

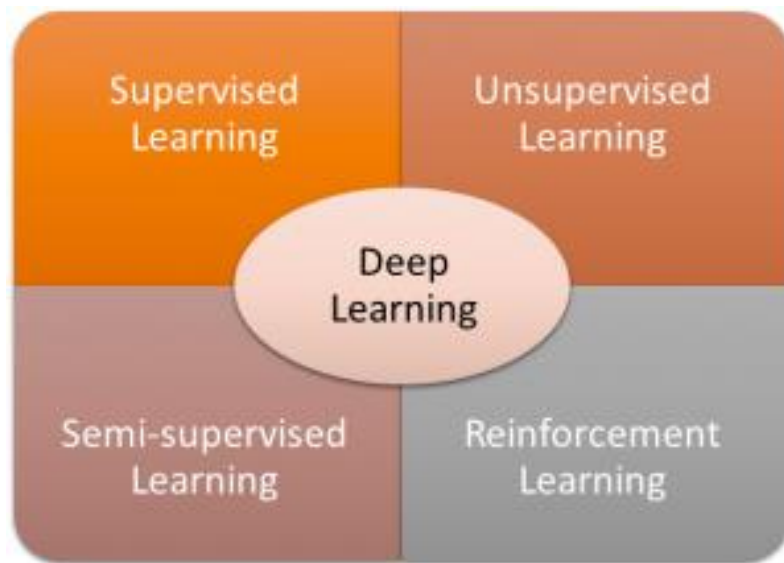
Anatomy and Arithmetic of Convolutional Neural Networks Part I

Odessa, 2018

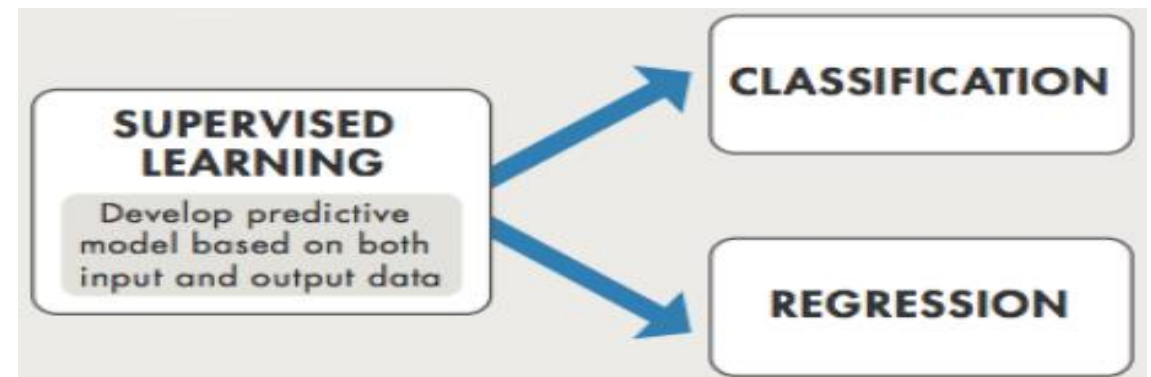
Styles of Learning

Supervised	Unsupervised	Semi-Supervised	Reinforcement
<ul style="list-style-type: none"> Data has known labels 	<ul style="list-style-type: none"> No labels Focus on finding patterns and gaining insight from the data 	<ul style="list-style-type: none"> Labels known for a subset of data A blend of supervised and unsupervised 	<ul style="list-style-type: none"> Focus on making decisions based on previous experience

(text and layout were taken from blogs.sas.com)



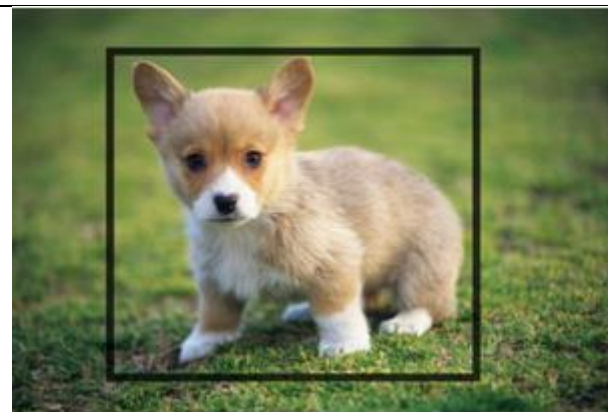
(leonardoaraujosantos.gitbooks.io)



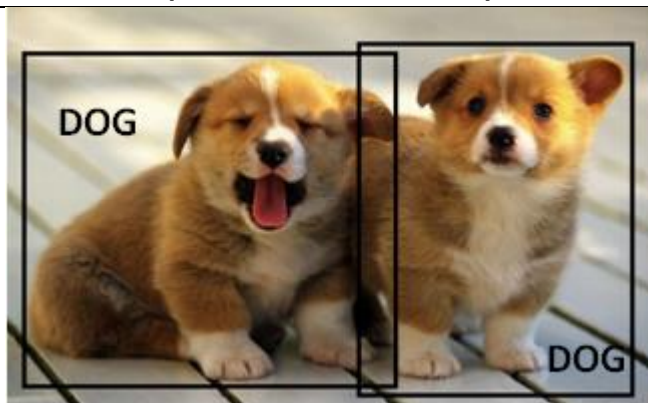
(vitalflux.com)



(classification)



(classification + localization)



(object detection)






(object segmentation)

Semantic Segmentation = Pixel-wise Classification



Loss Functions – 1

Bounding box	Facial landmarks	Pose estimation
 <p>(coursera.org)</p>	 <p>(stackoverflow.com)</p>	 <p>(cs231n.stanford.edu)</p>
(c_x, c_y, w, h)	$(x_1, y_1, \dots, x_M, y_M)$	$(x_1, y_1, \dots, x_M, y_M)$

$$\frac{1}{m} \sum_{i=1}^m \left(\sum_{j=1}^n L(s_j, \hat{s}_j) \right) \rightarrow \min$$

$$L(s_j, \hat{s}_j) = \frac{1}{2} (s_j - \hat{s}_j)^2$$

Loss Functions – 2

loss function:

$$\begin{aligned}
 & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\
 & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[\left(\sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left(\sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\
 & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\
 & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2 \\
 & + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \quad (3)
 \end{aligned}$$

where $\mathbb{1}_i^{\text{obj}}$ denotes if object appears in cell i and $\mathbb{1}_{ij}^{\text{obj}}$ denotes that the j th bounding box predictor in cell i is “responsible” for that prediction.

Loss Functions – 3. Categorical cross-entropy

$$\hat{y} = s \quad \text{vs} \quad y = j$$

$$\mathbf{p} = (p_1, \dots, p_K) \quad \text{vs} \quad \mathbf{y} = \left(0, \dots, \underbrace{1}_j, \dots, 0 \right)$$

$$P(\mathbf{y}) = p_j = p_1^0 \dots p_j^1 \dots p_K^0 = p_1^{y_1} \dots p_K^{y_K}$$



one-hot encoding

$$p_1^{y_1} \dots p_K^{y_K} \rightarrow \max$$

$$y_1 \ln p_1 + \dots + y_K \ln p_K \rightarrow \max$$

$$0 < p_k < 1 \implies \ln p_k < 0 \implies -\ln p_k > 0$$

$$-y_1 \ln p_1 - \dots - y_K \ln p_K \rightarrow \min$$

$$\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K -y_k^{(i)} \ln p_k^{(i)} \rightarrow \min$$

Three Main Types of Classification

- **Binary**: two mutually exclusive classes (object belongs to one class)

$$\begin{aligned} y^{(i)} &= 1 \\ p^{(i)} &= 0.7 \end{aligned}$$

$$-y^{(i)} \ln p^{(i)} - (1 - y^{(i)}) \ln(1 - p^{(i)})$$

- **Multi-class**: any number of mutually exclusive (object belongs to one class)

$$\begin{aligned} y^{(i)} &= [0 \quad 1 \quad 0] \\ p^{(i)} &= [0.3 \quad 0.6 \quad 0.1] \end{aligned}$$

$$\sum_{k=1}^K -y_k^{(i)} \ln p_k^{(i)}$$

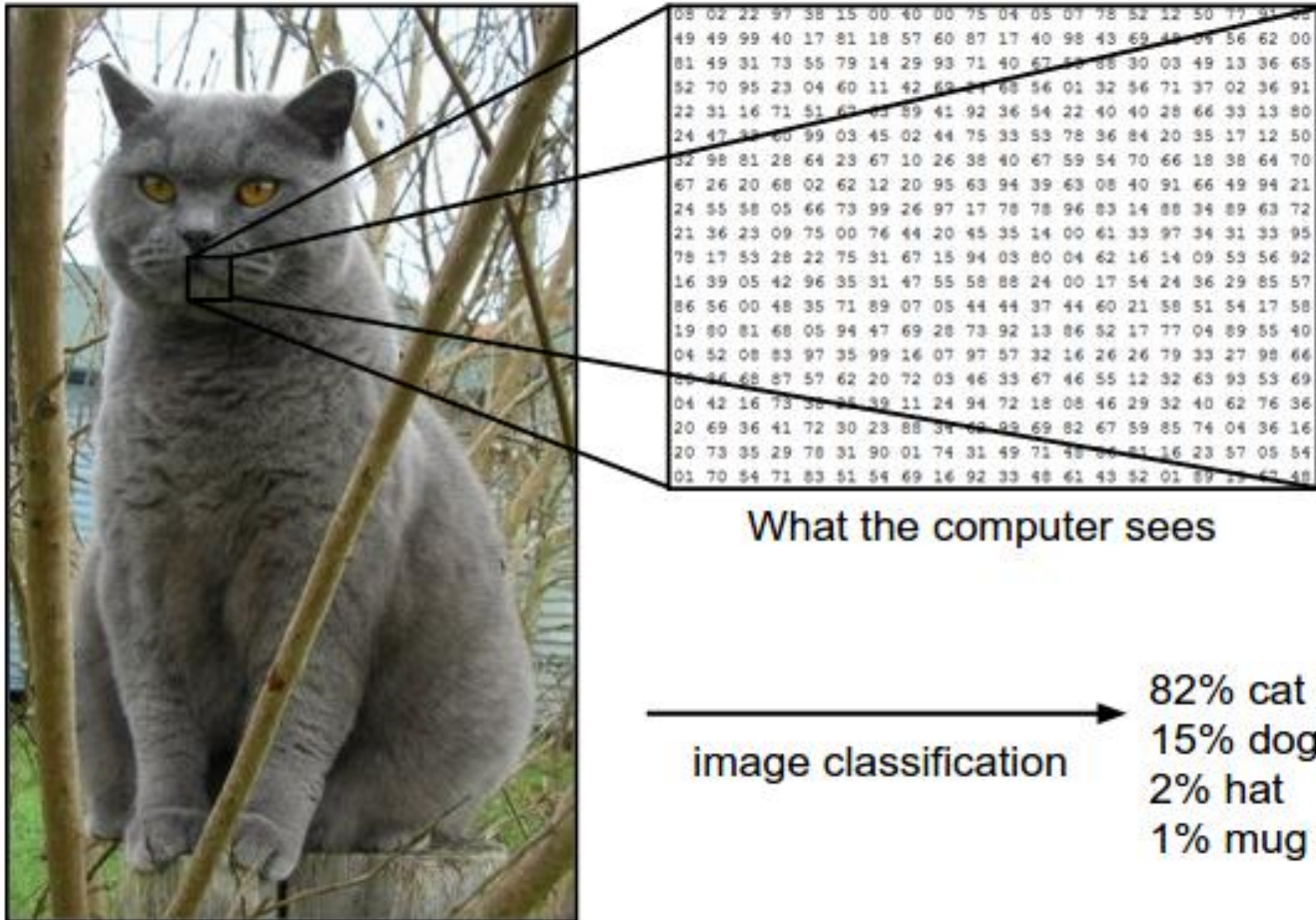
- **Multi-label**: any number of independent classes (object can belong to many classes)

$$\begin{aligned} y^{(i)} &= [1 \quad 1 \quad ? \quad 0] \\ p^{(i)} &= [0.5 \quad 0.7 \quad 0.4 \quad 0.2] \end{aligned}$$

$$\sum_{\substack{k=1 \\ y_k^{(i)} \neq ?}}^K \left[-y_k^{(i)} \ln p_k^{(i)} - (1 - y_k^{(i)}) \ln (1 - p_k^{(i)}) \right]$$



Why images are hard? – 1



What the computer sees

image classification

82% cat
15% dog
2% hat
1% mug

Why images are hard? – 2

Viewpoint variation



Scale variation



Deformation



Occlusion



Illumination conditions



Background clutter

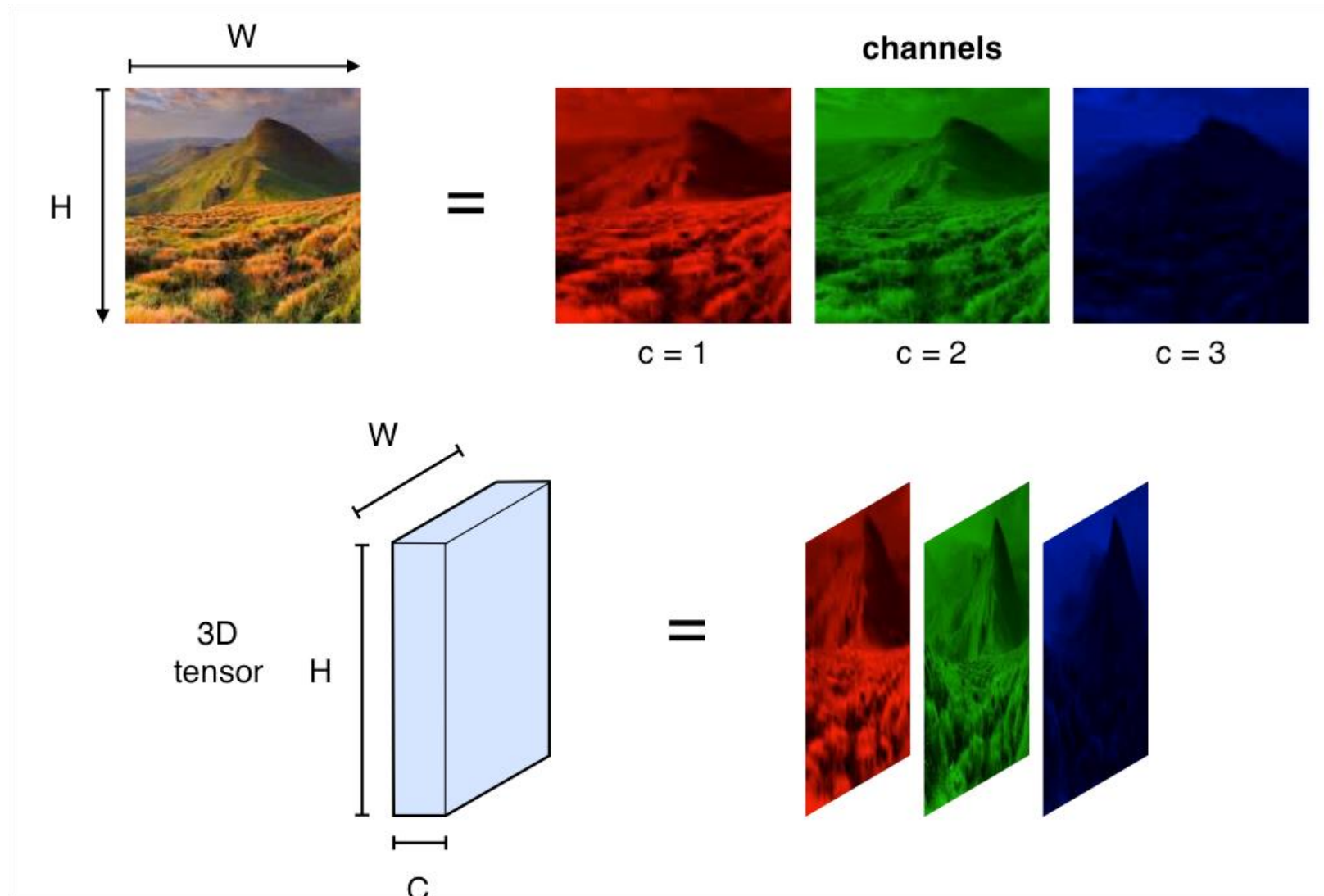


Intra-class variation



(cs231n.github.io)

What is an image – 2



Meet the Tensors

't'
'e'
'n'
's'
'o'
'r'

***Tensor of
dimension[1]***

3	1	4	1
5	9	2	6
5	3	5	8
9	7	9	3
2	3	8	4
6	2	6	4

***Tensor of
dimensions[2]***

2	1	2	1
2	4	9	4
2	5	6	2
7	7	3	2

***Tensor of
dimensions[3]***

(Edureka)

Common Operations on Tensors

- Extending scalar binary operations:

$$(A \star B)_\alpha = A_\alpha \star B_\alpha$$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} + \begin{bmatrix} 2 & 0 \\ 6 & 6 \end{bmatrix} = \begin{bmatrix} 3 & 2 \\ 9 & 10 \end{bmatrix}$$

- Applying scalar functions:

$$g: \mathbb{R} \rightarrow \mathbb{R}, \quad (g(A))_\alpha = g(A_\alpha)$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}, \quad \sigma\left(\begin{bmatrix} 0 & 1 \\ 2 & 3 \end{bmatrix}\right) = \begin{bmatrix} 0.5000 & 0.7311 \\ 0.8808 & 0.9526 \end{bmatrix}$$

- Flatten (reshaping to 1d-tensor):

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \rightarrow [1 \quad 2 \quad 3 \quad 4]$$

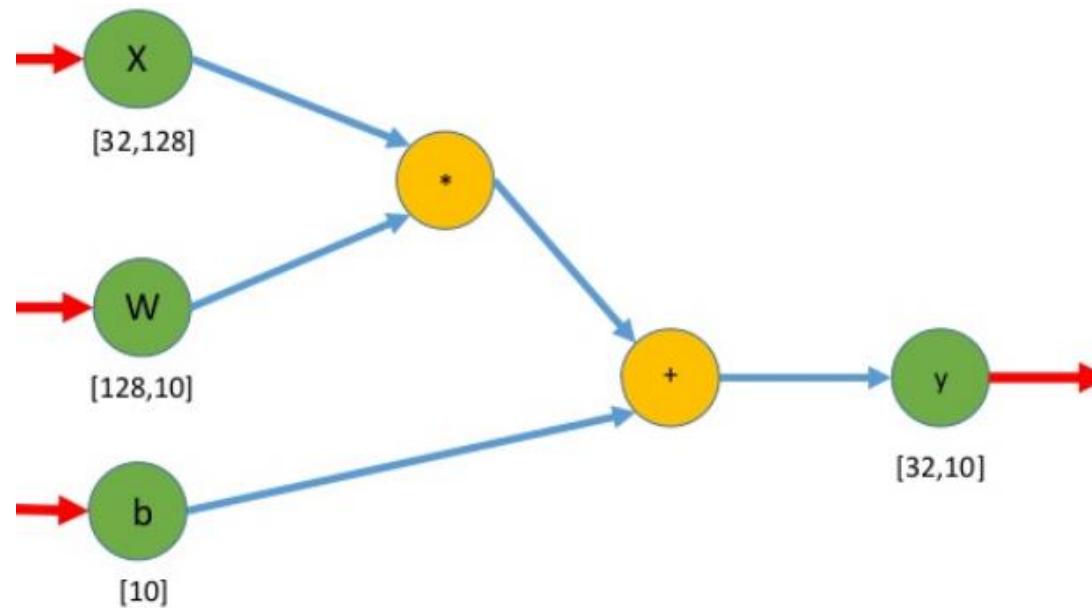
- Concatenation:

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \\ 10 & 11 & 12 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 & 7 & 8 & 9 \\ 4 & 5 & 6 & 10 & 11 & 12 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \\ 10 & 11 & 12 \end{bmatrix}$$

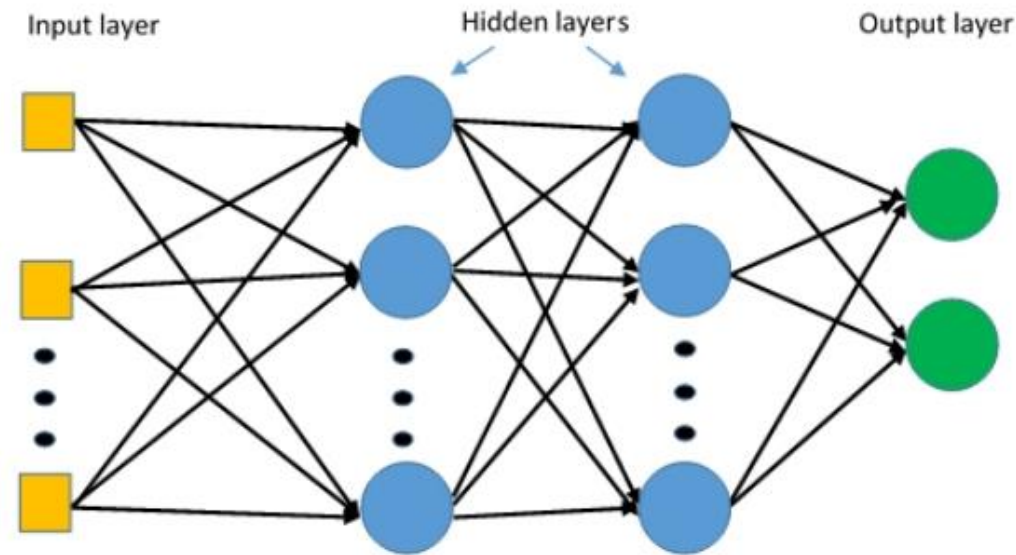
Modern Neural Network (MLP, CNN, RNN) == Graph of Tensors



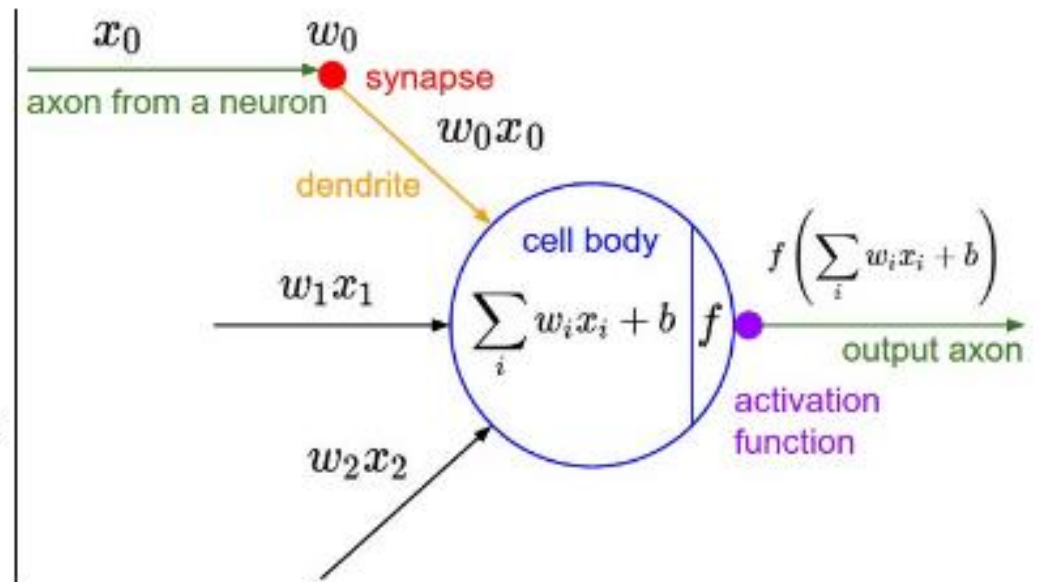
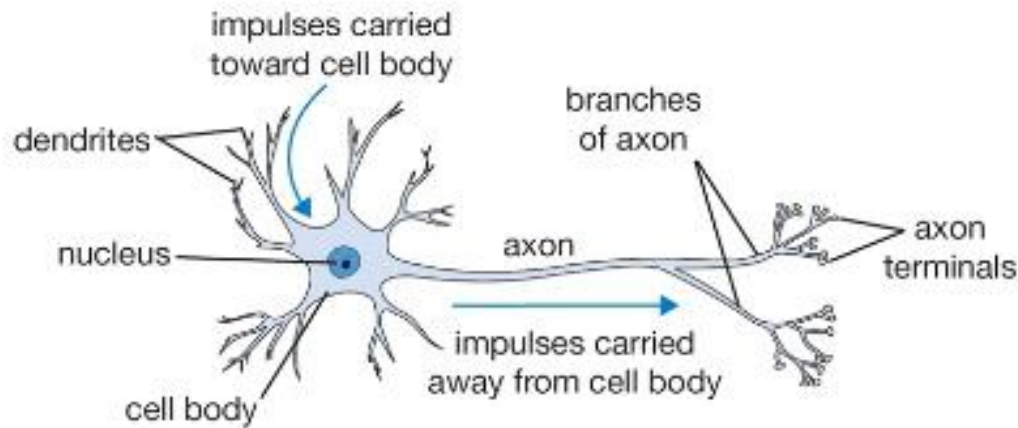
(Alessio Tonioni on slideshare.net)

$$Forward(\theta) = f(Parents(\theta))$$

$$Backward(\theta) = \frac{\partial E}{\partial \theta} = \sum_{c \in Children(\theta)} \left(\frac{\partial E}{\partial c} \cdot \frac{\partial c}{\partial \theta} \right)$$



(packtpub.com)

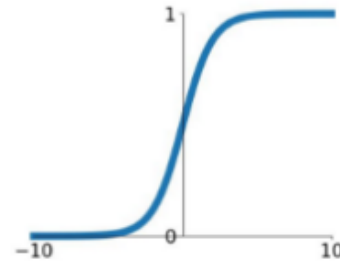


(cs231n.github.io)

Activation Functions

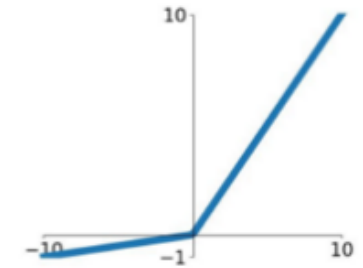
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



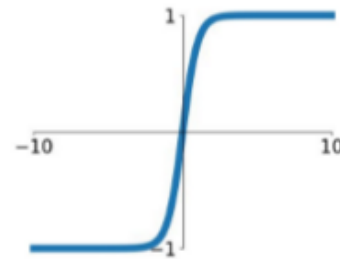
Leaky ReLU

$$\max(0.1x, x)$$



tanh

$$\tanh(x)$$

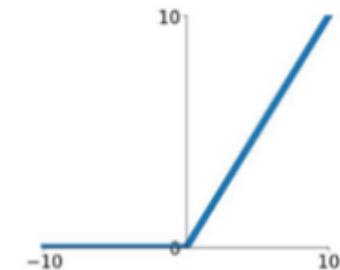


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

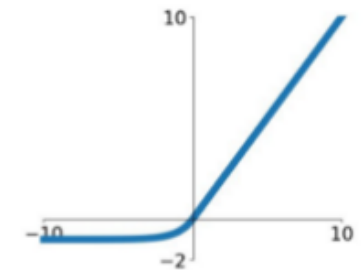
ReLU

$$\max(0, x)$$



ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



(Shruti Jadon on Medium)

arXiv.org > cs > arXiv:1710.05941

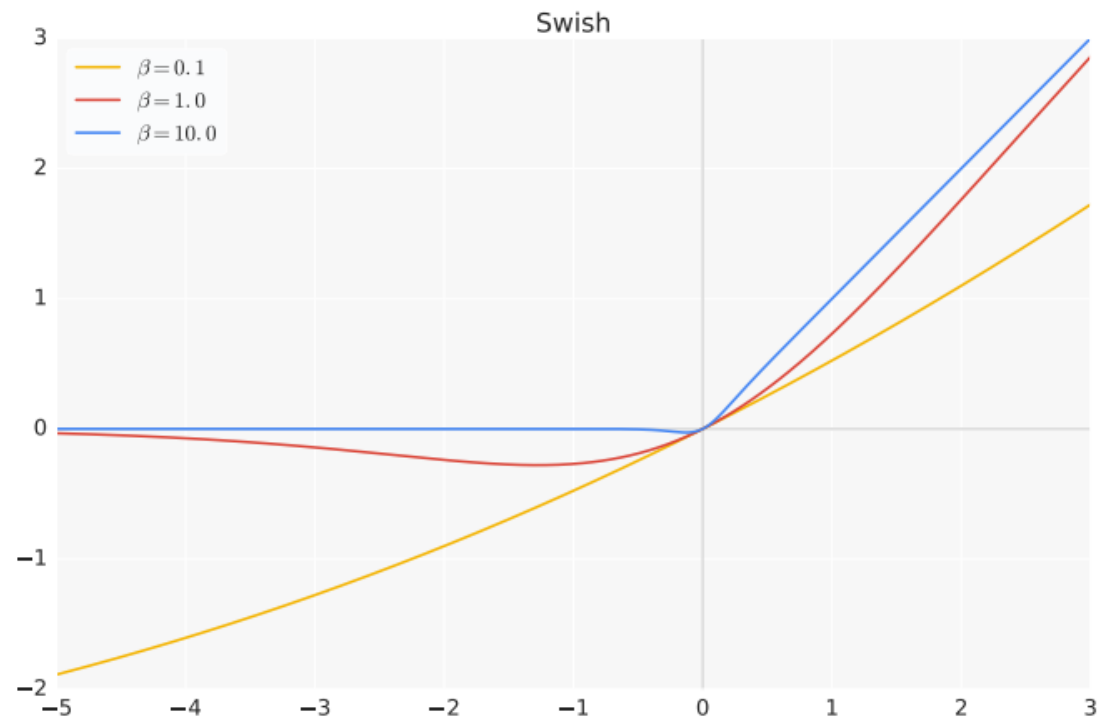
Computer Science > Neural and Evolutionary Computing

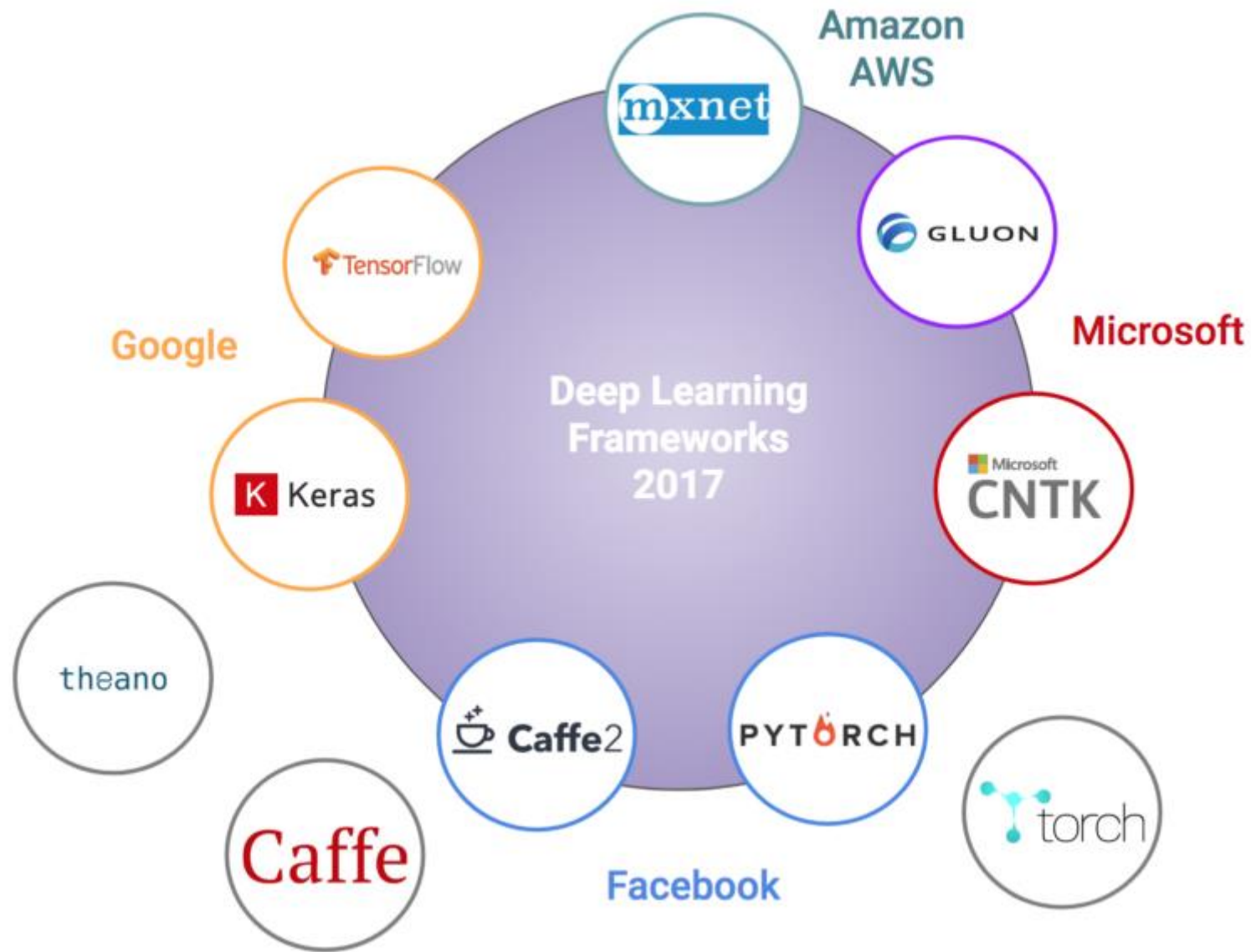
Searching for Activation Functions

Prajit Ramachandran, Barret Zoph, Quoc V. Le

(Submitted on 16 Oct 2017 (v1), last revised 27 Oct 2017 (this version, v2))

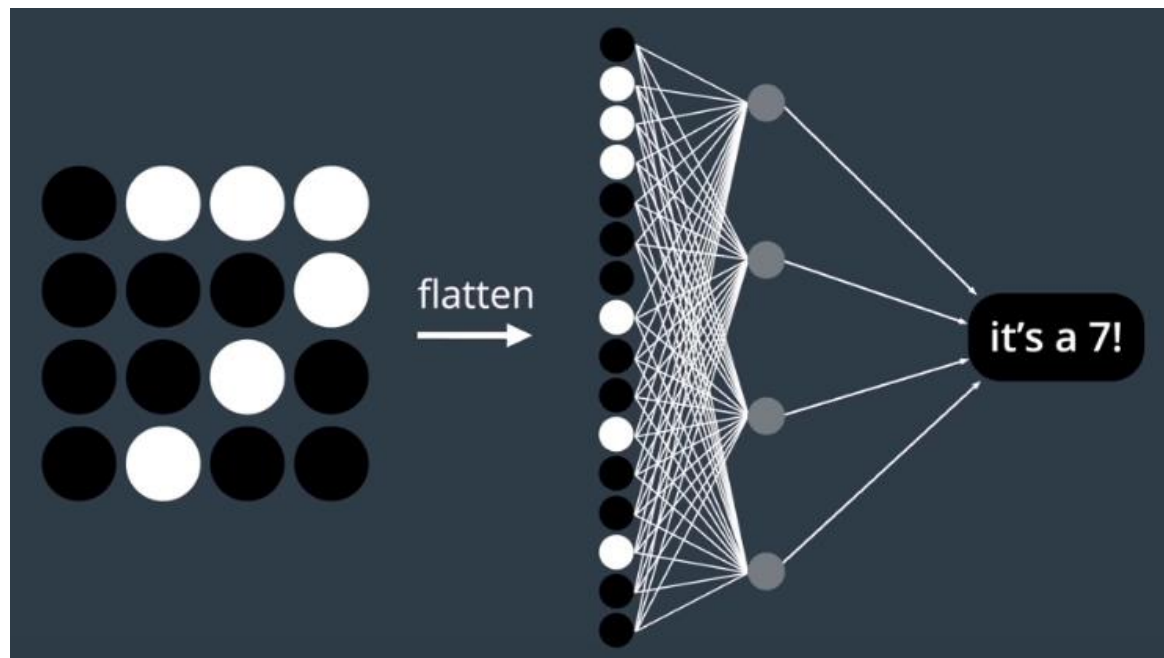
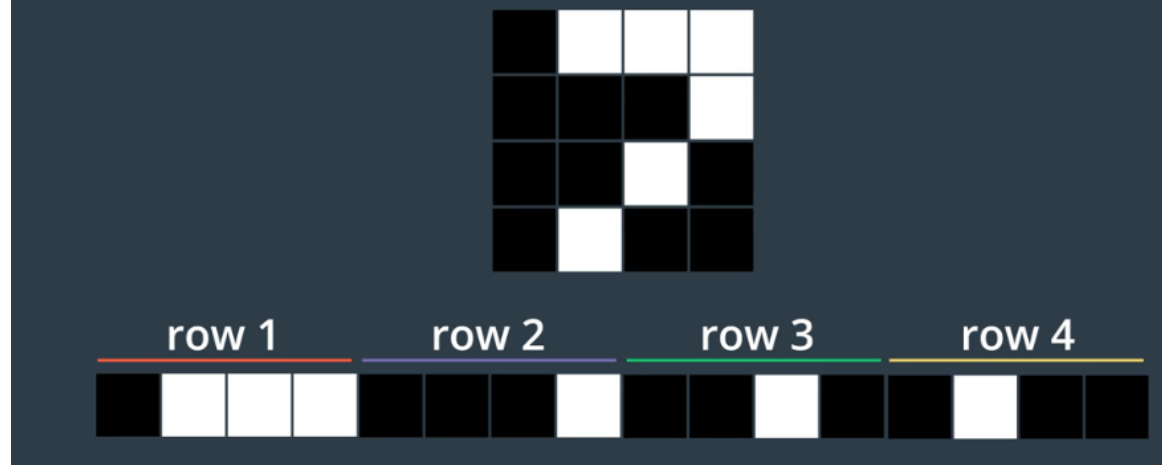
$$\text{Swish}(x) = x \cdot \sigma(\beta x) = \frac{x}{1 + e^{-\beta x}}$$





(towardsdatascience.com)

Flattening



(pictures with gray background were taken from Udacity)

Softmax – 1

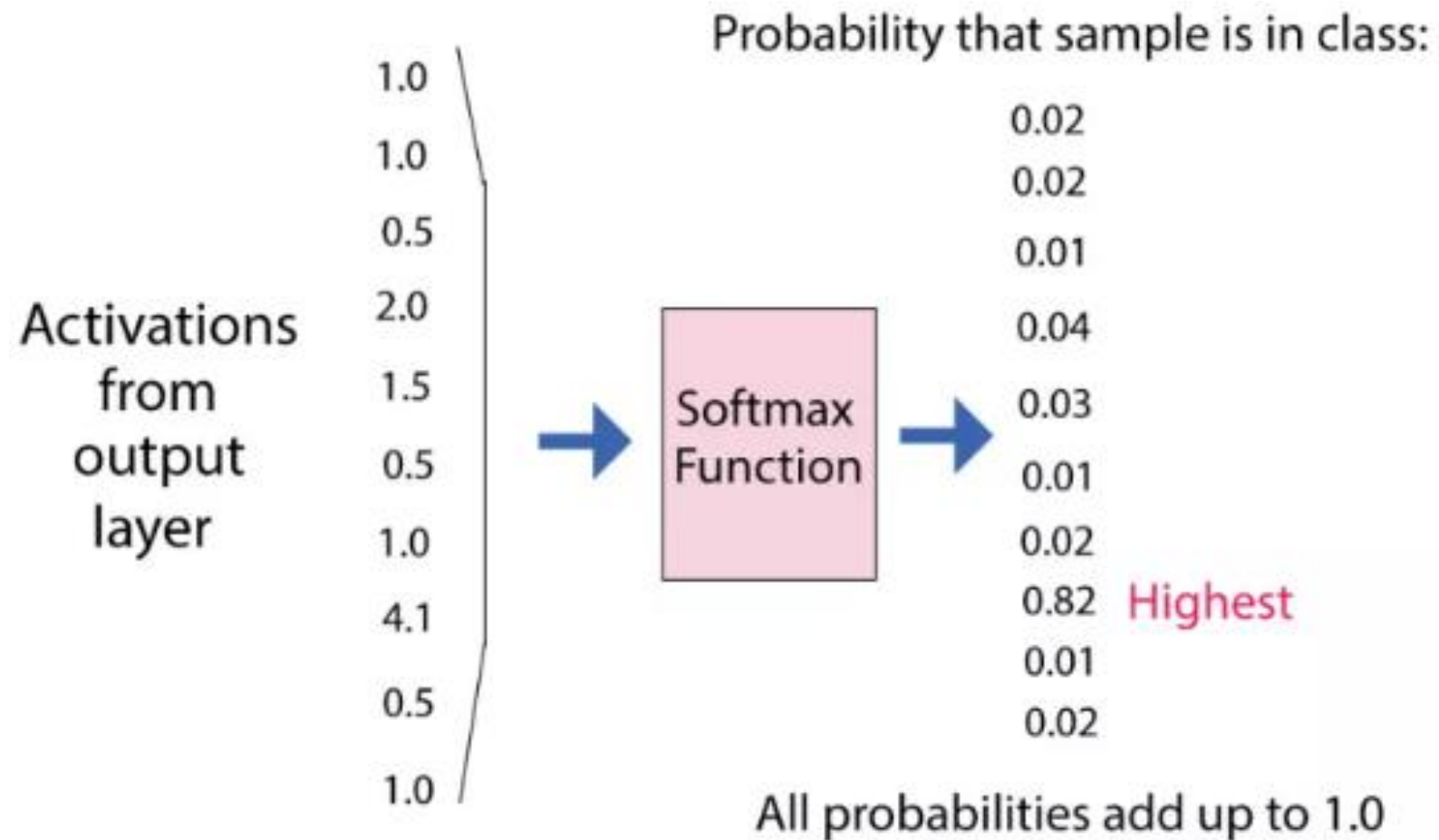
$$\begin{bmatrix} z_1 \\ \vdots \\ z_K \end{bmatrix} \mapsto \begin{bmatrix} \frac{z_1}{\sum e^{z_i}} \\ \vdots \\ \frac{z_K}{\sum z_i} \end{bmatrix}, \quad \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix} \mapsto \begin{bmatrix} 0.5 \\ -0.5 \\ 1.0 \end{bmatrix}$$

$$\begin{bmatrix} z_1 \\ \vdots \\ z_K \end{bmatrix} \mapsto \begin{bmatrix} e^{z_1} \\ \vdots \\ e^{z_K} \end{bmatrix} \mapsto \begin{bmatrix} \frac{e^{z_1}}{\sum e^{z_i}} \\ \vdots \\ \frac{e^{z_K}}{\sum e^{z_i}} \end{bmatrix} = \begin{bmatrix} p_1 \\ \vdots \\ p_K \end{bmatrix}$$

$$0 < p_i < 1, \quad \sum_{i=1}^K p_i = 1$$

$$z_i < z_j \quad \Rightarrow \quad p_i < p_j$$

Softmax – 2



(<http://principlesofdeeplearning.com>)

```

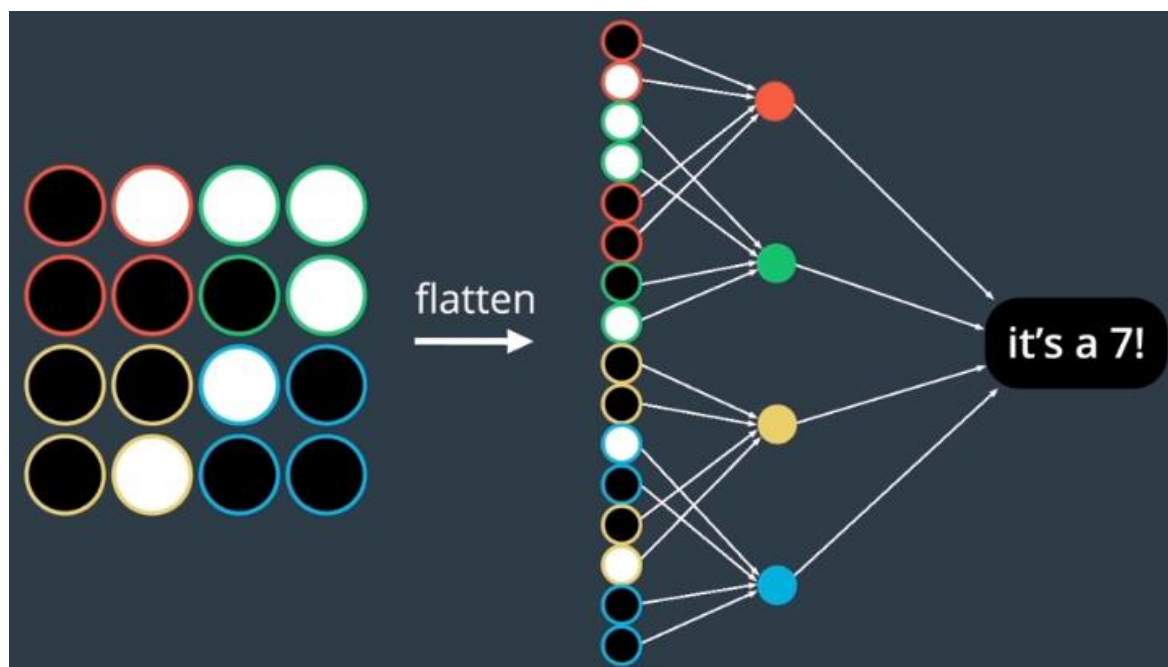
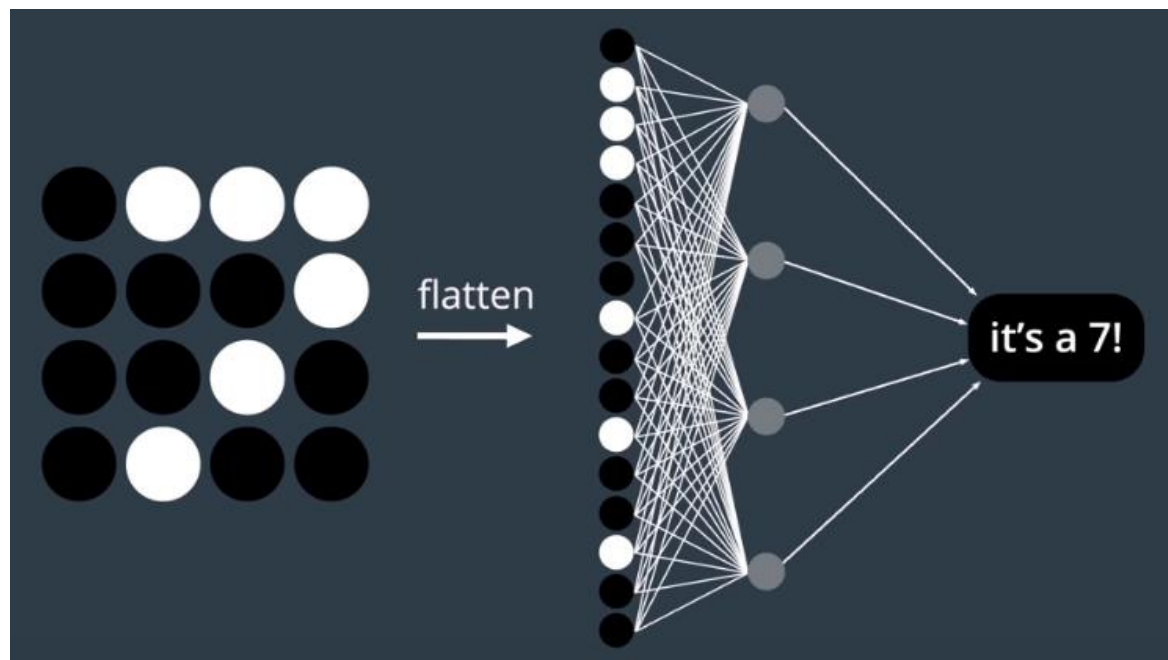
model = Sequential()
model.add(Dense(512, activation='relu', input_shape=(784,)))
model.add(Dropout(0.2))
model.add(Dense(512, activation='relu'))
model.add(Dropout(0.2))
model.add(Dense(num_classes, activation='softmax'))
model.summary()

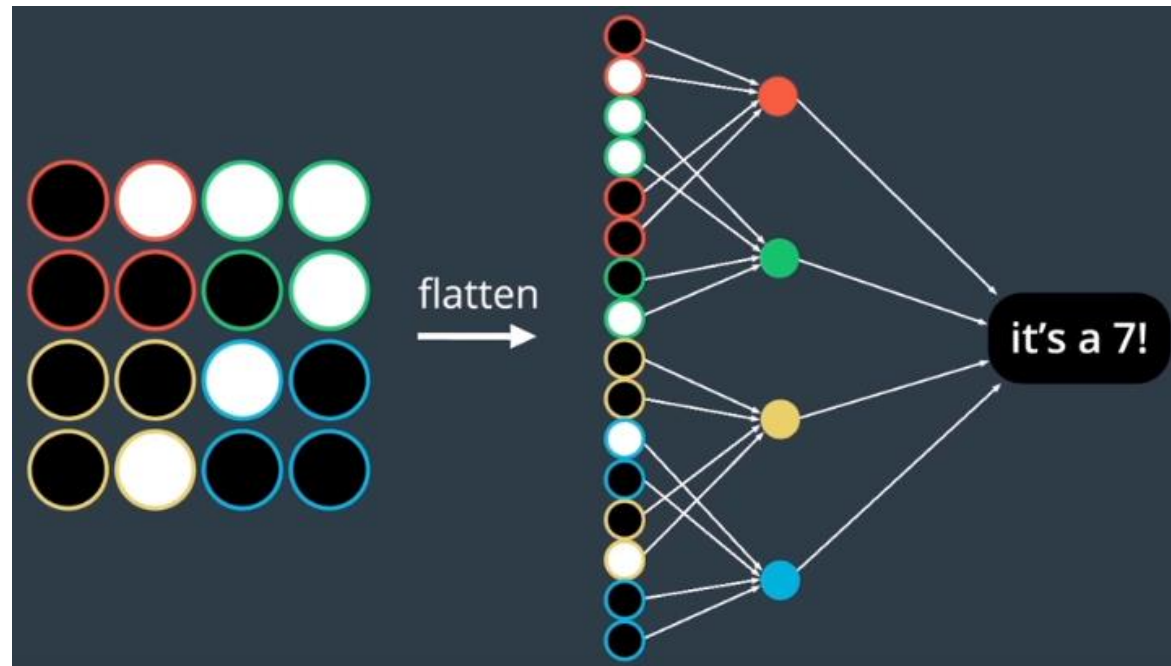
```

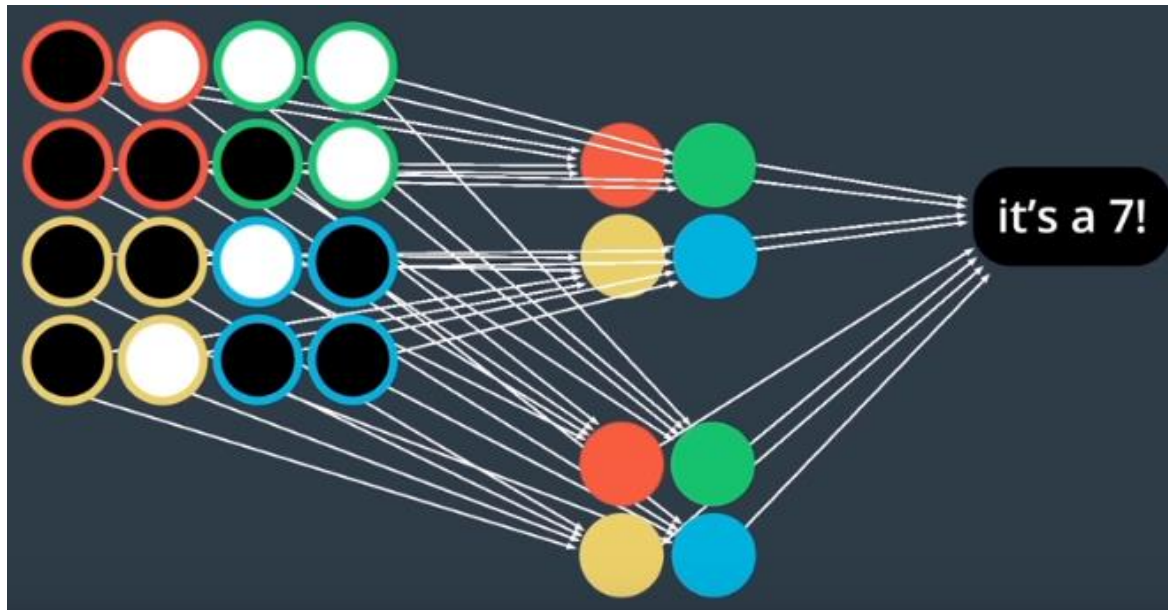
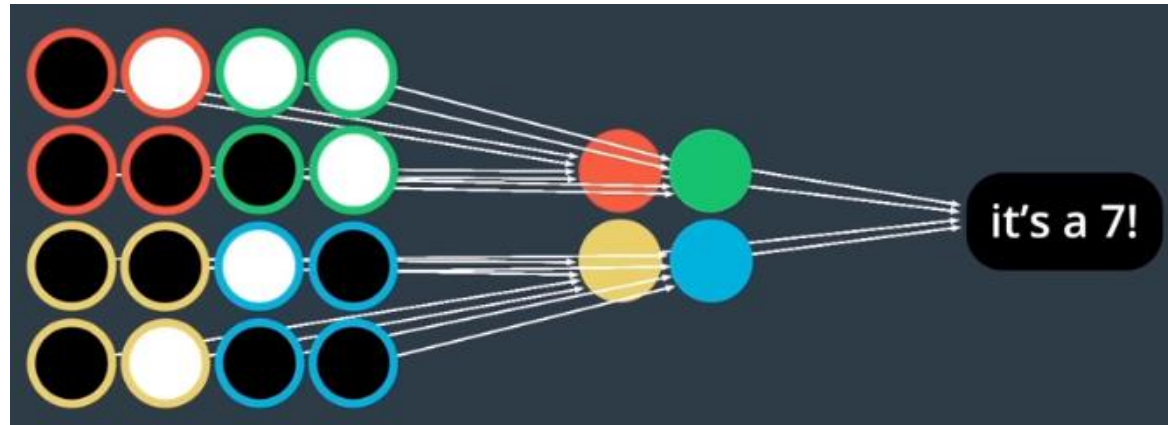
Layer (type)	Output Shape	Param #
=====	=====	=====
dense_1 (Dense)	(None, 512)	401920
dropout_1 (Dropout)	(None, 512)	0
dense_2 (Dense)	(None, 512)	262656
dropout_2 (Dropout)	(None, 512)	0
dense_3 (Dense)	(None, 10)	5130
=====	=====	=====
Total params: 669,706		
Trainable params: 669,706		
Non-trainable params: 0		

Layer (type)	Output Shape	Param #
=====	=====	=====
dense_1 (Dense)	(None, 512)	401920
dropout_1 (Dropout)	(None, 512)	0
dense_2 (Dense)	(None, 512)	262656
dropout_2 (Dropout)	(None, 512)	0
dense_3 (Dense)	(None, 10)	5130
=====	=====	=====
Total params: 669,706		
Trainable params: 669,706		
Non-trainable params: 0		

Input	$28 \cdot 28 = 784$
Hidden 1 (dense_1)	$(784 + 1) \cdot 512 = 401\,920$
Hidden 2 (dense_2)	$(512 + 1) \cdot 512 = 262\,656$
Output (dense_3)	$(512 + 1) \cdot 10 = 5\,130$







Feature is a specific 2D structure in the image such as a blob, corner or an edge than can be described in a local neighborhood by its appearance information.

Descriptor is a vector that contains local appearance information.

(Pablo F. Alcantarilla)



(<http://cis.eecs.qmul.ac.uk>)



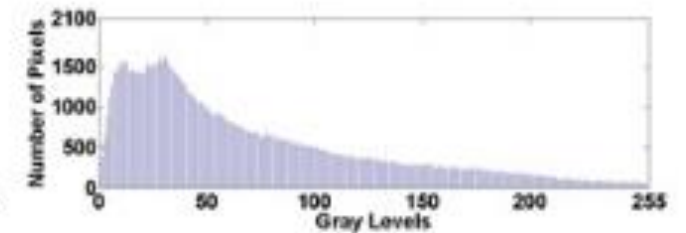
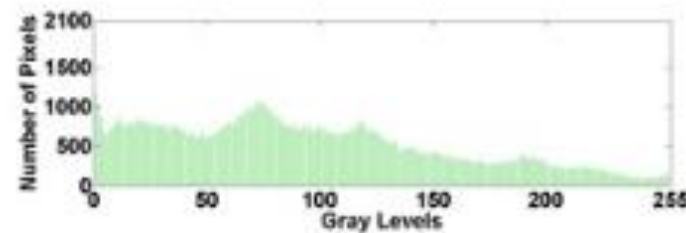
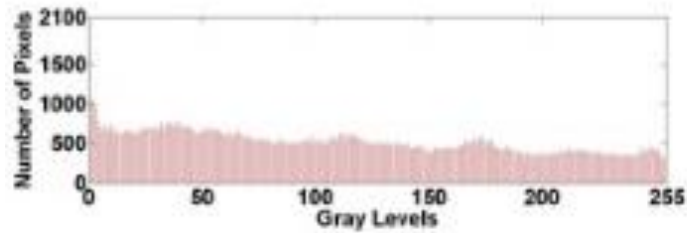
Red



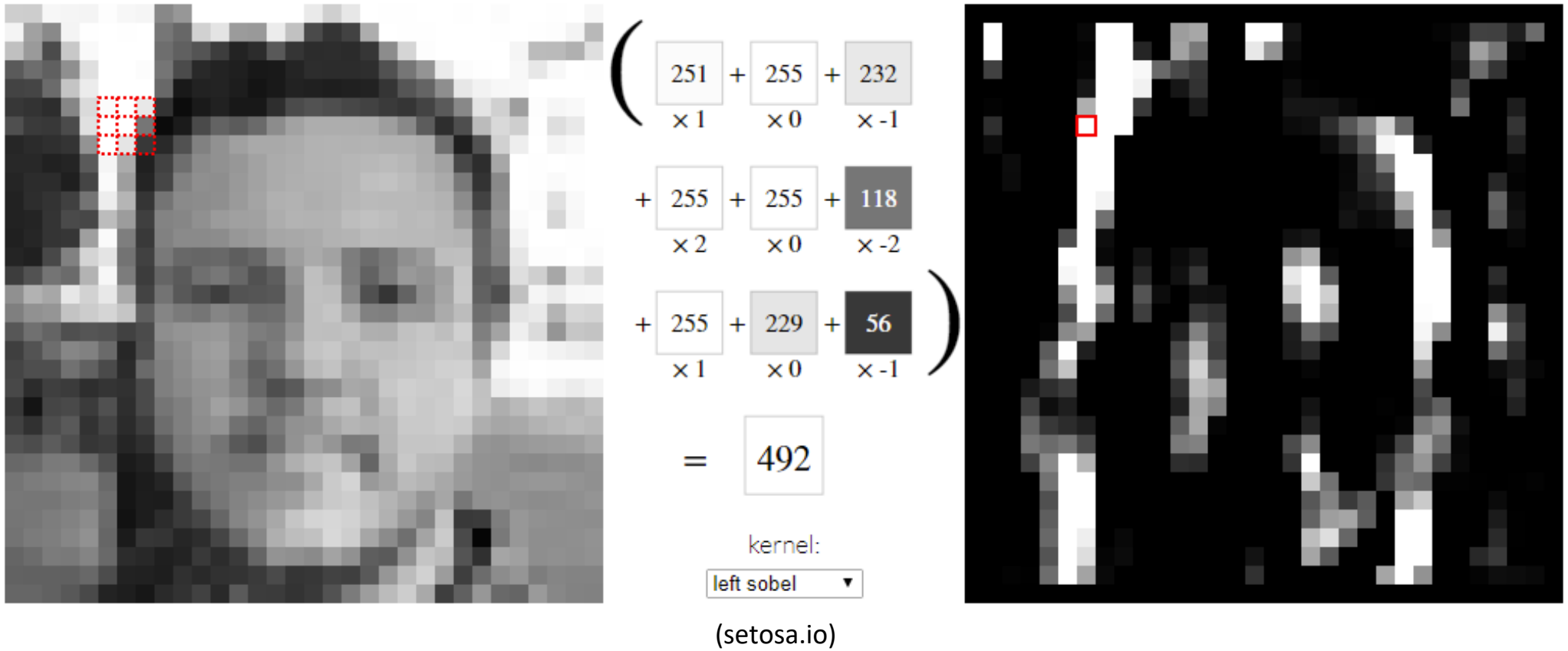
Green

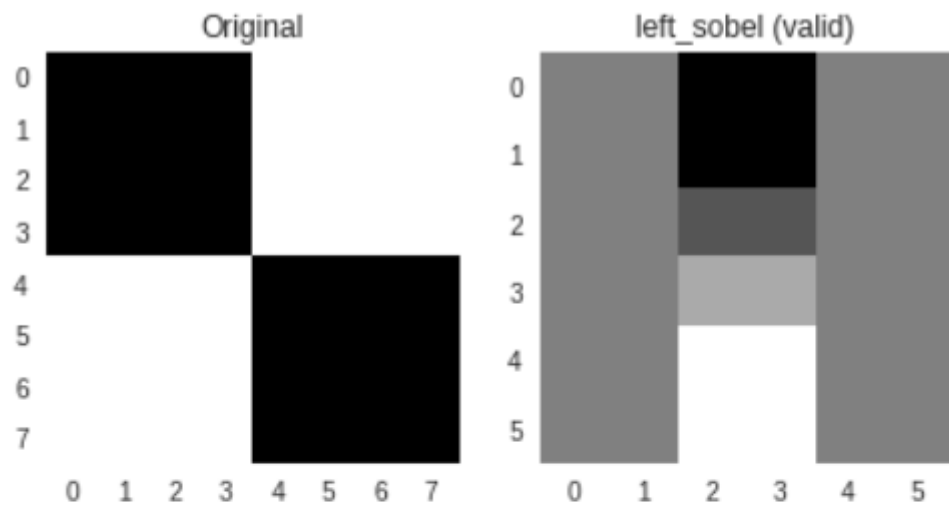


Blue

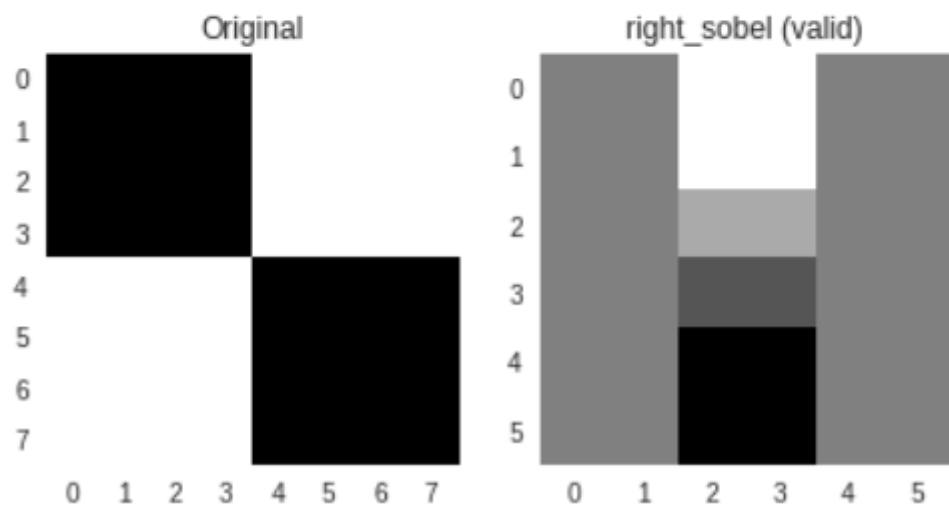


(allaboutcircuits.com)

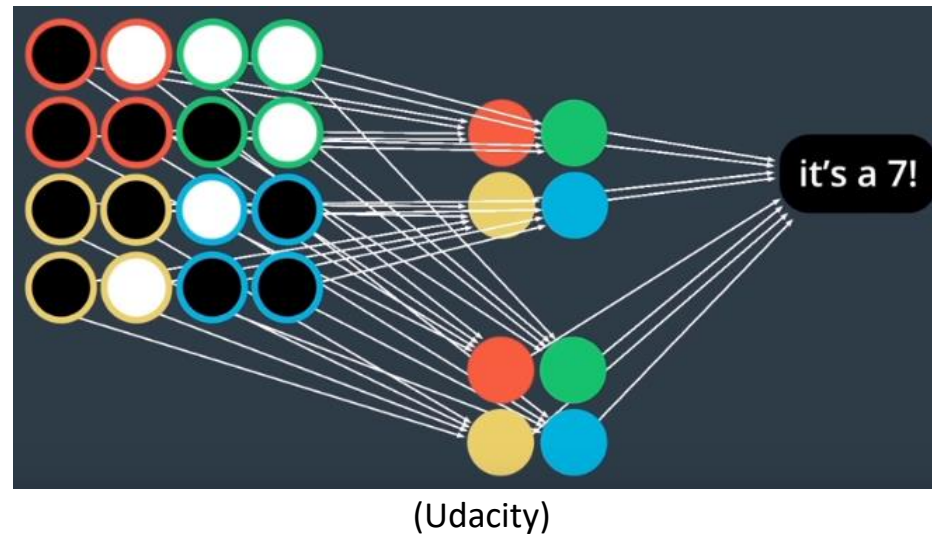




$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$



$$\begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$



- **Parameter sharing.** A feature detector (such as vertical edge detector) that's useful in one part on the image is *probably* useful in another part of the image.
- **Sparsity of connections.** In each layer each output value depends on a small number of inputs.

(Coursera)