partial J / \partial b_m = \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) V_{mk} z_{nm} (1 - z_{nm})$$

Vectorize it

- The easiest part to do first is the derivative of the activation function which is just element-wise multiplication

$$\partial J / \partial W_{dm} = \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) V_{mk} \boxed{z_{nm}(1 - z_{nm})} x_{nd}$$

$$Z' = Z \odot (1 - Z)$$

Consider the shapes

$$\partial J/\partial W_{dm} = \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) V_{mk} Z_{nm} (1 - z_{nm}) x_{nd}$$

$N \times K$
 $M \times K$

$$(T - Y)_{N \times K} V_{K \times M}^T \rightarrow N \times M$$

Consider the shapes

$$\frac{\partial J}{\partial W_{dm}} = \sum_{k=1}^K \sum_{n=1}^N \left(t_{nk} - y_{nk} \right) V_{nk} Z_{nm} (1 - Z_{nm}) x_{nd}$$

$N \times M$
 $N \times M$

$$\left[(T - Y) V^T \right]_{N \times M} \odot Z'_{N \times M}$$

Consider the shapes

$$\frac{\partial J}{\partial W}_{dm} = \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) V_{mk} z_{nm} (1 - z_{nm}) x_{nd}$$

$N \times M$
 $N \times D$

$$\nabla_W J = X^T \{ [(T - Y) V^T] \odot Z \odot (1 - Z) \}$$

$D \times M$

With a generic activation function

- Sigmoid

$$\nabla_W J = X^T \left\{ \left[(T - Y) V^T \right] \odot Z \odot (1 - Z) \right\}$$

- Generic

$$\nabla_W J = X^T \left\{ \left[(T - Y) V^T \right] \odot Z' \right\}$$

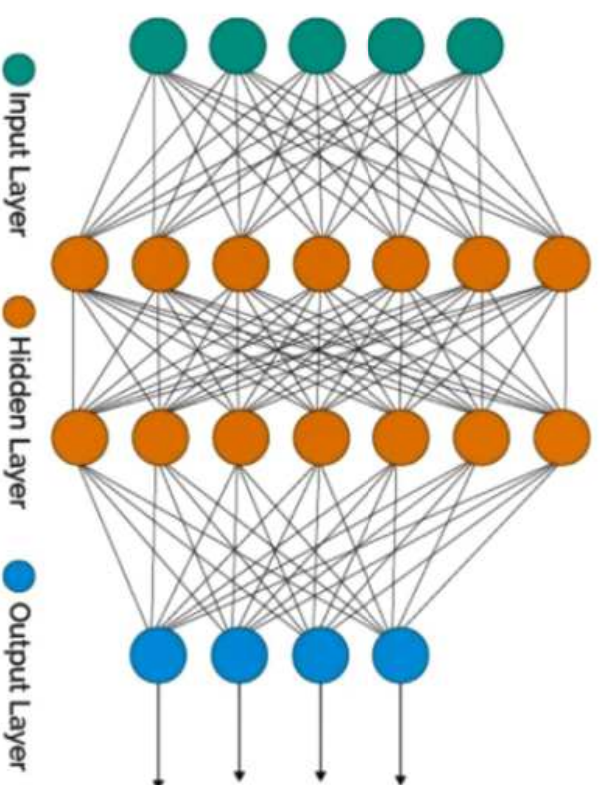
Vectorized bias term gradient

```
grad_b = np.sum ( (T-Y) .dot (V.T) * Z * (1 - Z) , axis=0)
```

- Note: the same thing we had for W, except without X (or equivalently, if X was just 1)

More Layers!

- Can we work our way up to an arbitrary number of layers?
- In order to help us more clearly see the patterns, let's first consider a neural network with 3 layers (of weights)



New Symbols

- We are running out of letters, let's use numbered superscripts instead

$$a^{(1)} = W^{(1)T} x + b^{(1)}$$

$$z^{(1)} = \sigma(a^{(1)})$$

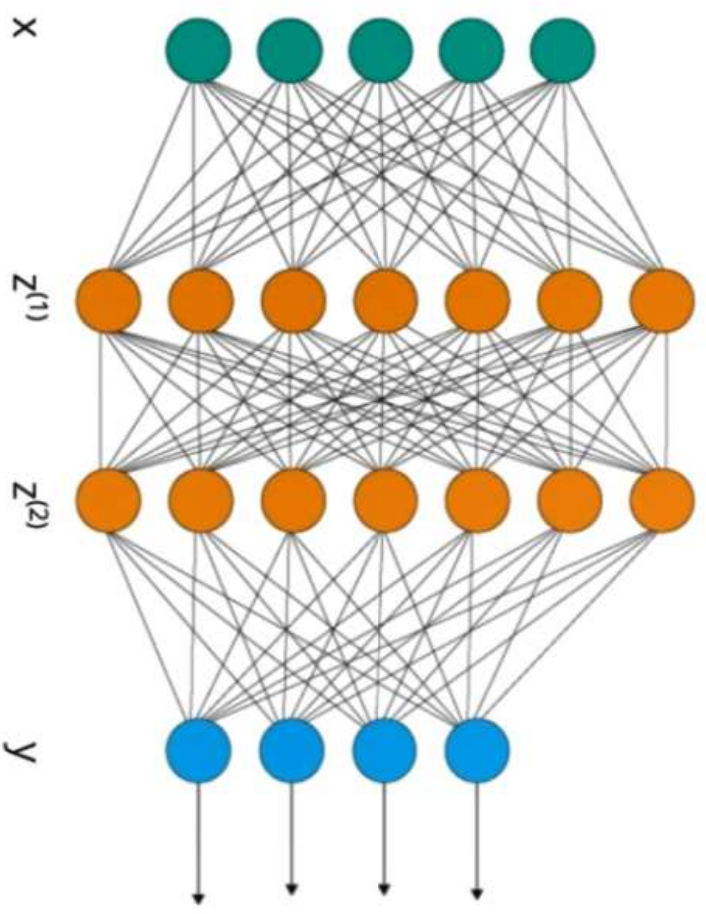
$$a^{(2)} = W^{(2)T} z^{(1)} + b^{(2)}$$

$$z^{(2)} = \sigma(a^{(2)})$$

$$a^{(3)} = W^{(3)T} z^{(2)} + b^{(3)}$$

$$y = \text{softmax}(a^{(3)})$$

Think of $x = z^{(0)}$, $y = z^{(3)}$



New Symbols

- Superscript weights also - the general pattern is:
 $\text{size}(W^{(l)}) = M^{(l-1)} \times M^{(l)}$

$$\text{size}(W^{(1)}) = D \times M^{(1)}, \text{size}(b^{(1)}) = M^{(1)}$$

$$\text{size}(W^{(2)}) = M^{(1)} \times M^{(2)}, \text{size}(b^{(2)}) = M^{(2)}$$

$$\text{size}(W^{(3)}) = M^{(2)} \times K, \text{size}(b^{(3)}) = K$$

$$\text{Think of } D = M^{(0)}, K = M^{(3)}$$

Gradients in last layer

- No different than what we saw before - just with new symbols!

$$\frac{\partial J}{\partial W^{(3)}_{m^{(2)}k}} = \sum_{n=1}^N \sum_{k'=1}^K \frac{\frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a^{(3)}_{nk}} \frac{\partial a^{(3)}_{nk}}{\partial W^{(3)}_{m^{(2)}k}}}{}$$

$$\frac{\partial J}{\partial b^{(3)}_k} = \sum_{n=1}^N \sum_{k'=1}^K \frac{\frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a^{(3)}_{nk}} \frac{\partial a^{(3)}_{nk}}{\partial b^{(3)}_k}}{}$$

Gradients in last layer

- Important: the answer hasn't changed, only the symbols are different

$$\partial J / \partial W^{(3)}_{m^{(2)}k} = \sum_{n=1}^N (t_{nk} - y_{nk}) z^{(2)}_{nm^{(2)}}$$

$$\partial J / \partial b^{(3)}_k = \sum_{n=1}^N (t_{nk} - y_{nk})$$

Gradients in the 2nd layer

- More terms (due to the chain rule) - but still: the answer hasn't changed
- Remember: law of total derivatives

$$\frac{\partial J}{\partial W^{(2)}}_{m^{(1)}m^{(2)}} = \sum_{k=1}^K \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a^{(3)}_{nk}} \frac{\partial a^{(3)}_{nk}}{\partial z^{(2)}_{nm^{(2)}}} \frac{\partial z^{(2)}_{nm^{(2)}}}{\partial a^{(2)}_{nm^{(2)}}} \frac{\partial a^{(2)}_{nm^{(2)}}}{\partial W^{(2)}_{m^{(1)}m^{(2)}}}$$

$$\frac{\partial J}{\partial b^{(2)}}_{m^{(2)}} = \sum_{k=1}^K \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a^{(3)}_{nk}} \frac{\partial a^{(3)}_{nk}}{\partial z^{(2)}_{nm^{(2)}}} \frac{\partial z^{(2)}_{nm^{(2)}}}{\partial a^{(2)}_{nm^{(2)}}} \frac{\partial a^{(2)}_{nm^{(2)}}}{\partial b^{(2)}_{m^{(2)}}}$$

Gradients in the 2nd layer

- Again assuming sigmoid activation function
- Since there's another layer behind us, where we saw x before we now see $z^{(1)}$

$$\frac{\partial J}{\partial W^{(2)}}_{m^{(1)}m^{(2)}} = \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) W^{(3)}_{m^{(2)}k} z^{(2)}_{nm^{(2)}} (1 - z^{(2)}_{nm^{(2)}}) z^{(1)}_{nm^{(1)}}$$

$$\frac{\partial J}{\partial b^{(2)}}_{m^{(2)}} = \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) W^{(3)}_{m^{(2)}k} z^{(2)}_{nm^{(2)}} (1 - z^{(2)}_{nm^{(2)}})$$

Gradients for the 1st layer

- The largest # of terms from the chain rule we've seen so far
- Law of total derivatives: must sum over k and $m^{(2)}$
- "Sum over any index that doesn't appear on LHS"

$$\frac{\partial J}{\partial W^{(1)}} \frac{dm^{(1)}}{dm^{(1)}} = \sum_{m^{(2)}=1}^{M^{(2)}} \sum_{k=1}^K \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a^{(3)}_{nk}} \frac{\partial a^{(3)}_{nk}}{\partial z^{(2)}_{nm^{(2)}}} \frac{\partial z^{(2)}_{nm^{(2)}}}{\partial a^{(2)}_{nm^{(2)}}} \frac{\partial a^{(2)}_{nm^{(2)}}}{\partial z^{(1)}_{nm^{(1)}}} \frac{\partial z^{(1)}_{nm^{(1)}}}{\partial a^{(1)}_{nm^{(1)}}} \frac{\partial a^{(1)}_{nm^{(1)}}}{\partial W^{(1)}} \frac{dm^{(1)}}{dm^{(1)}}$$

$$\frac{\partial J}{\partial b^{(1)}} \frac{dm^{(1)}}{m^{(1)}} = \sum_{m^{(2)}=1}^{M^{(2)}} \sum_{k=1}^K \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a^{(3)}_{nk}} \frac{\partial a^{(3)}_{nk}}{\partial z^{(2)}_{nm^{(2)}}} \frac{\partial z^{(2)}_{nm^{(2)}}}{\partial a^{(2)}_{nm^{(2)}}} \frac{\partial a^{(2)}_{nm^{(2)}}}{\partial z^{(1)}_{nm^{(1)}}} \frac{\partial z^{(1)}_{nm^{(1)}}}{\partial a^{(1)}_{nm^{(1)}}} \frac{\partial a^{(1)}_{nm^{(1)}}}{\partial b^{(1)}} \frac{dm^{(1)}}{m^{(1)}}$$

Gradients for the 1st layer

- 7 functions → 7 derivatives - go through this slowly by yourself

$$a^{(1)} = W^{(1)T} x + b^{(1)}$$

$$z^{(1)} = \sigma(a^{(1)})$$

$$a^{(2)} = W^{(2)T} z^{(1)} + b^{(2)}$$

$$z^{(2)} = \sigma(a^{(2)})$$

$$a^{(3)} = W^{(3)T} z^{(2)} + b^{(3)}$$

$$y = \text{softmax}(a^{(3)})$$

$$J = \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log y_{nk}$$

$$\frac{\partial J}{\partial W^{(1)}} \frac{dm^{(1)}}{dn^{(1)}} = \sum_{m^{(2)}=1}^{M^{(2)}} \sum_{k=1}^K \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a^{(3)}_{nk}} \frac{\partial a^{(3)}_{nk}}{\partial z^{(2)}_{nm^{(2)}}} \frac{\partial z^{(2)}_{nm^{(2)}}}{\partial a^{(2)}_{nm^{(2)}}} \frac{\partial a^{(2)}_{nm^{(2)}}}{\partial z^{(1)}_{nm^{(1)}}} \frac{\partial z^{(1)}_{nm^{(1)}}}{\partial a^{(1)}_{nm^{(1)}}} \frac{\partial a^{(1)}_{nm^{(1)}}}{\partial W^{(1)}_{nm^{(1)}}} \frac{dm^{(1)}}{dn^{(1)}}$$

$$\frac{\partial J}{\partial b^{(1)}} \frac{dm^{(1)}}{dn^{(1)}} = \sum_{m^{(2)}=1}^{M^{(2)}} \sum_{k=1}^K \sum_{n=1}^N \sum_{k'=1}^K \frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a^{(3)}_{nk}} \frac{\partial a^{(3)}_{nk}}{\partial z^{(2)}_{nm^{(2)}}} \frac{\partial z^{(2)}_{nm^{(2)}}}{\partial a^{(2)}_{nm^{(2)}}} \frac{\partial a^{(2)}_{nm^{(2)}}}{\partial z^{(1)}_{nm^{(1)}}} \frac{\partial z^{(1)}_{nm^{(1)}}}{\partial a^{(1)}_{nm^{(1)}}} \frac{\partial a^{(1)}_{nm^{(1)}}}{\partial b^{(1)}_{nm^{(1)}}} \frac{dm^{(1)}}{dn^{(1)}}$$

Gradients for the 1st layer

- We've solved all of these derivatives already
- Same form but different subscripts
- E.g. $W^{(1)T}x$ and $W^{(2)T}z^{(1)}$ have the same form
- For compactness, I've replaced activation function derivative with z'

$$\frac{\partial J}{\partial W^{(1)}}_{dm^{(1)}} = \sum_{m^{(2)=1}}^{M^{(2)}} \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) W^{(3)}_{m^{(2)}k} z^{(2)'}_{nm^{(2)}} W^{(2)}_{m^{(1)}m^{(2)}} z^{(1)'}_{nm^{(1)}} x_{nd}$$

$$\frac{\partial J}{\partial b^{(1)}}_{m^{(1)}} = \sum_{m^{(2)=1}}^{M^{(2)}} \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) W^{(3)}_{m^{(2)}k} z^{(2)'}_{nm^{(2)}} W^{(2)}_{m^{(1)}m^{(2)}} z^{(1)'}_{nm^{(1)}}$$

Patterns to notice

- First layer gradient had 7 derivatives
- Second layer gradient had 5 derivatives
- Output layer gradient had 3 derivatives
- At each layer, we add 2 functions:
 - Linear transformation
 - Activation function

$$a^{(l)} = W^{(l)T} z^{(l-1)} + b^{(l)}$$

$$z^{(l)} = \sigma(a^{(l)})$$

Patterns to notice

- The further we go back, the more common terms we encounter

$$\frac{\partial J}{\partial W^{(3)}}_{m^{(2)}k} = \sum_{n=1}^N \sum_{k'=1}^K \boxed{\frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a^{(3)}_{nk}}} \frac{\partial a^{(3)}_{nk}}{\partial W^{(3)}_{m^{(2)}k}}$$

$$\frac{\partial J}{\partial W^{(2)}}_{m^{(1)}m^{(2)}} = \sum_{k=1}^K \sum_{n=1}^N \sum_{k'=1}^K \boxed{\frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a^{(3)}_{nk}}} \frac{\partial a^{(3)}_{nk}}{\partial z^{(2)}_{nm^{(2)}}} \frac{\partial z^{(2)}_{nm^{(2)}}}{\partial a^{(2)}_{nm^{(2)}}} \frac{\partial a^{(2)}_{nm^{(2)}}}{\partial W^{(2)}_{m^{(1)}m^{(2)}}}$$

$$\frac{\partial J}{\partial W^{(1)}}_{dm^{(1)}} = \sum_{m^{(2)}=1}^{M^{(2)}} \sum_{k=1}^K \sum_{n=1}^N \sum_{k'=1}^K \boxed{\frac{\partial J}{\partial y_{nk'}} \frac{\partial y_{nk'}}{\partial a^{(3)}_{nk}}} \frac{\partial a^{(3)}_{nk}}{\partial z^{(2)}_{nm^{(2)}}} \frac{\partial z^{(2)}_{nm^{(2)}}}{\partial a^{(2)}_{nm^{(2)}}} \frac{\partial a^{(2)}_{nm^{(2)}}}{\partial z^{(1)}_{nm^{(1)}}} \frac{\partial z^{(1)}_{nm^{(1)}}}{\partial a^{(1)}_{nm^{(1)}}} \frac{\partial a^{(1)}_{nm^{(1)}}}{\partial W^{(1)}_{dm^{(1)}}}$$

Patterns to notice

- The same pattern seems to be followed: input to the layer multiplied by some other term
- Let's call it "delta" for convenience

$$\nabla_{W^{(l)}} J = Z^{(l-1)T} \delta^{(l)}$$

Patterns to notice

- Let's consider the weight updates for a 3-layer neural network to "guess" what δ should be

$$\partial J / \partial W^{(3)}_{m^{(2)}k} = \sum_{n=1}^N (t_{nk} - y_{nk}) z_{nm^{(2)}}^{(2)}$$

$$\partial J / \partial W^{(2)}_{m^{(1)}m^{(2)}} = \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) W^{(3)}_{m^{(2)}k} z_{nm^{(2)}}^{(2)'} z_{nm^{(1)}}^{(1)}$$

$$\partial J / \partial W^{(1)}_{dm^{(1)}} = \sum_{m^{(2)}=1}^{M^{(2)}} \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) W^{(3)}_{m^{(2)}k} z_{nm^{(2)}}^{(2)'} W^{(2)}_{m^{(1)}m^{(2)}} z_{nm^{(1)}}^{(1)'} x_{nd}$$

Delta

$$\frac{\partial J}{\partial W^{(3)}}_{m^{(2)}k} = \sum_{n=1}^N (t_{nk} - y_{nk}) z_{nm^{(2)}}^{(2)} \delta_{nk}^{(3)} = t_{nk} - y_{nk}$$

$$\frac{\partial J}{\partial W^{(2)}}_{m^{(1)}m^{(2)}} = \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) W^{(3)}_{m^{(2)}k} z_{nm^{(2)}}^{(2)'} z_{nm^{(1)}}^{(1)}$$

$$\frac{\partial J}{\partial W^{(1)}}_{dm^{(1)}} = \sum_{m^{(2)}=1}^{M^{(2)}} \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) W^{(3)}_{m^{(2)}k} z_{nm^{(2)}}^{(2)'} W^{(2)}_{m^{(1)}m^{(2)}} z_{nm^{(1)}}^{(1)'} x_{nd}$$

Delta

$$\delta_{nm^{(2)}}^{(2)} = \sum_{k=1}^K (t_{nk} - y_{nk}) W_{m^{(2)}k}^{(3)} z_{nm^{(2)}}^{(2)'}$$

$$\partial J/\partial W_{m^{(2)}k}^{(3)} = \sum_{n=1}^N (t_{nk} - y_{nk}) z_{nm^{(2)}}^{(2)}$$

$$\partial J/\partial W_{m^{(1)}m^{(2)}}^{(2)} = \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) W_{m^{(2)}k}^{(3)} z_{nm^{(2)}}^{(2)'} z_{nm^{(1)}}^{(1)}$$

$$\partial J/\partial W_{dm^{(1)}}^{(1)} = \sum_{m^{(2)}=1}^{M^{(2)}} \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) W_{m^{(2)}k}^{(3)} z_{nm^{(2)}}^{(2)'} W_{m^{(1)}m^{(2)}}^{(2)} z_{nm^{(1)}}^{(1)'} x_{nd}$$

Delta

$$\delta_{nm^{(1)}}^{(1)} = \sum_{m^{(2)}=1}^{M^{(2)}} \sum_{k=1}^K (t_{nk} - y_{nk}) W^{(3)}_{m^{(2)}k} z^{(2)'}_{nm^{(2)}} W^{(2)}_{m^{(1)}m^{(2)}} z^{(1)'}_{nm^{(1)}}$$

$$\partial J / \partial W^{(3)}_{m^{(2)}k} = \sum_{n=1}^N (t_{nk} - y_{nk}) z^{(2)}_{nm^{(2)}}$$

$$\partial J / \partial W^{(2)}_{m^{(1)}m^{(2)}} = \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) W^{(3)}_{m^{(2)}k} z^{(2)'}_{nm^{(2)}} z^{(1)}_{nm^{(1)}}$$

$$\partial J / \partial W^{(1)}_{dm^{(1)}} = \sum_{m^{(2)}=1}^{M^{(2)}} \sum_{k=1}^K \sum_{n=1}^N (t_{nk} - y_{nk}) W^{(3)}_{m^{(2)}k} z^{(2)'}_{nm^{(2)}} W^{(2)}_{m^{(1)}m^{(2)}} z^{(1)'}_{nm^{(1)}} x_{nd}$$

Delta

Can you see the pattern?

$$\delta_{nk}^{(3)} = t_{nk} - y_{nk}$$

$$\delta_{nm^{(2)}}^{(2)} = \sum_{k=1}^K (t_{nk} - y_{nk}) W_{m^{(2)}k}^{(3)} z_{nm^{(2)}}^{(2)'} \quad$$

$$\delta_{nm^{(1)}}^{(1)} = \sum_{m^{(2)}=1}^{M^{(2)}} \sum_{k=1}^K (t_{nk} - y_{nk}) W_{m^{(2)}k}^{(3)} z_{nm^{(2)}}^{(2)'} W_{m^{(1)}m^{(2)}}^{(2)} z_{nm^{(1)}}^{(1)'} \quad$$

Delta

- Define delta in terms of delta

$$\delta_{nk}^{(3)} = t_{nk} - y_{nk}$$

$$\delta_{nm^{(2)}}^{(2)} = \sum_{k=1}^K \delta_{nk}^{(3)} W_{m^{(2)}k}^{(2)} z_{nm^{(2)}}^{(2)'}$$

$$\delta_{nm^{(1)}}^{(1)} = \sum_{m^{(2)}=1}^{M^{(2)}} \delta_{nm^{(2)}}^{(2)} W_{m^{(1)}m^{(2)}}^{(2)} z_{nm^{(1)}}^{(1)'}$$

Delta Recursion

$$\delta_{nk}^{(L)} = t_{nk} - y_{nk}$$

$$\delta_{nm^{(l)}}^{(l)} = \sum_{m^{(l+1)}=1}^{M^{(l+1)}} \delta_{nm^{(l+1)}}^{(l+1)} W_{m^{(l)}m^{(l+1)}}^{(l+1)} Z_{nm^{(l)}}^{(l)'} , for \, l = 1 \dots L - 1$$

Vectorized Delta

$$\delta^{(L)} = T - Y$$

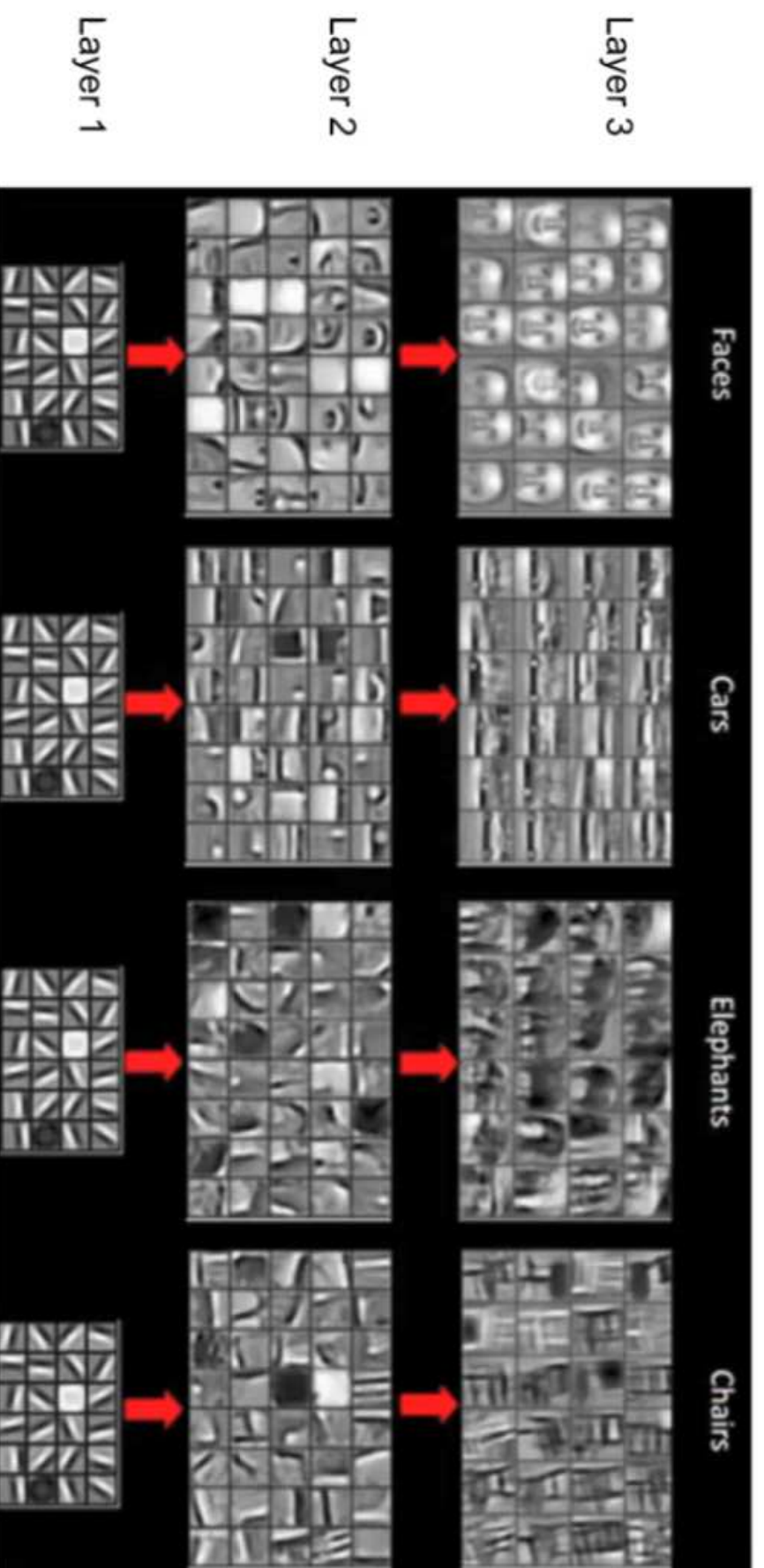
$$\delta^{(l)} = (\delta^{(l+1)} W^{(l+1)T}) \odot Z^{(l)'}, \text{ for } l = 1 \dots L - 1$$

Generic Weight Update

$$\nabla_{W^{(l)}} J = Z^{(l-1)T} \delta^{(l)}$$

$$\nabla_{b^{(l)}} J = \text{sum}(\delta^{(l)}, \text{axis} = 0) = \delta^{(l)T} \mathbf{1}_N$$

What has a CNN learned?



منابع

- An efficient and improved scheme for handwritten digit recognition based on convolutional neural network Saqib Ali¹ · Zeeshan Shaukat¹ · Muhammad Azeem¹ · Zareen Sakhawat¹ · Tariq Mahmood² · Khalil ur Rehman¹ @ Springer Nature Switzerland AG 2019
- Back Propagation Neural Networks Article in Substance Use & Misuse · February 1998