

Introduction to Deep Learning

Shangsong Liang
Sun Yat-sen University

DL is providing breakthrough results in speech recognition and image classification ...

From this Hinton et al 2012 paper:

<http://static.googleusercontent.com/media/research.google.com/en//pubs/archive/38131.pdf>

modeling technique	#params [10 ⁶]	WER		task	hours of training data	DNN-HMM	GMM-HMM with same data	GMM-HMM with more data
		Hub5'00-SWB	RT03S-FSH					
GMM, 40 mix DT 309h SI	29.4	23.6	27.4	Switchboard (test set 1)	309	18.5	27.4	18.6 (2000 hrs)
				Switchboard (test set 2)	309	16.1	23.6	17.1 (2000 hrs)
NN 1 hidden-layer×4634 units	43.6	26.0	29.4	English Broadcast News	50	17.5	18.8	
+ 2×5 neighboring frames	45.1	22.4	25.7	Bing Voice Search	24	30.4	36.2	
DBN-DNN 7 hidden layers×2048 units	45.1	17.1	19.6	(Sentence error rates)				
+ updated state alignment	45.1	16.4	18.6	Google Voice Input	5,870	12.3		16.0 (>>5,870hrs)
+ sparsification	15.2 nz	16.1	18.5	Youtube	1,400	47.6	52.3	
GMM 72 mix DT 2000h SA	102.4	17.1	18.6					

go here: <http://yann.lecun.com/exdb/mnist/>

From here:

<http://people.idsia.ch/~juergen/cvpr2012.pdf>

Dataset	Best result of others [%]	MCDNN [%]	Relative improv. [%]
MNIST	0.39	0.23	41
NIST SD 19	see Table 4	see Table 4	30-80
HWDB1.0 on.	7.61	5.61	26
HWDB1.0 off.	10.01	6.5	35
CIFAR10	18.50	11.21	39
traffic signs	1.69	0.54	72
NORB	5.00	2.70	46

Deep Learning is providing breakthrough results in speech recognition, image classification, etc.



google inception network

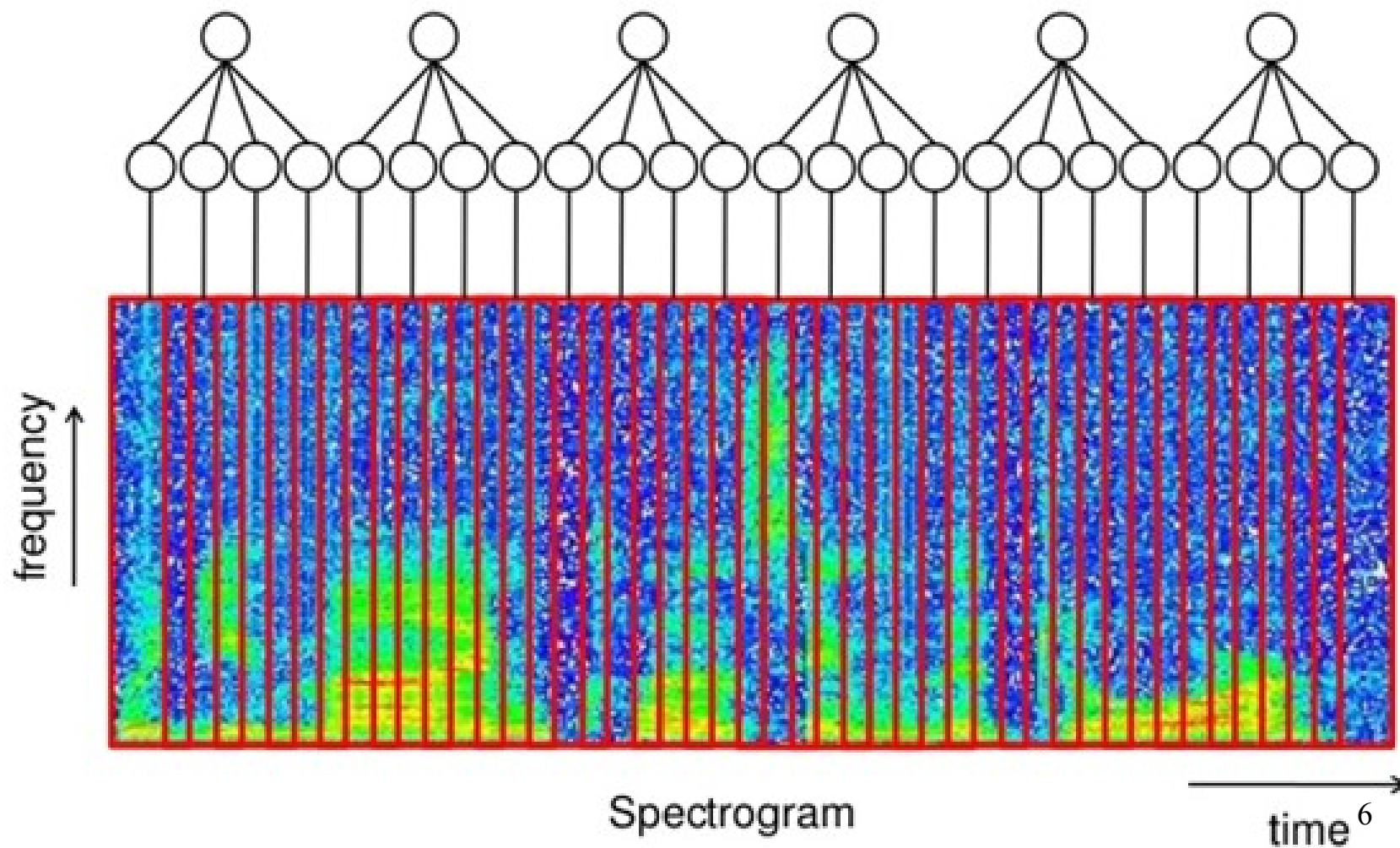
Examples from the test set (with the network's guesses)



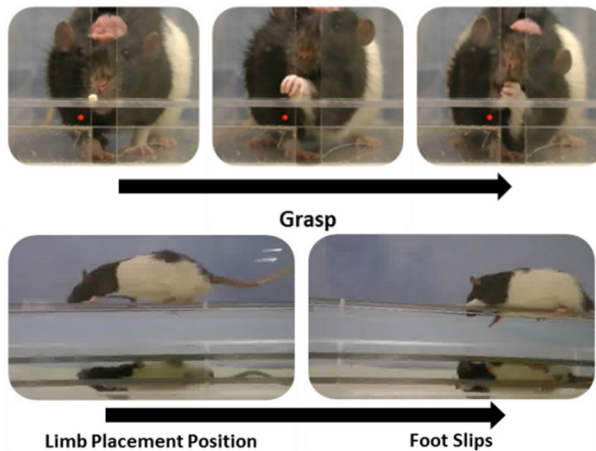
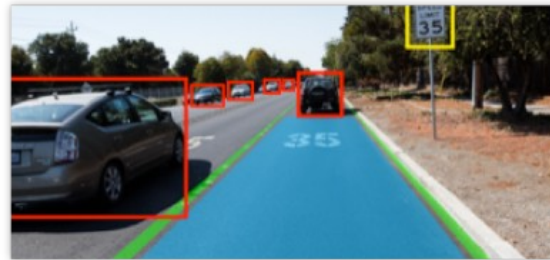
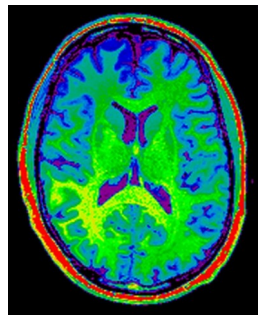
Video analyses and decision making



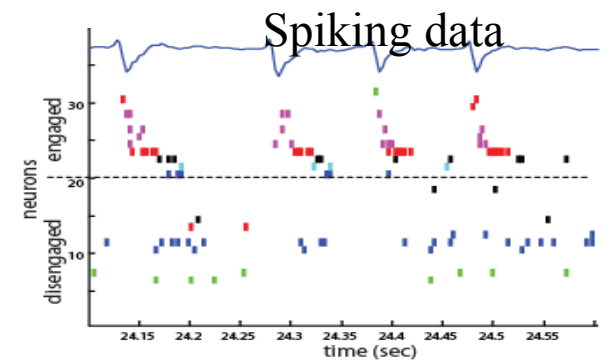
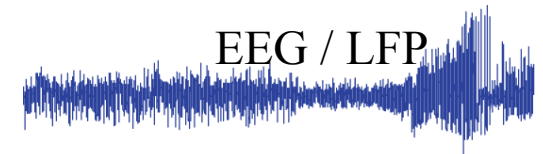
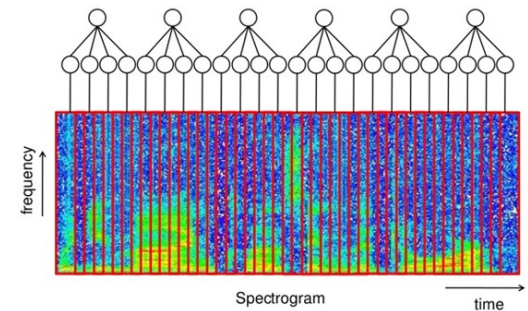
Speech recognition



Neuroscience data is similar to other types of data



Ryait et al. & Luczak



Images are Numbers



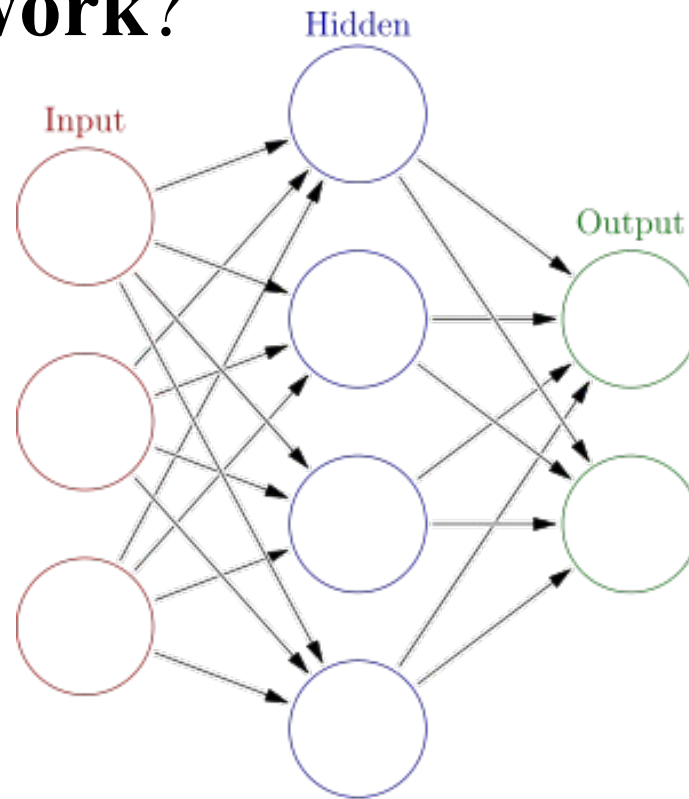
05	02	22	97	38	15	00	40	00	75	04	05	07	78	52	12	50	77	51	58
49	49	99	40	17	81	18	57	60	87	17	40	98	43	69	45	04	56	62	00
81	49	31	73	55	79	14	29	93	71	40	67	55	58	30	03	49	13	36	65
92	70	95	23	04	60	11	42	69	24	68	56	01	32	56	71	37	02	36	91
22	31	16	71	51	63	03	89	41	92	36	54	22	40	40	28	66	33	13	80
24	47	35	40	99	03	45	02	44	75	33	53	78	36	84	20	35	17	12	50
92	95	81	28	64	23	67	10	26	38	40	67	59	54	70	66	18	38	64	70
67	26	20	68	02	62	12	20	95	63	94	39	63	08	40	91	66	49	94	21
24	55	58	05	66	73	99	26	97	17	78	78	96	83	14	88	34	89	63	72
21	36	23	09	75	00	76	44	20	45	35	14	00	61	33	97	34	31	33	95
78	17	53	28	22	75	31	67	15	94	03	80	04	62	16	14	09	53	56	92
16	39	05	42	96	35	31	47	55	58	88	24	00	17	54	24	36	29	85	57
86	56	00	48	35	71	89	07	05	44	44	37	44	60	21	58	51	54	17	58
19	80	81	68	05	94	47	69	28	73	92	13	86	52	17	77	04	89	55	40
04	52	08	83	97	35	99	16	07	97	57	32	16	26	26	79	33	27	98	66
05	86	68	87	57	62	20	72	03	46	33	67	46	55	12	32	63	93	53	69
04	42	16	73	35	25	39	11	24	94	72	18	08	46	29	32	40	62	76	36
20	69	36	41	72	30	23	88	34	88	85	69	82	67	59	85	74	04	36	16
20	73	35	29	78	31	90	01	74	31	49	71	19	55	81	16	23	57	05	54
01	70	54	71	83	51	54	69	16	92	33	48	61	43	52	01	89	29	67	88

What the computer sees

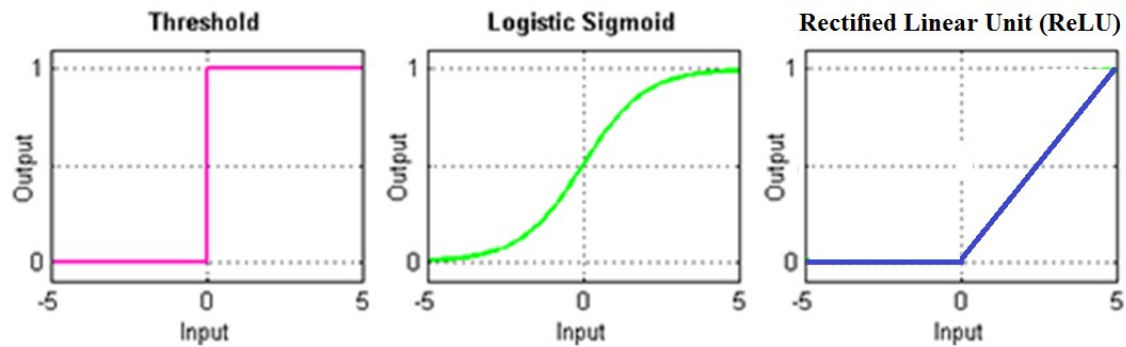
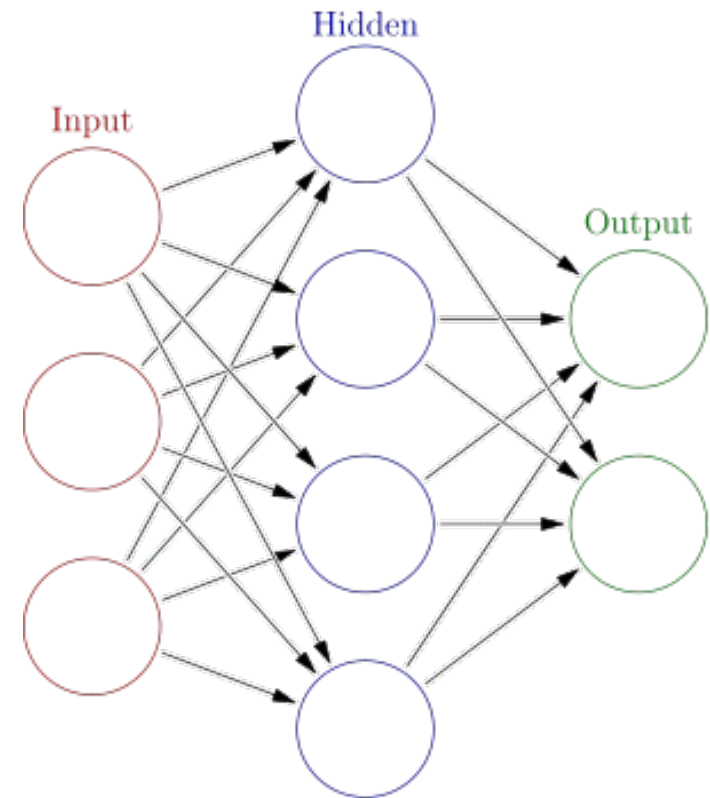
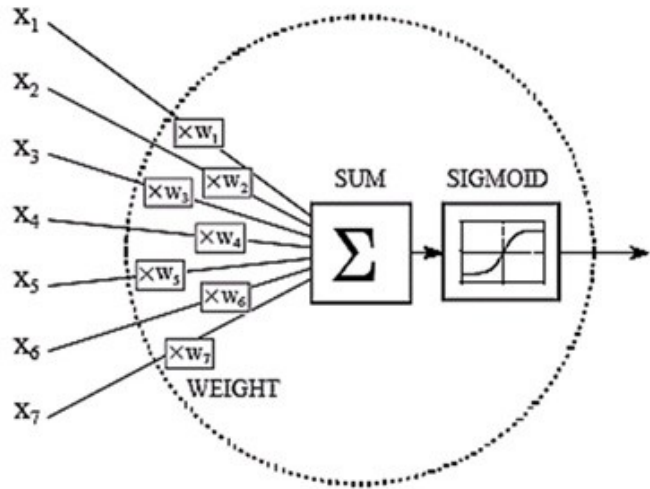
image classification →

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

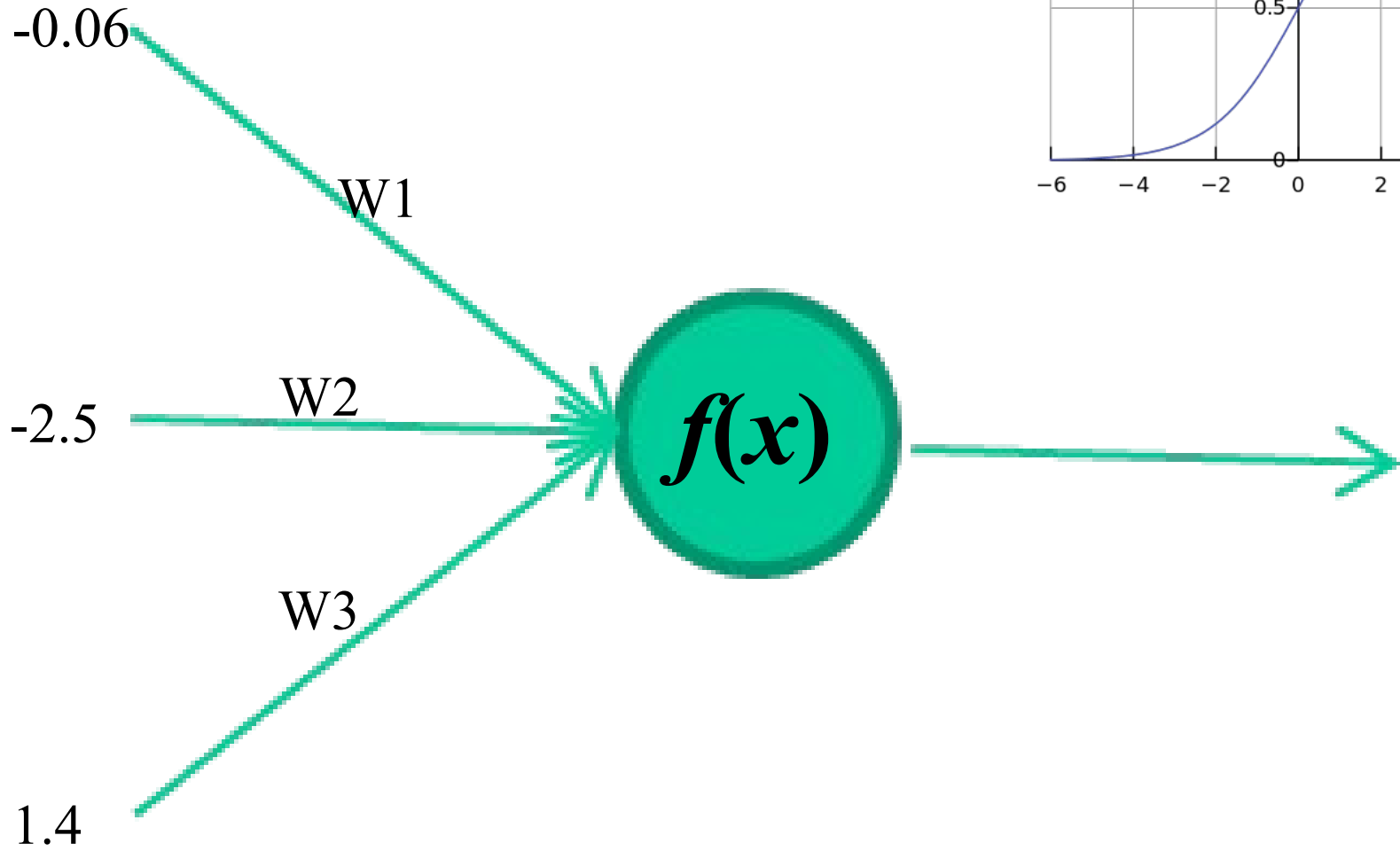
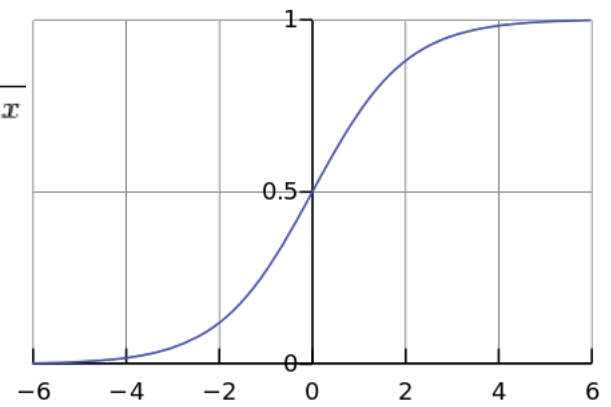
What is Artificial Neuronal Network?



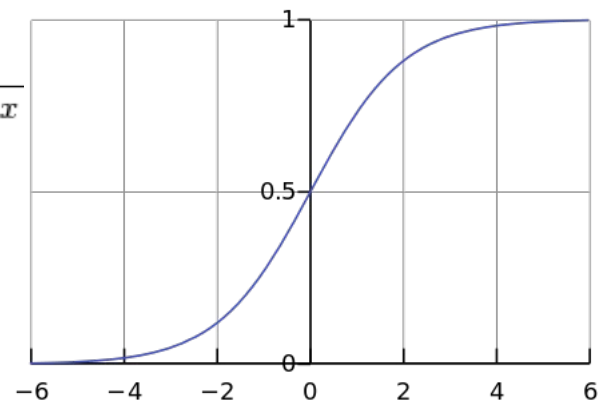
Artificial neural



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



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



-0.06

2.7

-2.5

-8.6

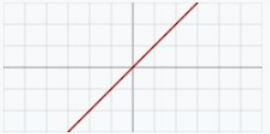


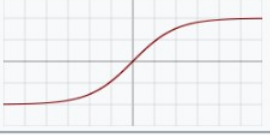

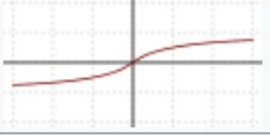
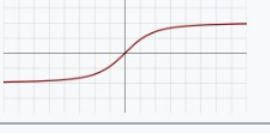
0.002

1.4

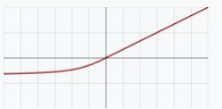



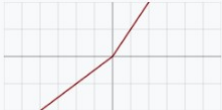


$f(x)$

$$x = -0.06 \times 2.7 + 2.5 \times 8.6 + 1.4 \times 0.002 = 21.34$$

Activation function

Identity		$f(x) = x$	$f'(x) = 1$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a. Sigmoid or Soft step)		$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$ ^[1]	$f'(x) = f(x)(1 - f(x))$
TanH		$f(x) = \tanh(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})}$	$f'(x) = 1 - f(x)^2$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$
ArSinH		$f(x) = \sinh^{-1}(x) = \ln(x + \sqrt{x^2 + 1})$	$f'(x) = \frac{1}{\sqrt{x^2 + 1}}$
ElliotSig ^{[9][10][11]} Softsign ^{[12][13]}		$f(x) = \frac{x}{1 + x }$	$f'(x) = \frac{1}{(1 + x)^2}$
Inverse square root unit (ISRU) ^[14]		$f(x) = \frac{x}{\sqrt{1 + \alpha x^2}}$	$f'(x) = \left(\frac{1}{\sqrt{1 + \alpha x^2}} \right)^3$

Activation function

Inverse square root linear unit (ISRLU) ^[14]		$f(x) = \begin{cases} \frac{x}{\sqrt{1+\alpha x^2}} & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \left(\frac{1}{\sqrt{1+\alpha x^2}} \right)^3 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Square Nonlinearity (SQNL) ^[11]		$f(x) = \begin{cases} 1 & : x > 2.0 \\ x - \frac{x^2}{4} & : 0 \leq x \leq 2.0 \\ x + \frac{x^2}{4} & : -2.0 \leq x < 0 \\ -1 & : x < -2.0 \end{cases}$	$f'(x) = 1 \mp \frac{x}{2}$
Rectified linear unit (ReLU) ^[15]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Bipolar rectified linear unit (BReLU) ^[16]		$f(x_i) = \begin{cases} ReLU(x_i) & \text{if } i \bmod 2 = 0 \\ -ReLU(-x_i) & \text{if } i \bmod 2 \neq 0 \end{cases}$	$f'(x_i) = \begin{cases} ReLU'(x_i) & \text{if } i \bmod 2 = 0 \\ -ReLU'(-x_i) & \text{if } i \bmod 2 \neq 0 \end{cases}$
Leaky rectified linear unit (Leaky ReLU) ^[17]		$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0.01 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Parametric rectified linear unit (PReLU) ^[18]		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Randomized leaky rectified linear unit (RRReLU) ^[19]		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases} \quad [3]$	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$

So, 1. **what exactly is deep learning ?**

And, 2. **why is it generally better** than other methods on image, speech and certain other types of data?

So, 1. **what exactly is deep learning ?**

And, 2. **why is it generally better** than other methods on image, speech and certain other types of data?

The short answers

- 1. ‘Deep Learning’ means using a neural network with several layers of nodes between input and output**
- 2. the series of layers between input & output do feature identification and processing in a series of stages, just as our brains seem to.**

hmmm... OK, but:

3. multilayer neural networks have been around for 25 years. What's actually new?

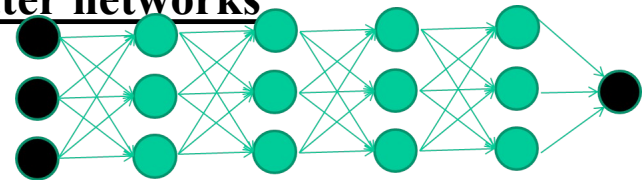
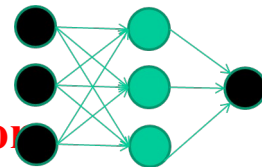
hmmm... OK, but:

3. multilayer neural networks have been around for 25 years. What's actually new?

we have always had good algorithms for learning the weights in networks with 1 hidden layer

but these algorithms are not good at learning the weights for networks with more hidden layers

what's new is: algorithms for training many-layer networks



longer answers

1. reminder/quick-explanation of how neural network weights are learned;
2. the idea of **unsupervised feature learning** (why ‘intermediate features’ are important for difficult classification tasks, and how NNs seem to naturally learn them)
3. The ‘breakthrough’ – the simple trick for training Deep neural networks

A dataset

<i>Fields</i>	<i>class</i>
---------------	--------------

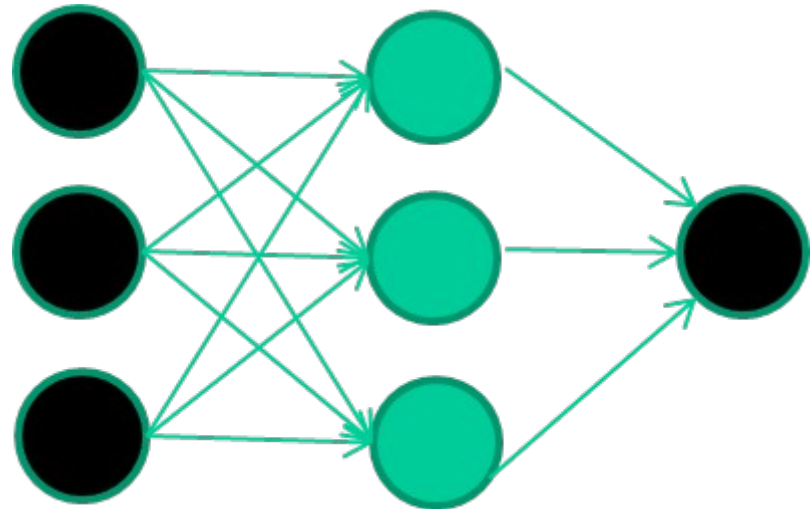
1.4 2.7 1.9	0
-------------	---

3.8 3.4 3.2	0
-------------	---

6.4 2.8 1.7	1
-------------	---

4.1 0.1 0.2	0
-------------	---

etc ...



Training the neural network

Fields ***class***

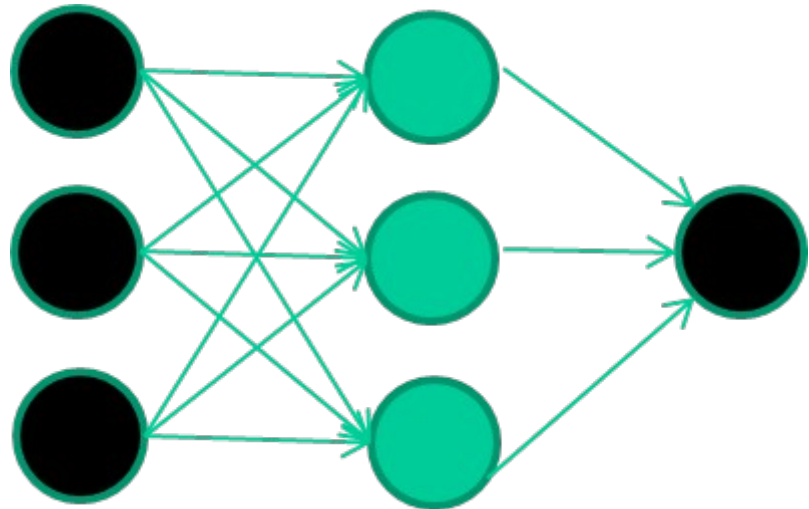
1.4 2.7 1.9 0

3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...



Training data

Fields ***class***

1.4 2.7 1.9 0

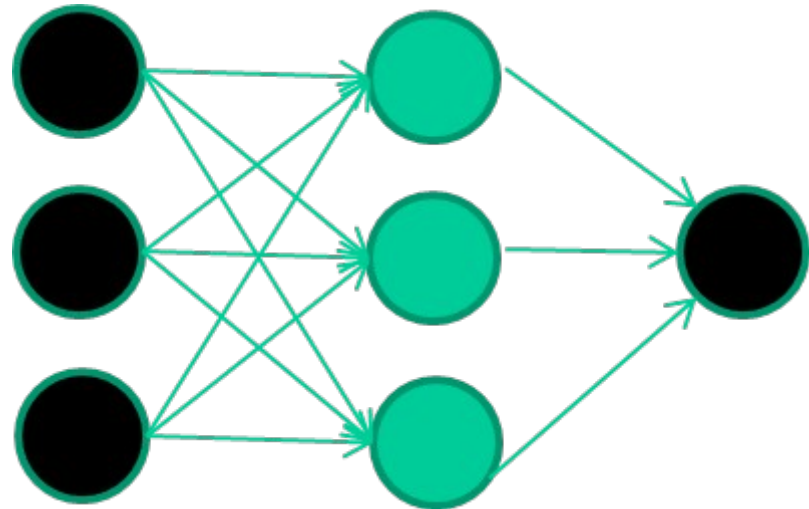
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Initialise with random weights



Training data

Fields ***class***

1.4 2.7 1.9 0

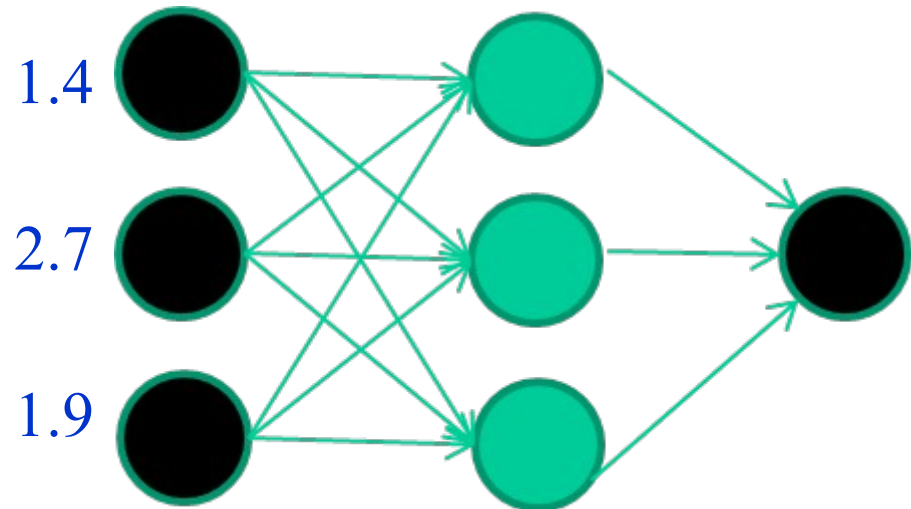
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Present a training pattern



Training data

Fields *class*

1.4 2.7 1.9 0

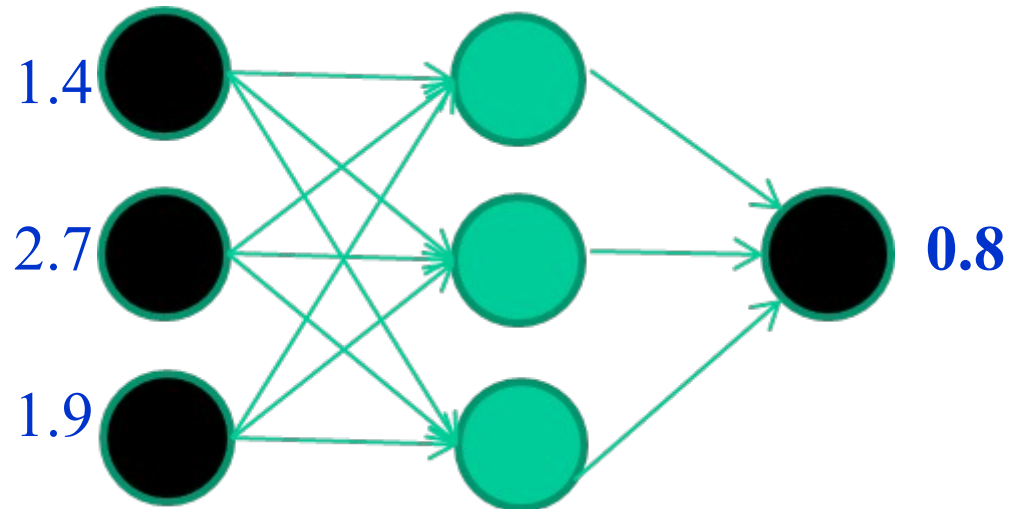
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Feed it through to get output



Training data

Fields *class*

1.4 2.7 1.9 0

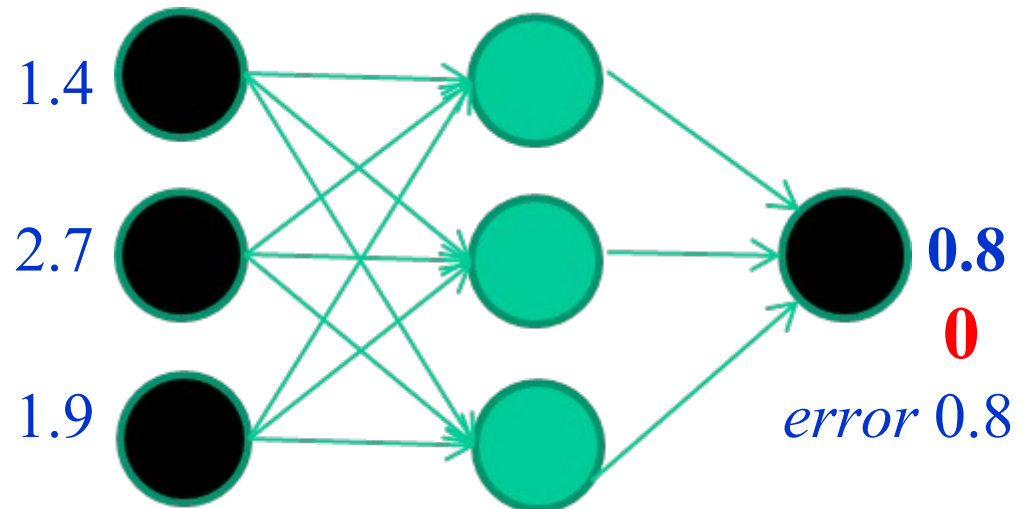
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Compare with target output



Training data

Fields *class*

1.4 2.7 1.9 0

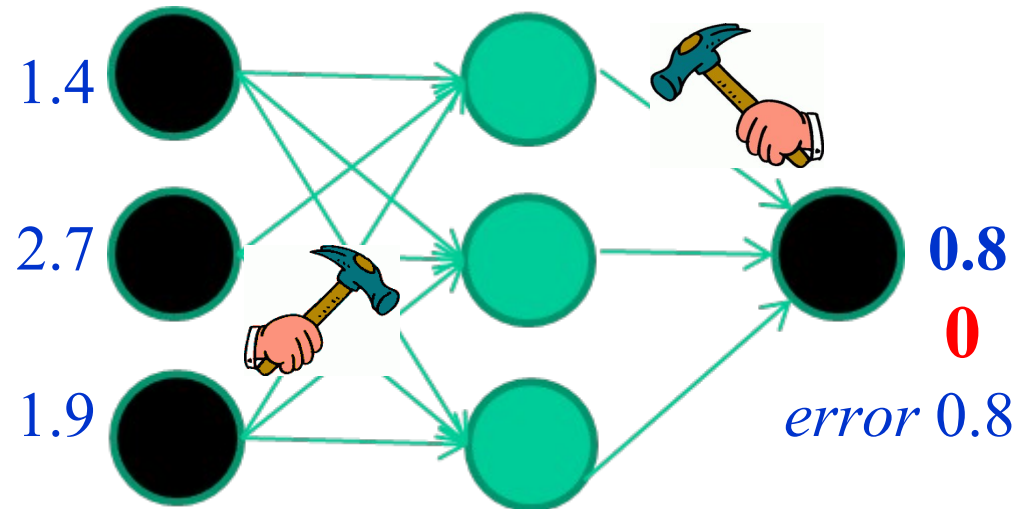
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Adjust weights based on error



Training data

Fields ***class***

1.4 2.7 1.9 0

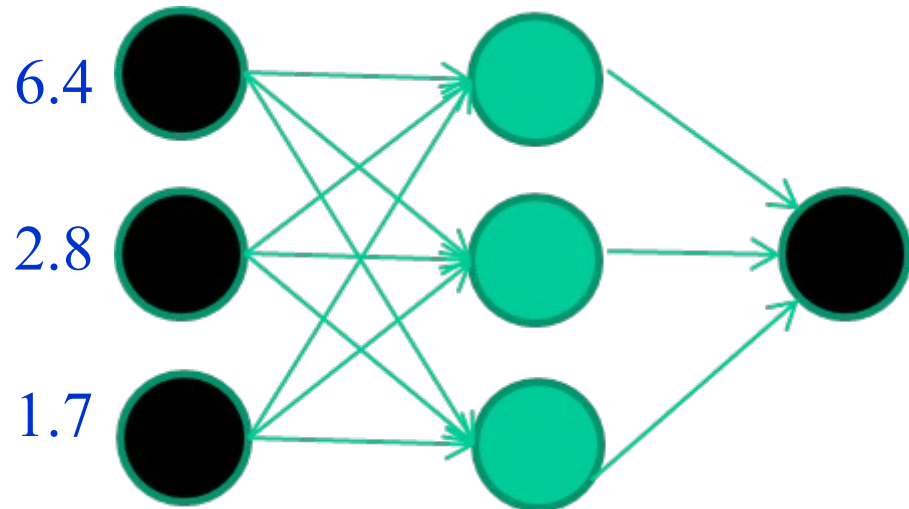
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Present a training pattern



Training data

Fields ***class***

1.4 2.7 1.9 0

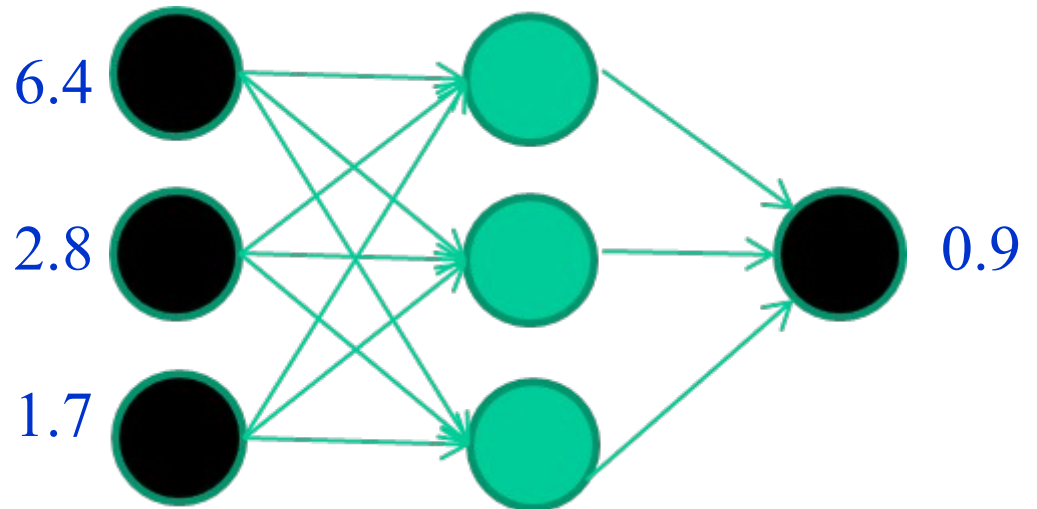
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Feed it through to get output



Training data

Fields ***class***

1.4 2.7 1.9 0

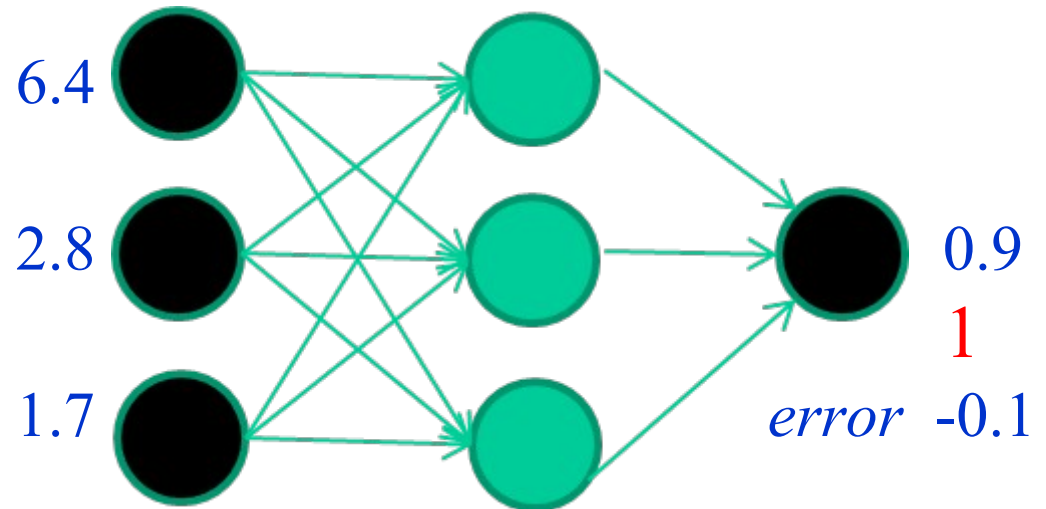
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Compare with target output



Training data

Fields ***class***

1.4 2.7 1.9 0

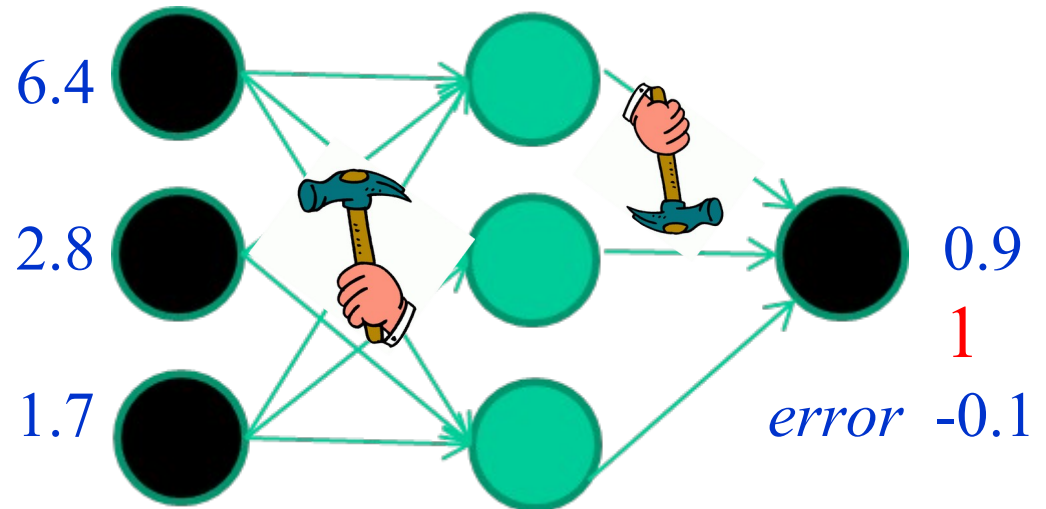
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

Adjust weights based on error



Training data

Fields ***class***

1.4 2.7 1.9 0

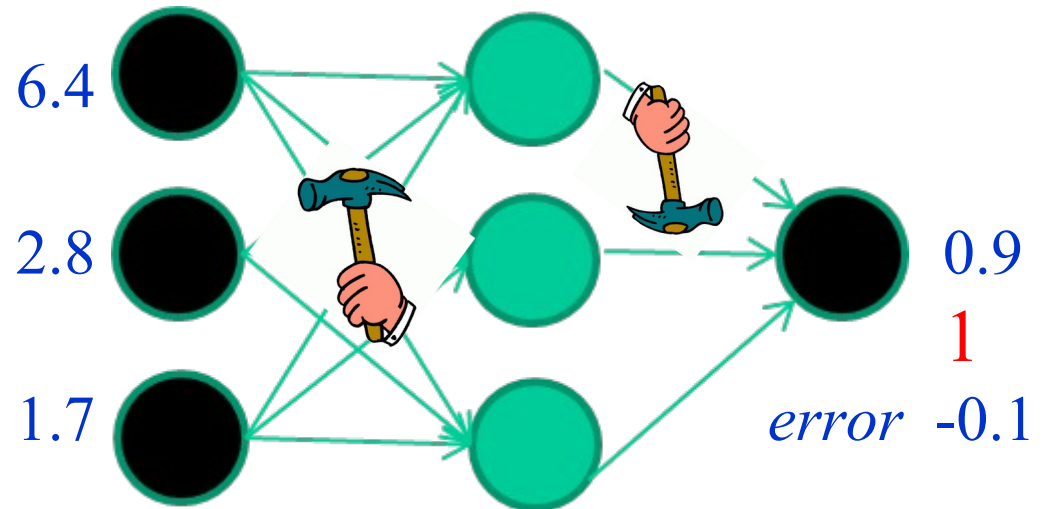
3.8 3.4 3.2 0

6.4 2.8 1.7 1

4.1 0.1 0.2 0

etc ...

And so on

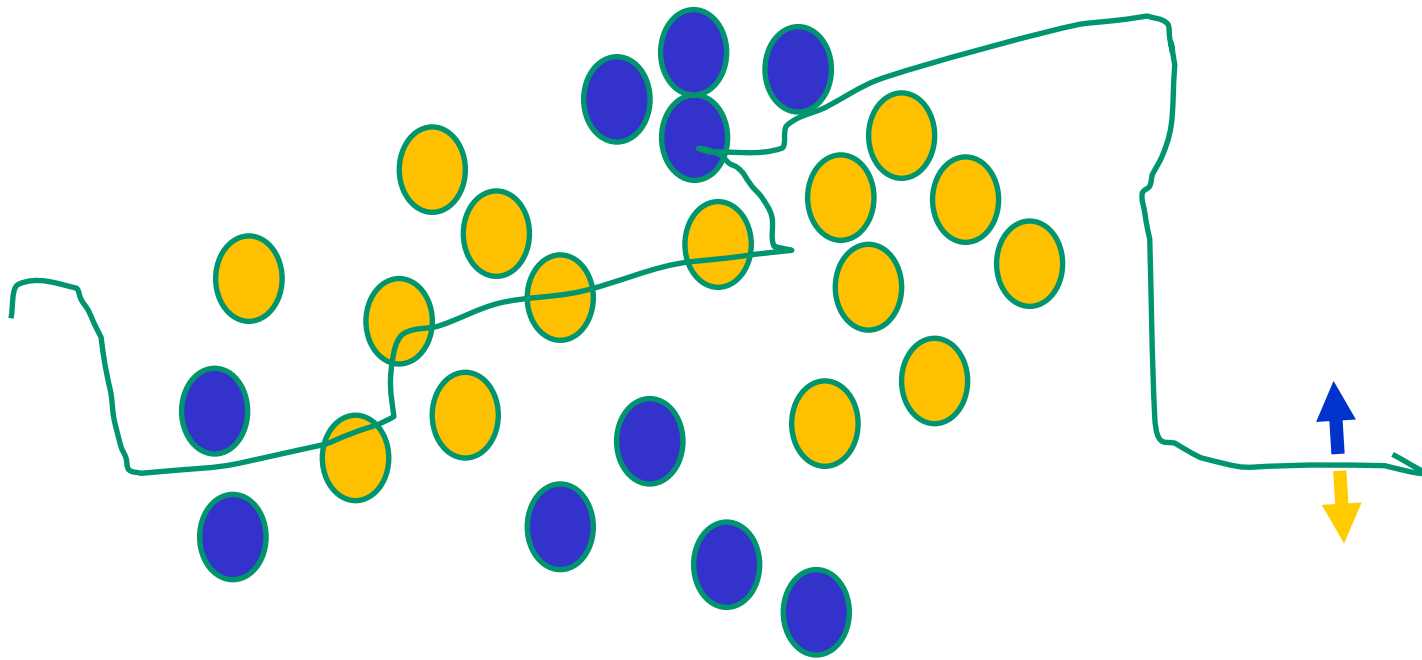


Repeat this thousands, maybe millions of times – each time taking a random training instance, and making slight weight adjustments

Algorithms for weight adjustment are designed to make changes that will reduce the error

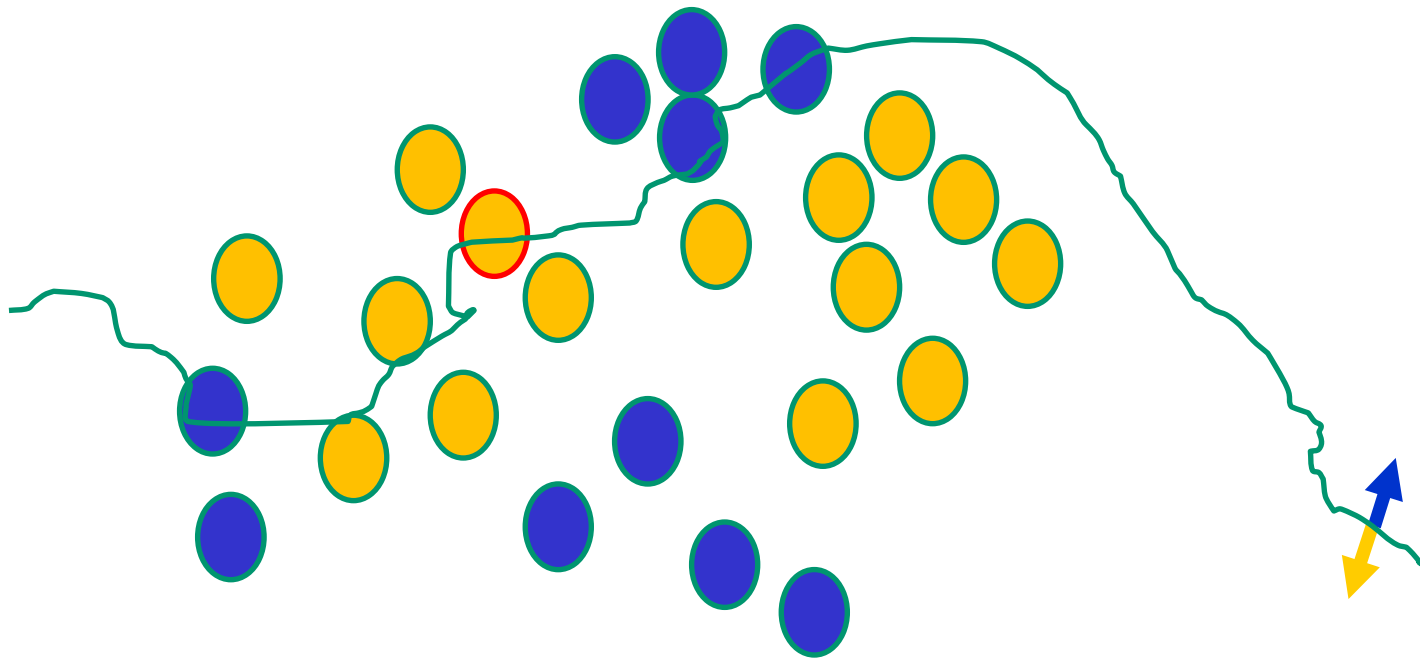
The decision boundary perspective...

Initial random weights



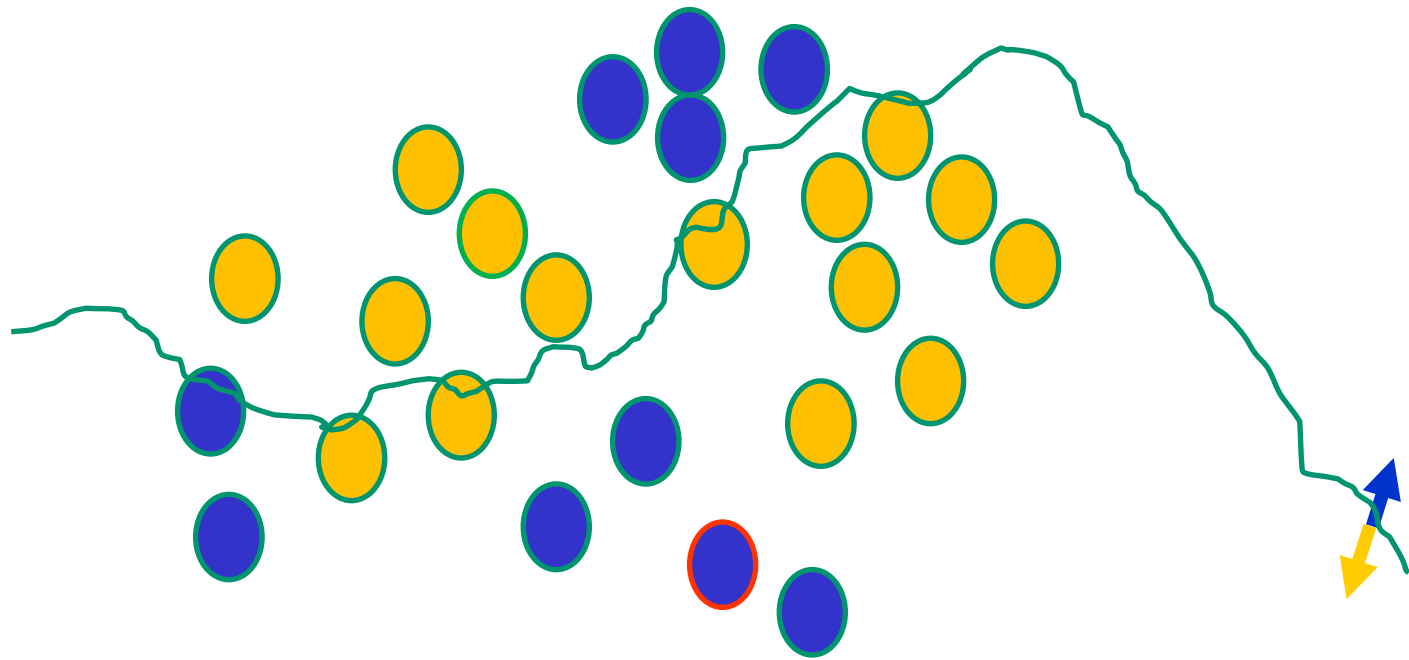
The decision boundary perspective...

Present a training instance / adjust the weights



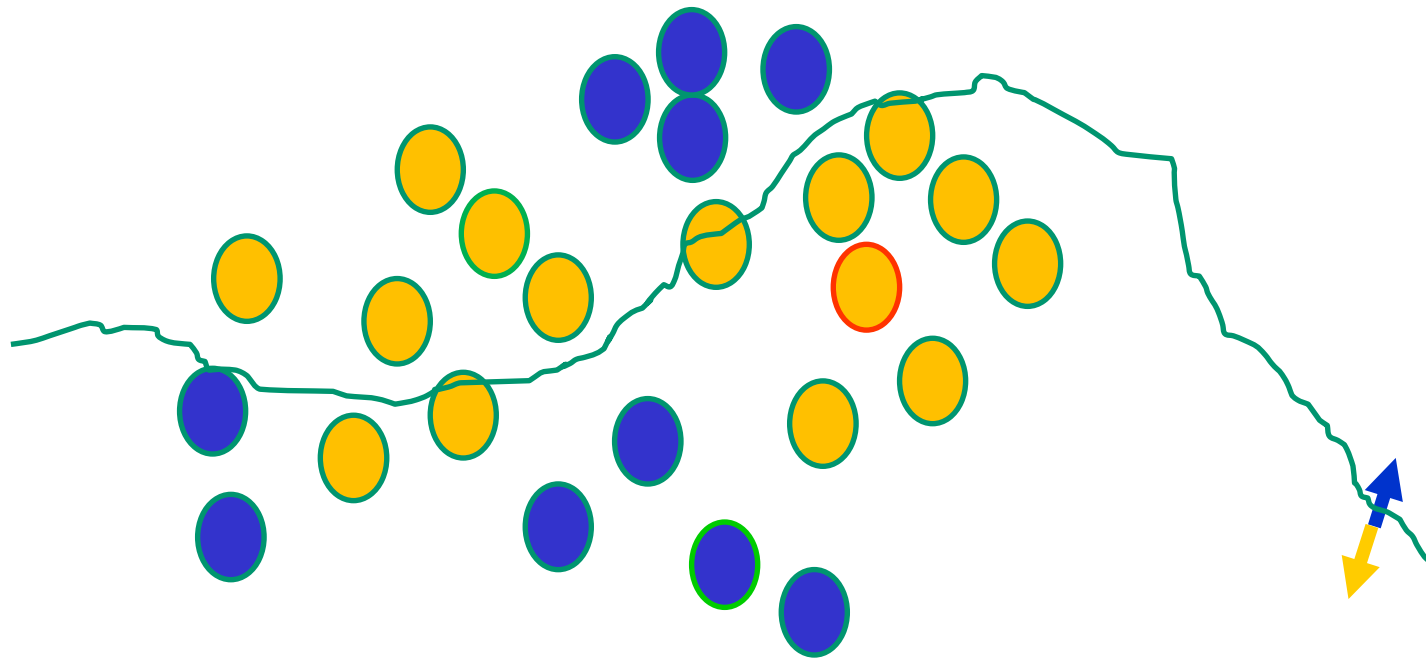
The decision boundary perspective...

Present a training instance / adjust the weights



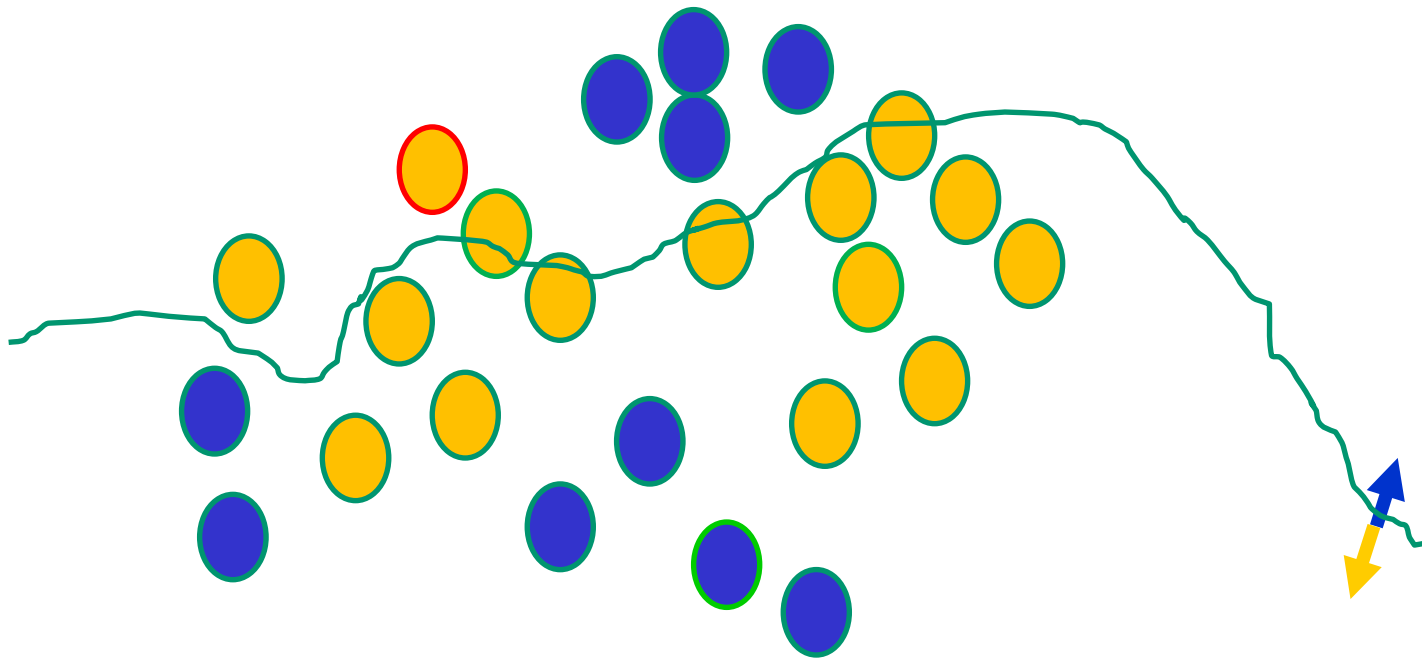
The decision boundary perspective...

Present a training instance / adjust the weights



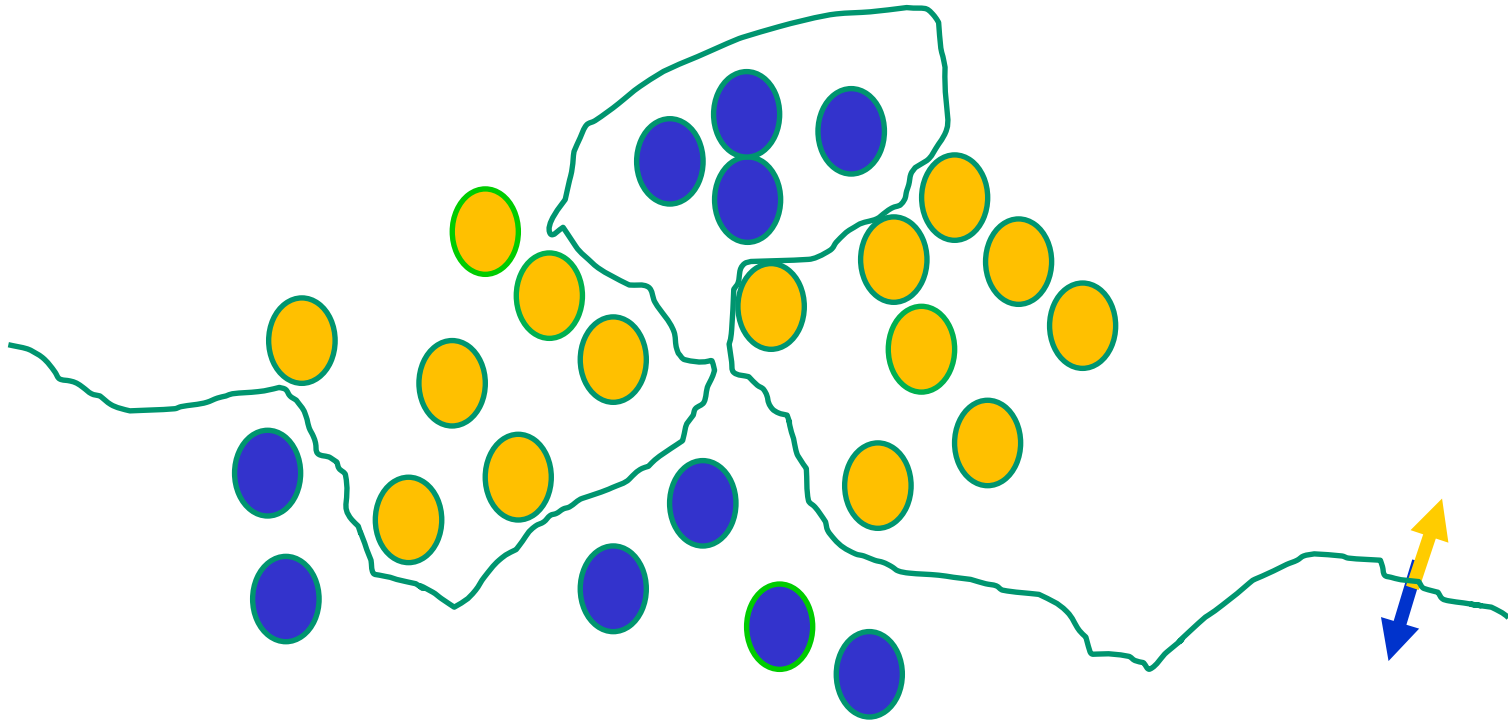
The decision boundary perspective...

Present a training instance / adjust the weights



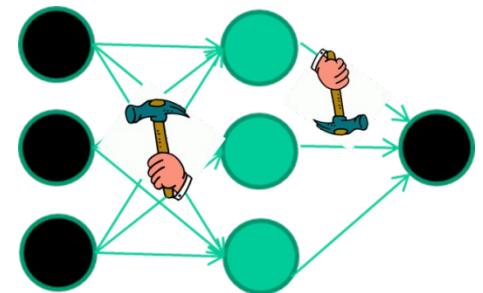
The decision boundary perspective...

Eventually



The point I am trying to make

- weight-learning algorithms for NNs are dumb
- they work by making thousands and thousands of tiny adjustments, each making the network do better at the most recent pattern, but perhaps a little worse on many others
- but, by dumb luck, eventually this tends to be good enough to learn effective classifiers for many real applications



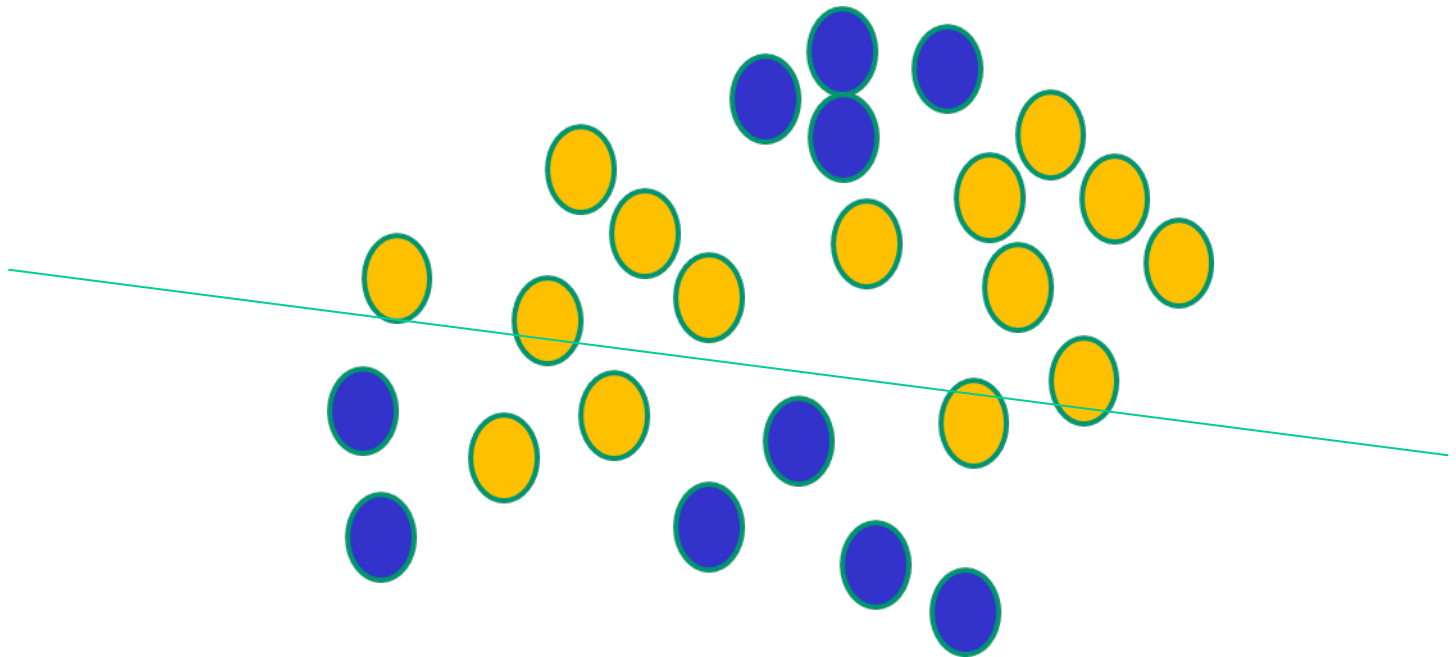
Some other points

Detail of a standard NN weight learning algorithm –
later

If $f(x)$ is **non-linear**, a network with 1 hidden layer can, in theory, learn perfectly any classification problem. A set of weights exists that can produce the targets from the inputs. The problem is finding them.

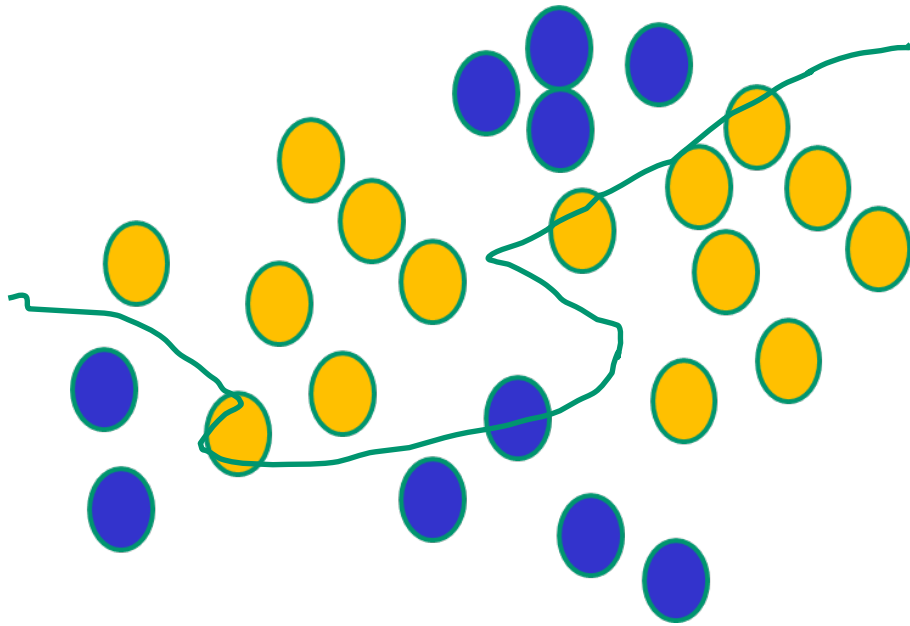
Some other ‘by the way’ points

If $f(x)$ is **linear**, the NN can **only** draw straight decision boundaries (even if there are many layers of units)



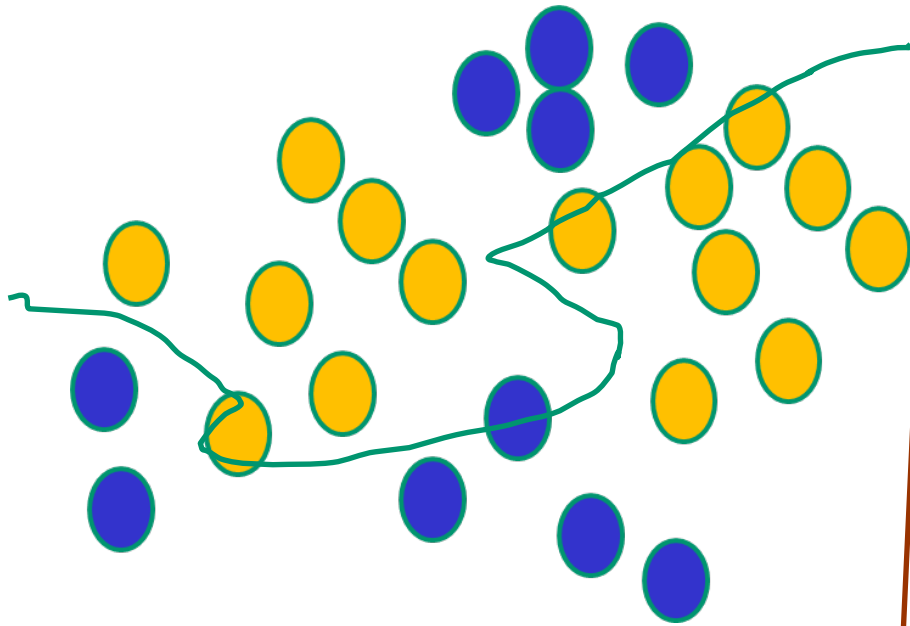
Some other ‘by the way’ points

NNs use nonlinear $f(x)$ so they
can draw complex boundaries,
but keep the data unchanged

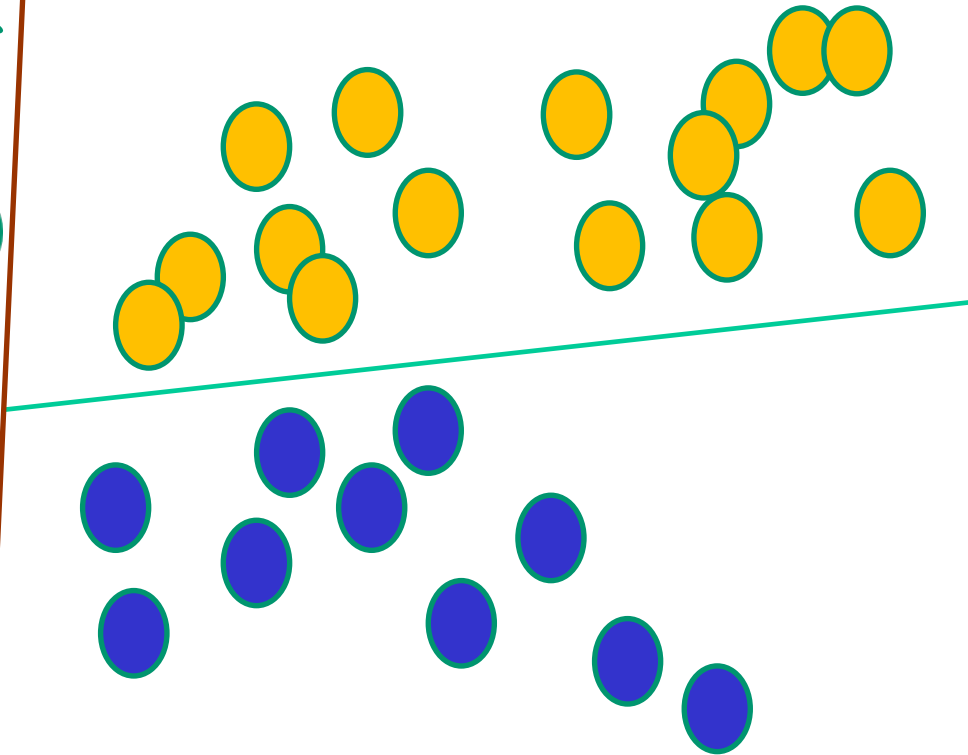


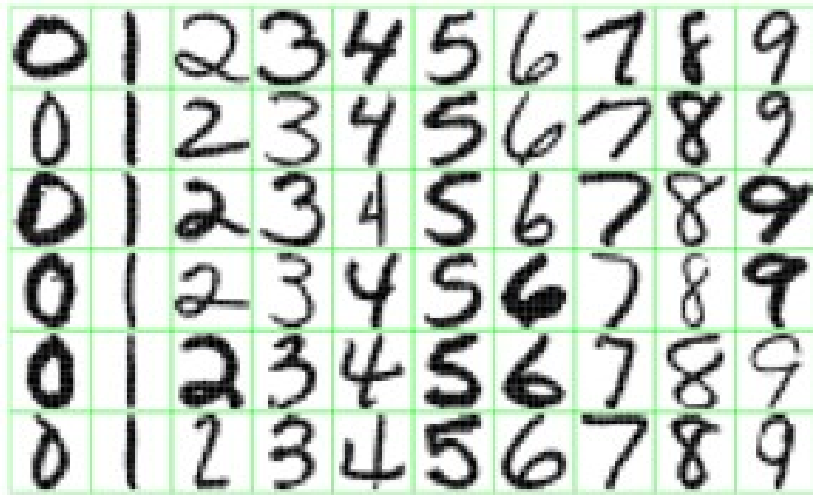
Some other ‘by the way’ points

NNs use nonlinear $f(x)$ so they can draw complex boundaries, but keep the data unchanged



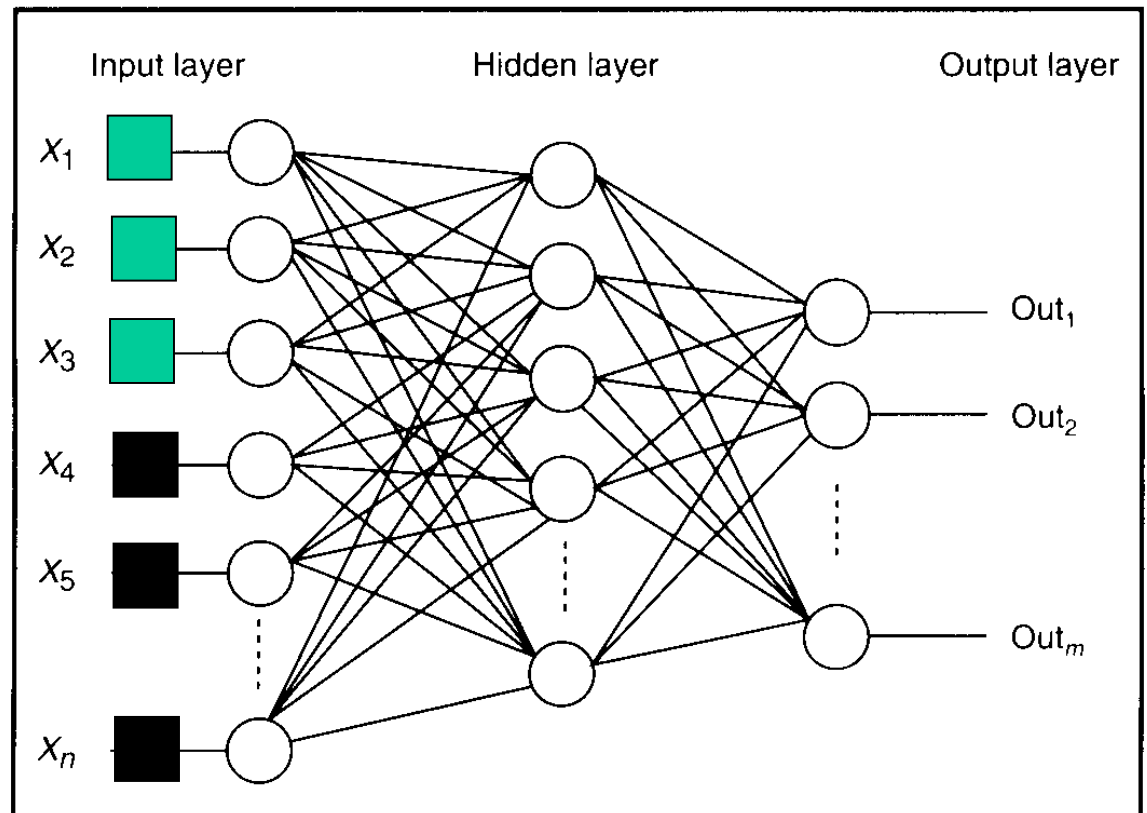
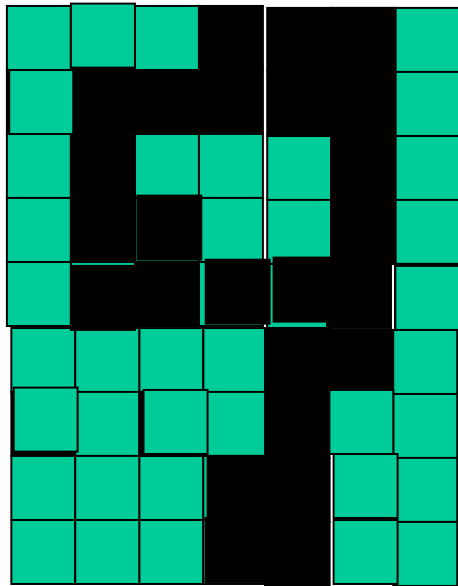
SVMs only draw straight lines, but they transform the data first in a way that makes that OK



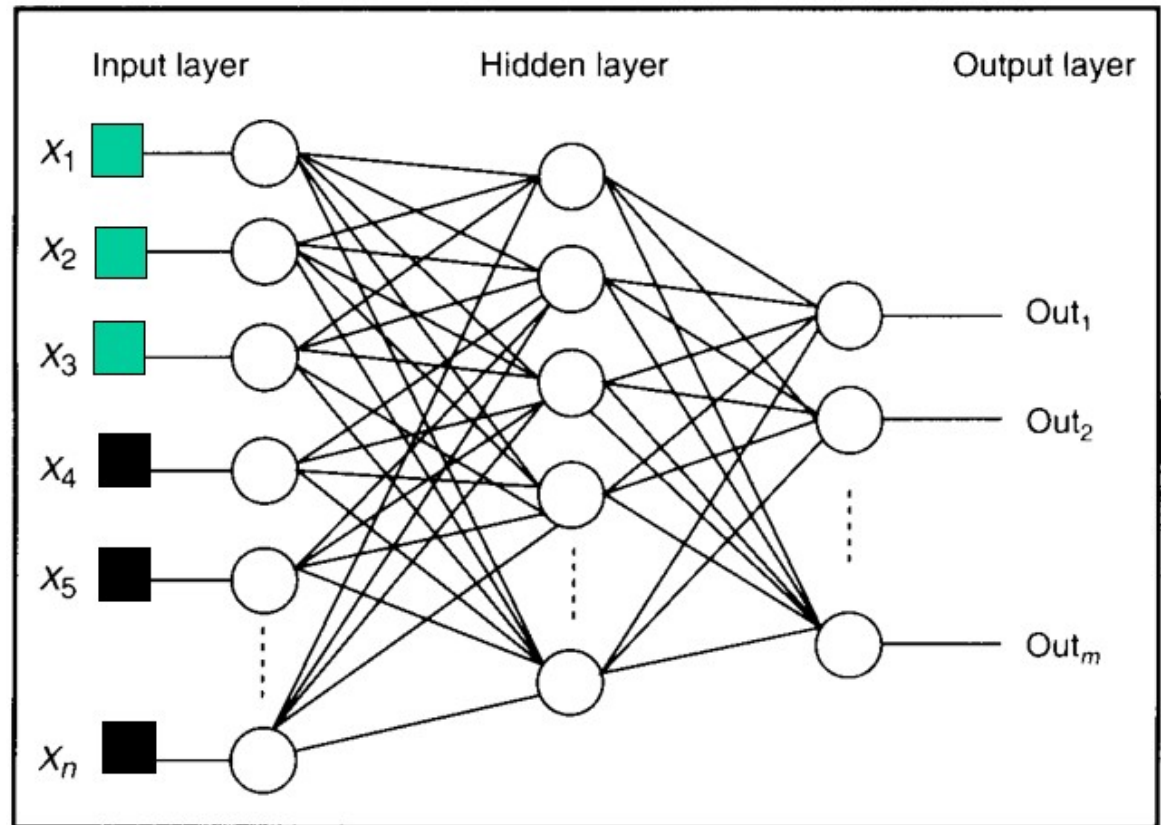
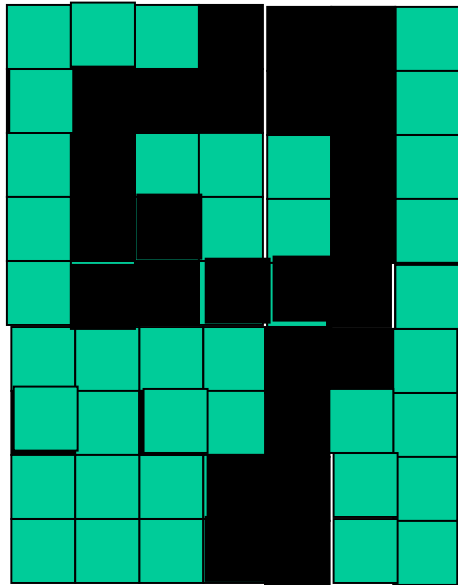


Feature detectors

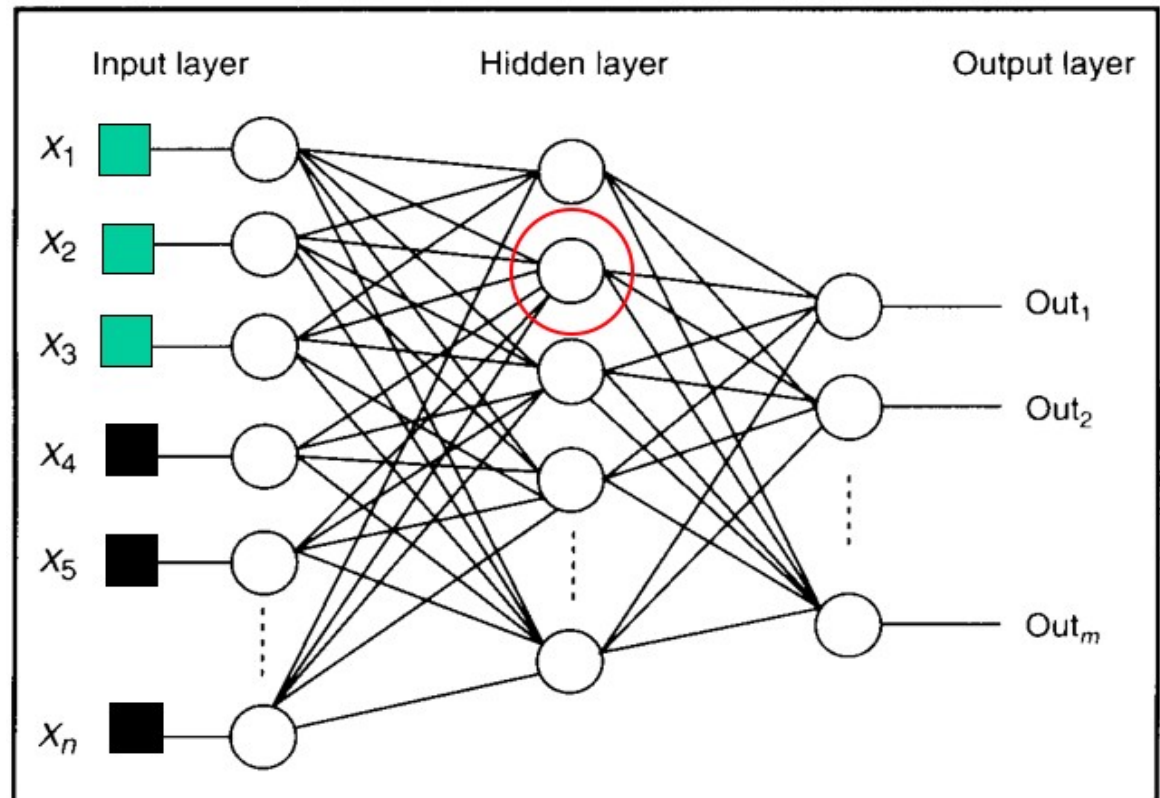
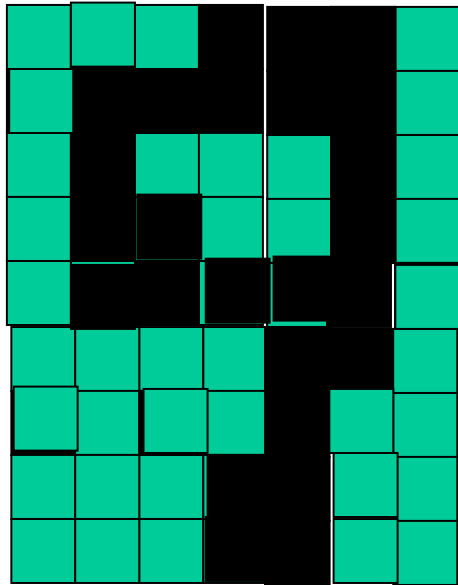
Figure 1.2: *Examples of handwritten digits postal envelopes.*



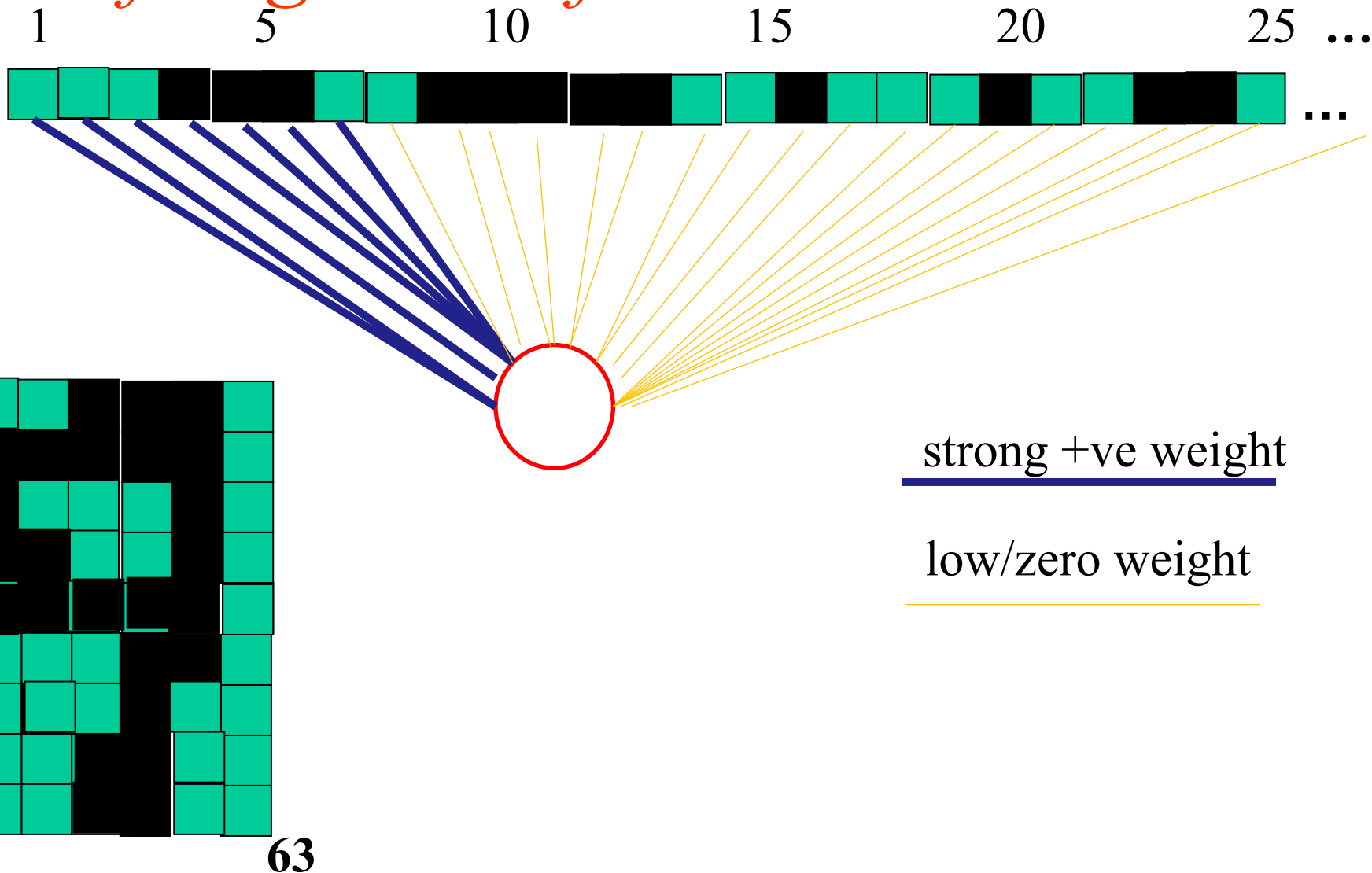
Feature detectors



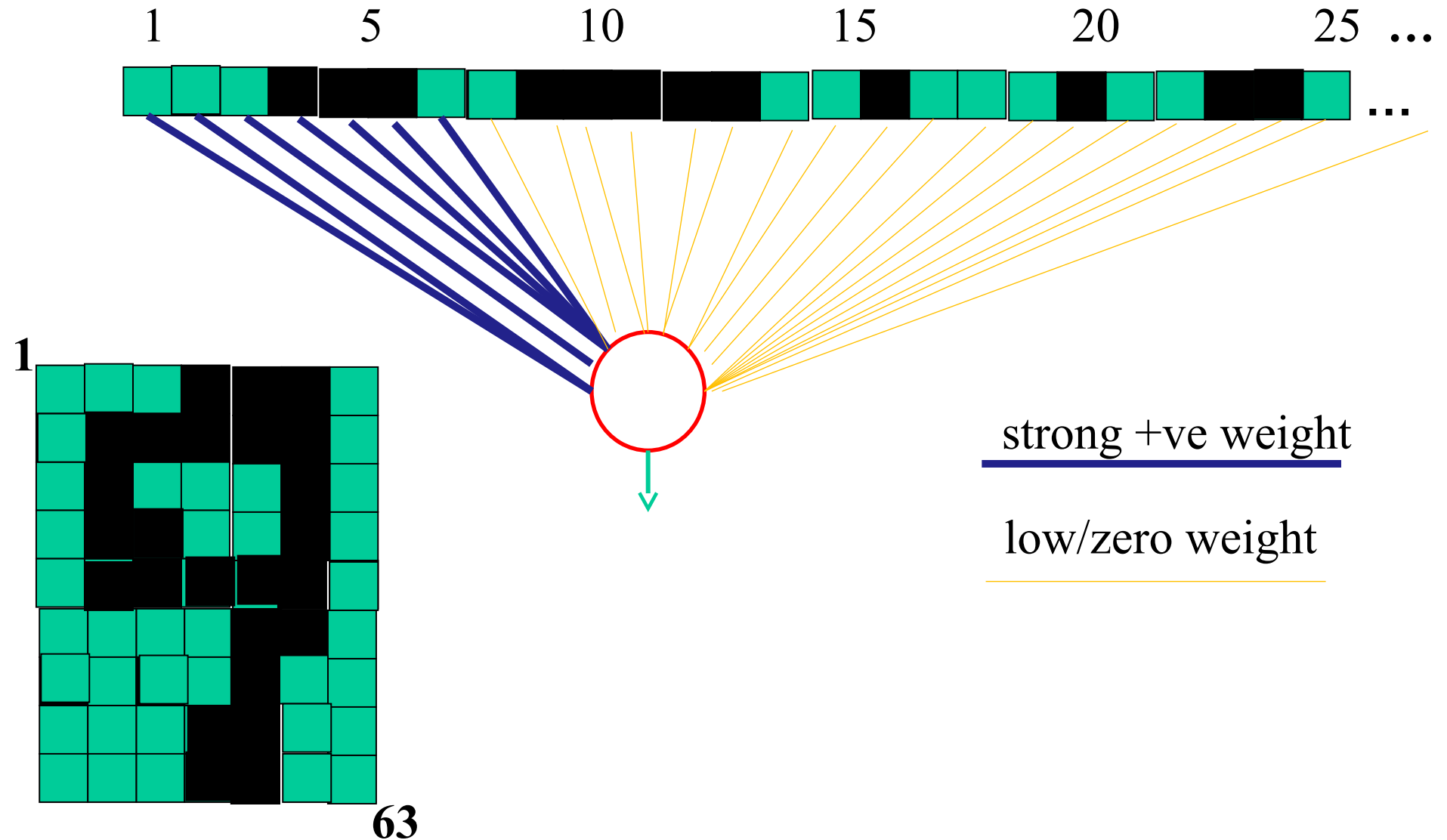
*what is this
unit doing?*



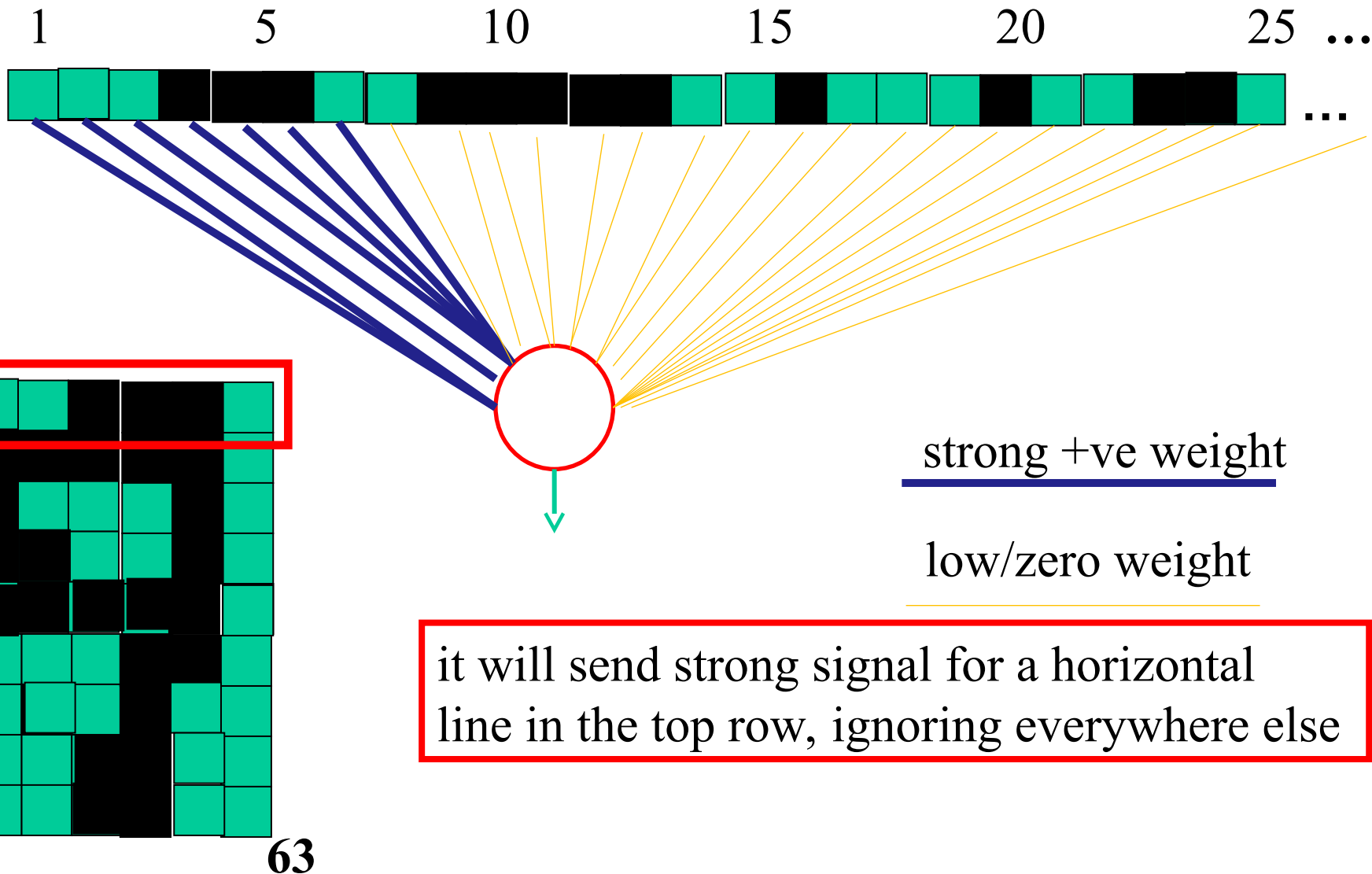
Hidden layer units become *self-organised feature detectors*



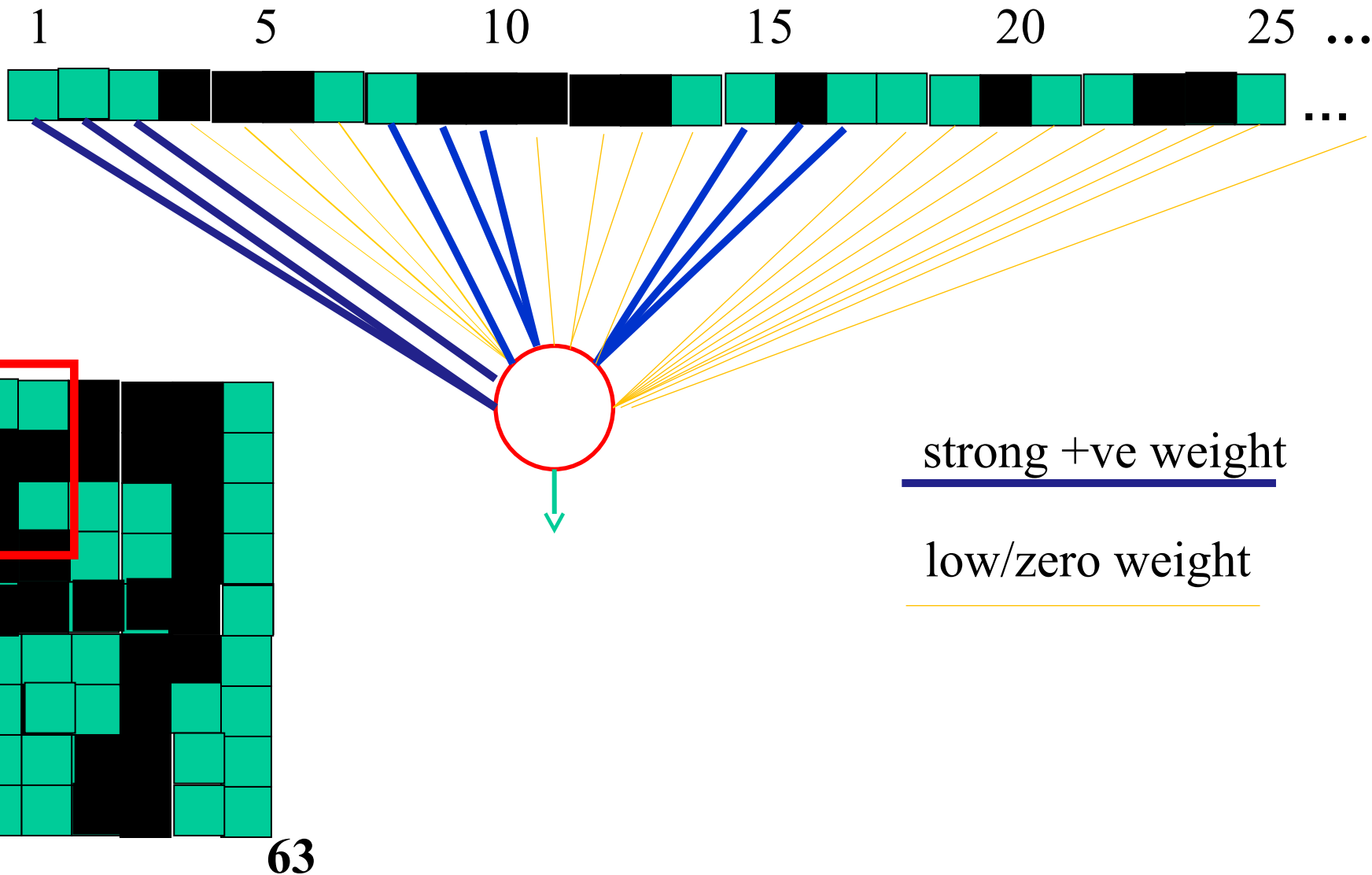
What does this unit detect?



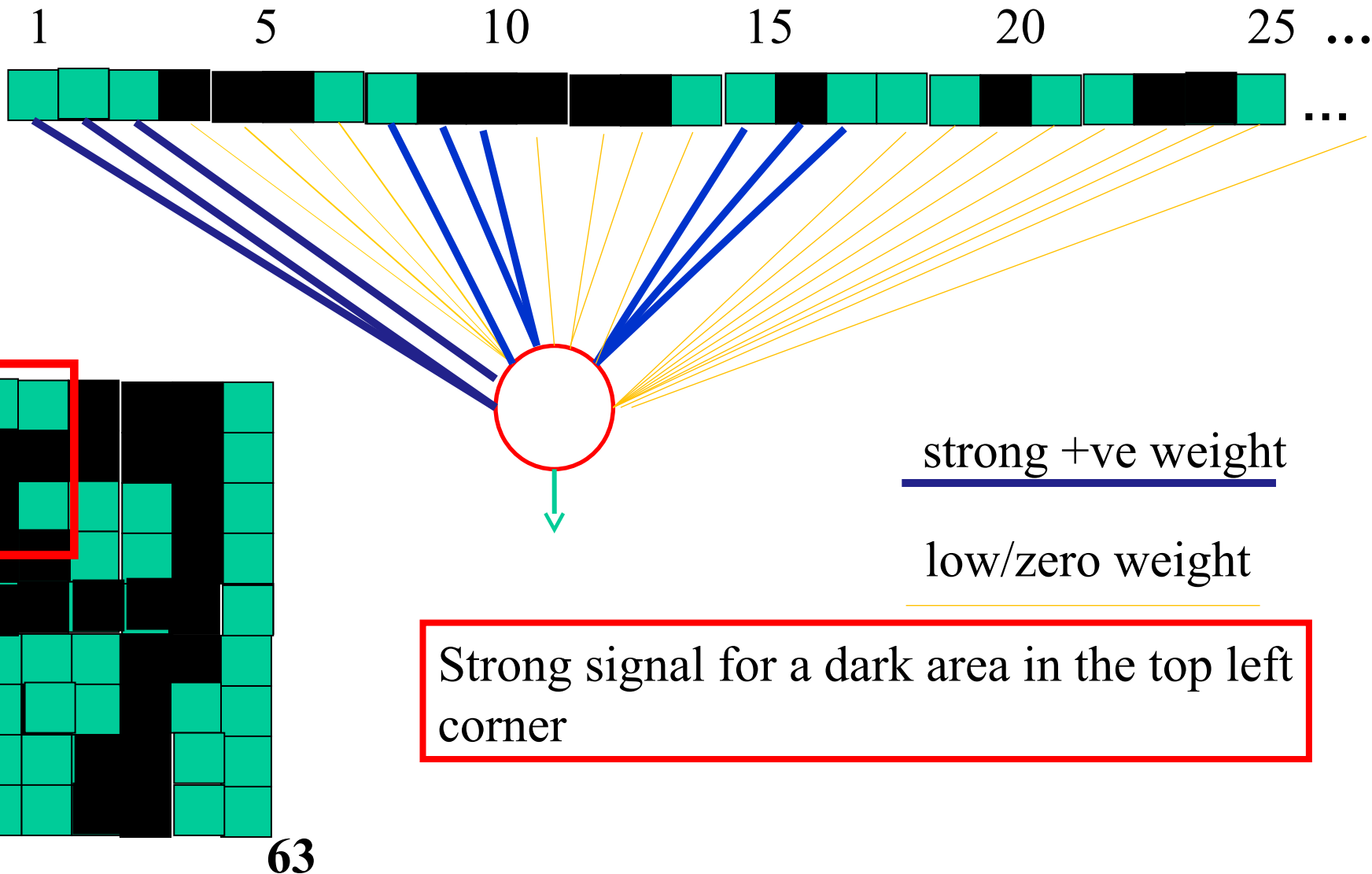
What does this unit detect?



What does this unit detect?



What does this unit detect?



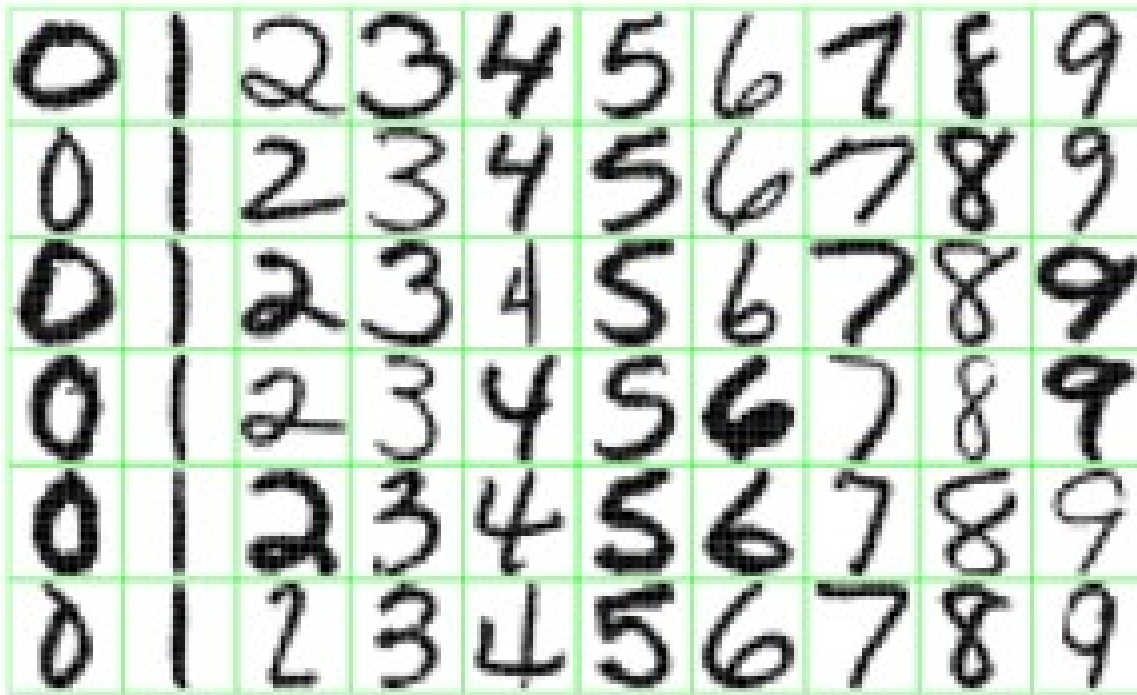


Figure 1.2: *Examples of handwritten digits from U.S. postal envelopes.*

What features might you expect a good NN to learn, when trained with data like this?

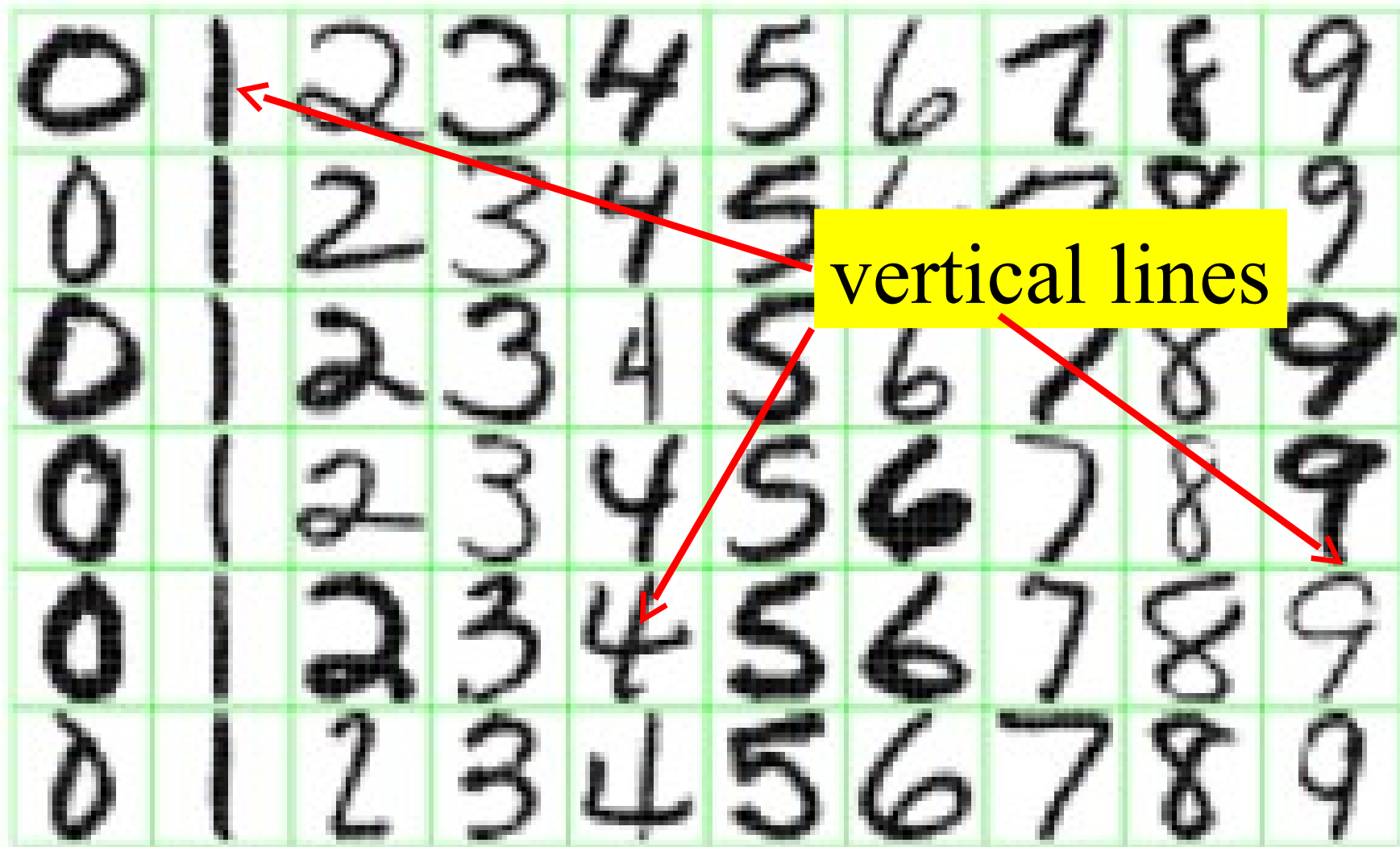


Figure 1.2: *Examples of handwritten digits from U.S. postal envelopes.*



Figure 1.2: *Examples of handwritten digits from U.S. postal envelopes.*

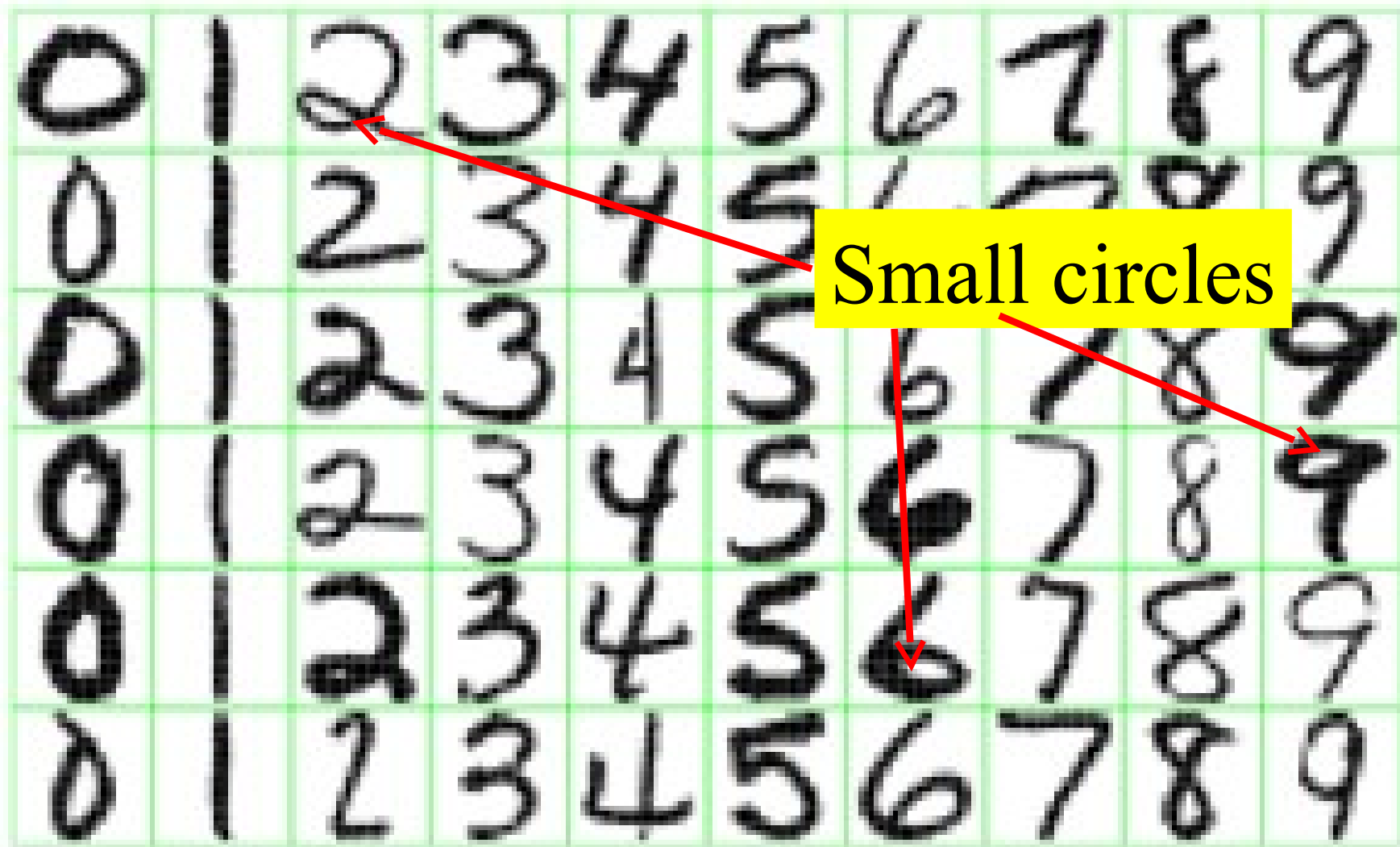
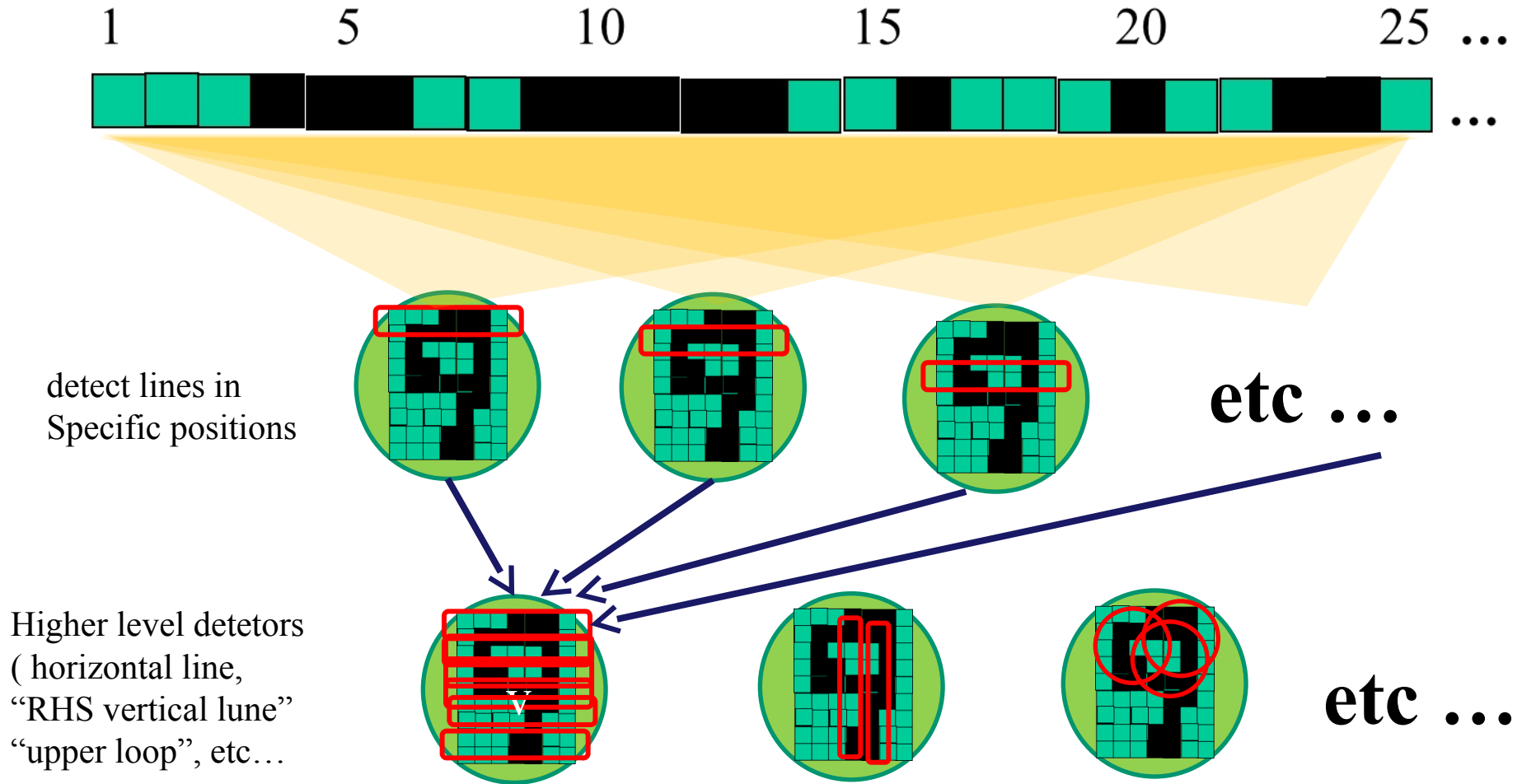


Figure 1.2: *Examples of handwritten digits from U.S. postal envelopes.*

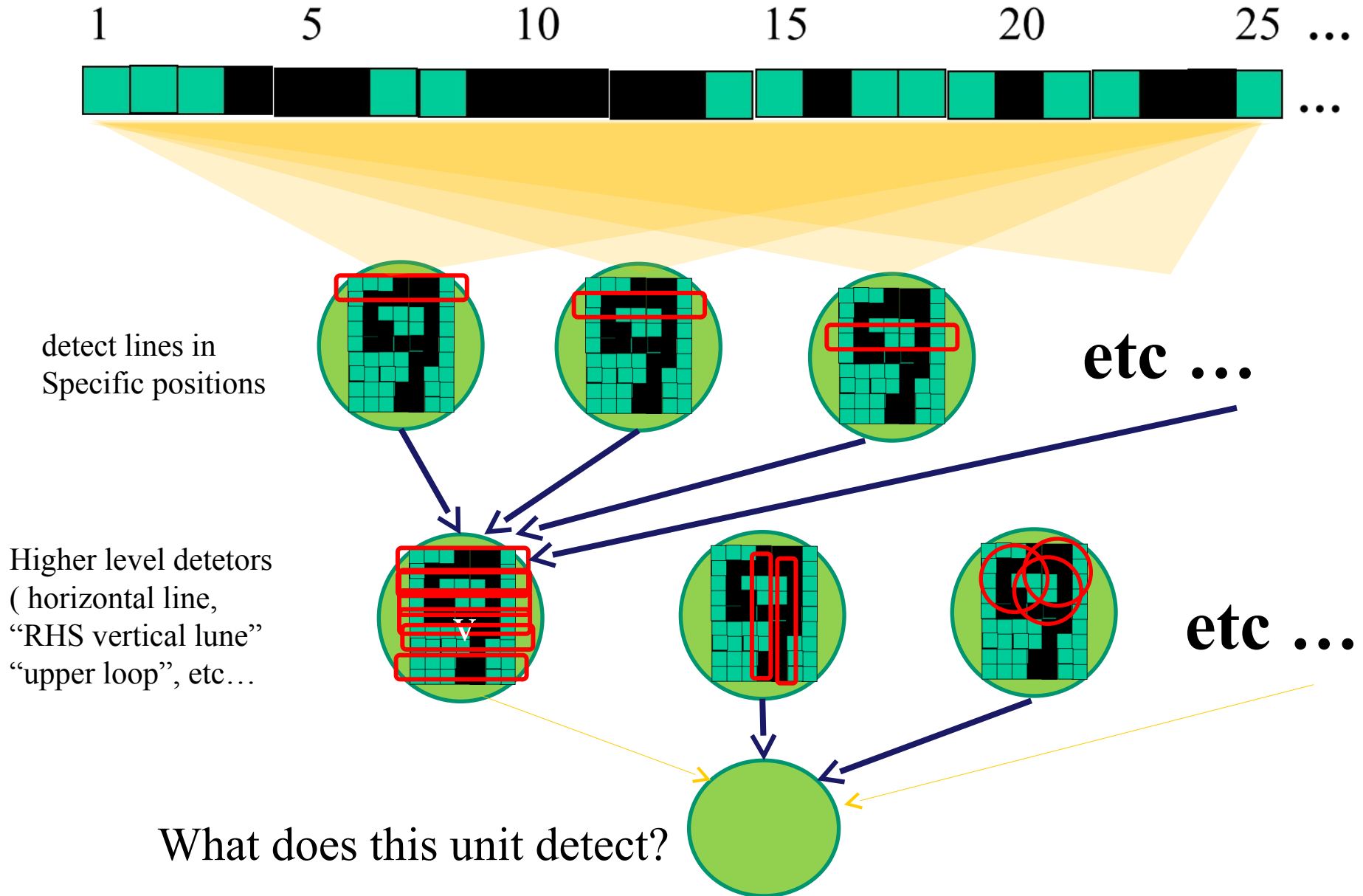


But what about position invariance ???
our example unit detectors were tied to
specific parts of the image

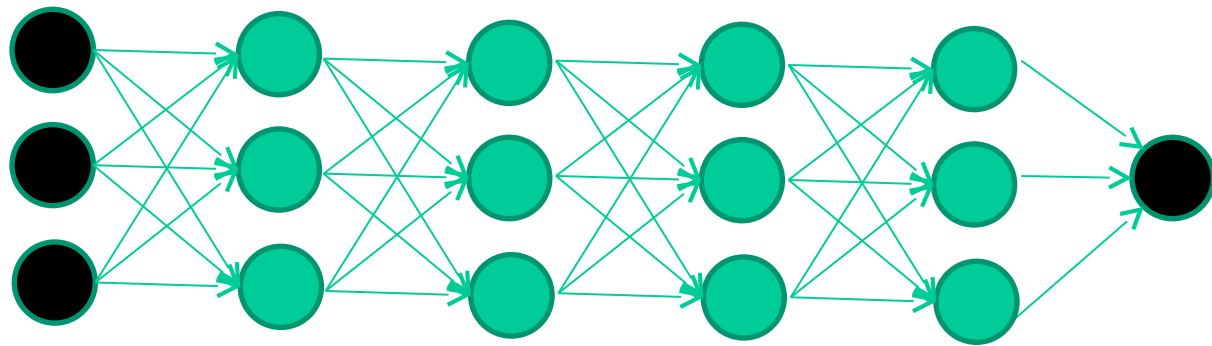
successive layers can learn higher-level features ...



successive layers can learn higher-level features ...

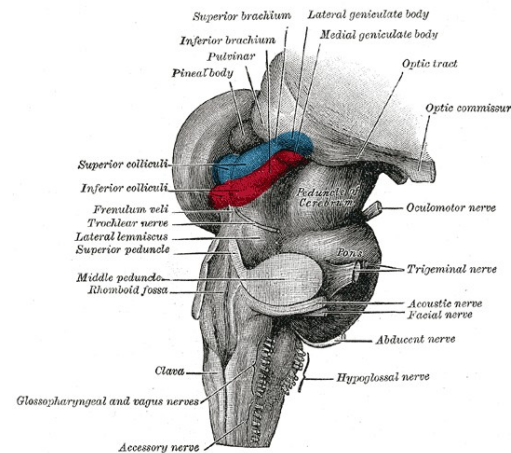
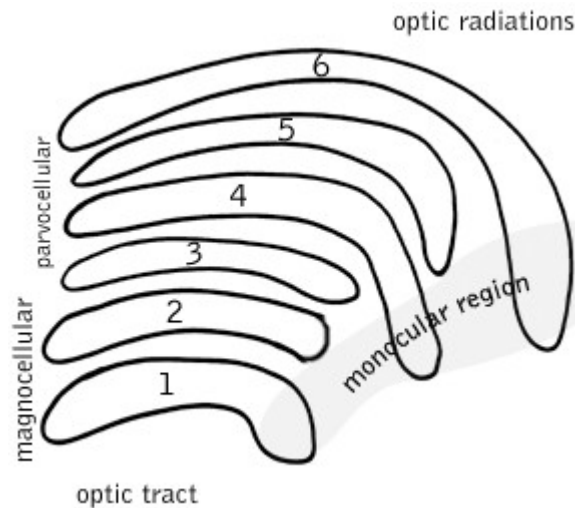


So: multiple layers make sense



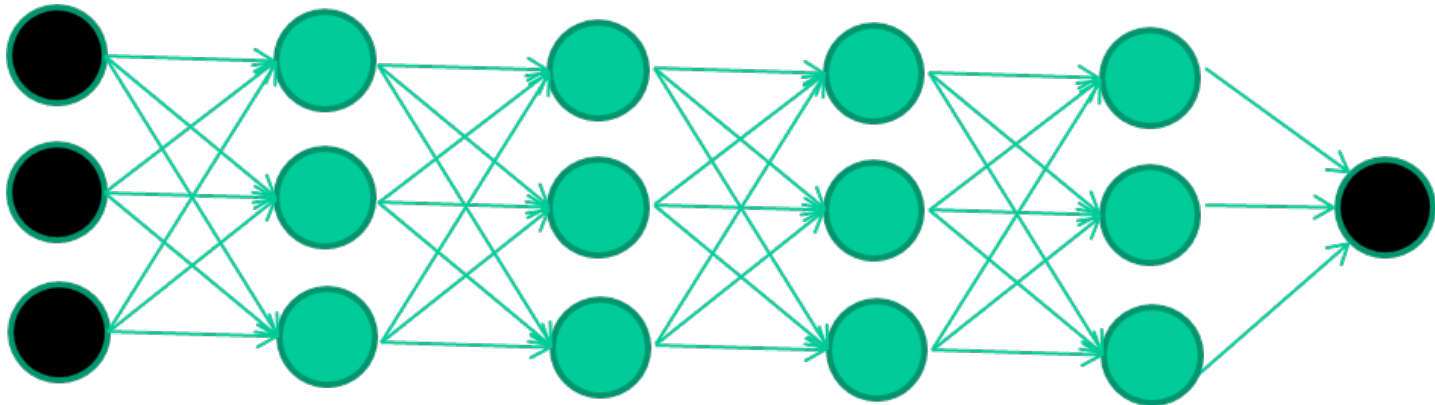
So: multiple layers make sense

Your brain works that way

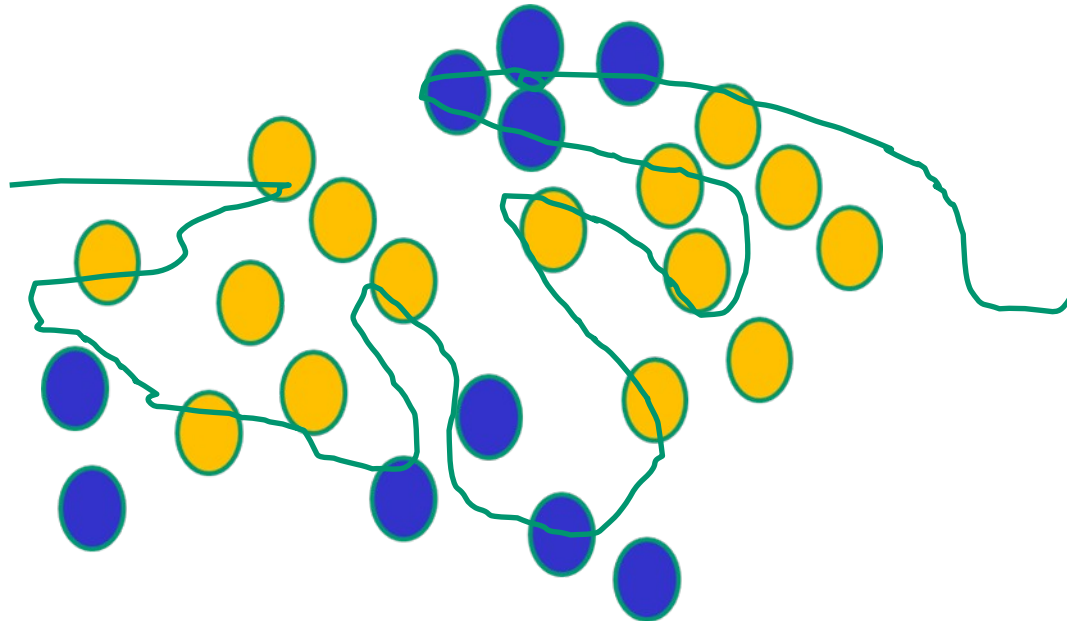
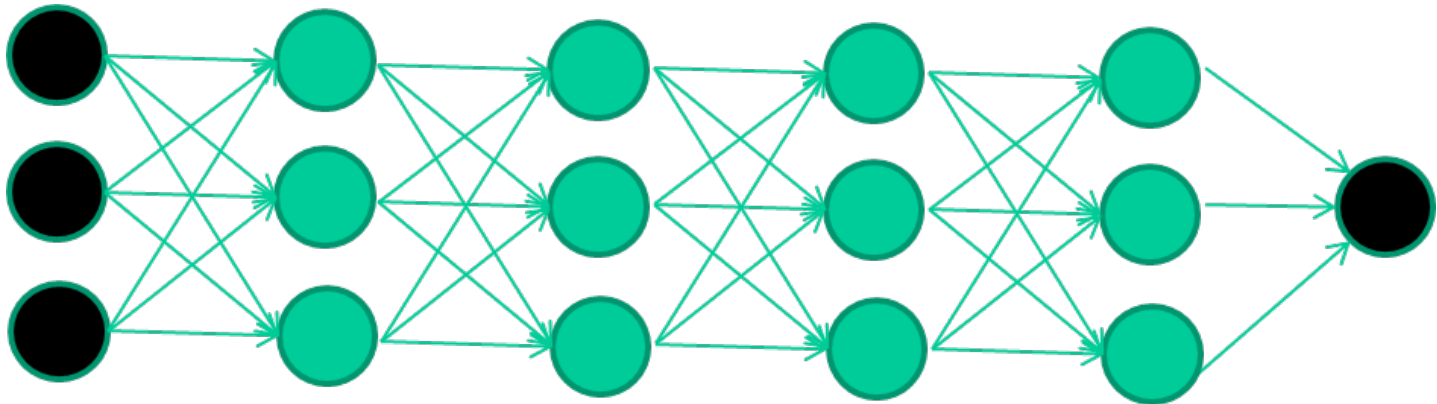


So: multiple layers make sense

Many-layer neural network architectures should be capable of learning the true underlying features and ‘feature logic’, and therefore generalise very well ...



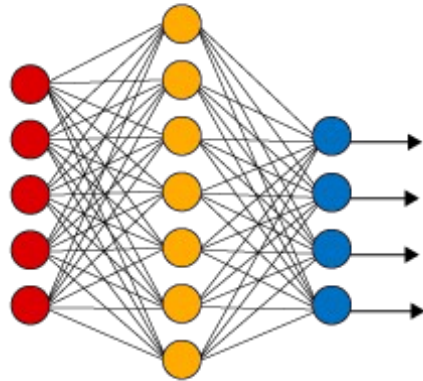
But, until very recently, our weight-learning algorithms simply did not work on multi-layer architectures



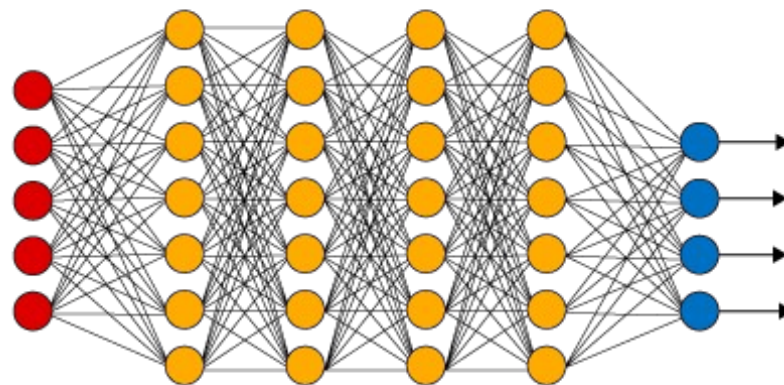
Along came deep learning ...

What is deep learning ?

Simple Neural Network



Deep Learning Neural Network



● Input Layer

● Hidden Layer

● Output Layer

A network with 1 hidden layer can, in theory, learn perfectly any classification problem. A set of weights exists that can produce the targets from the inputs. The problem is finding them.

visual routines

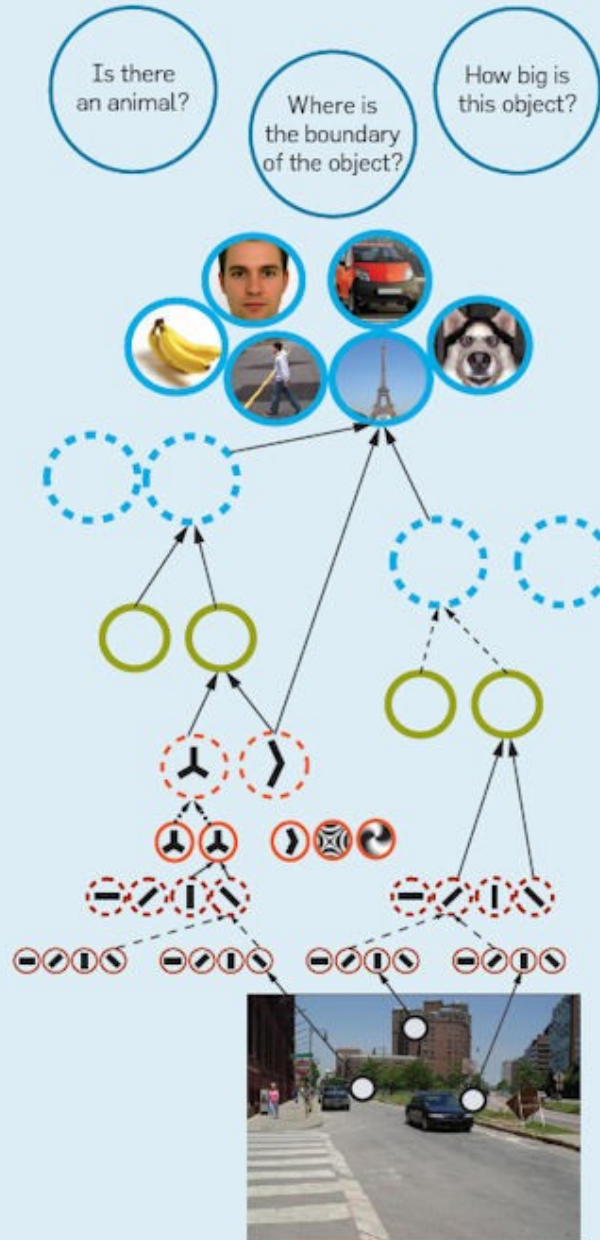
AIT

PIT

V2-V4

V1

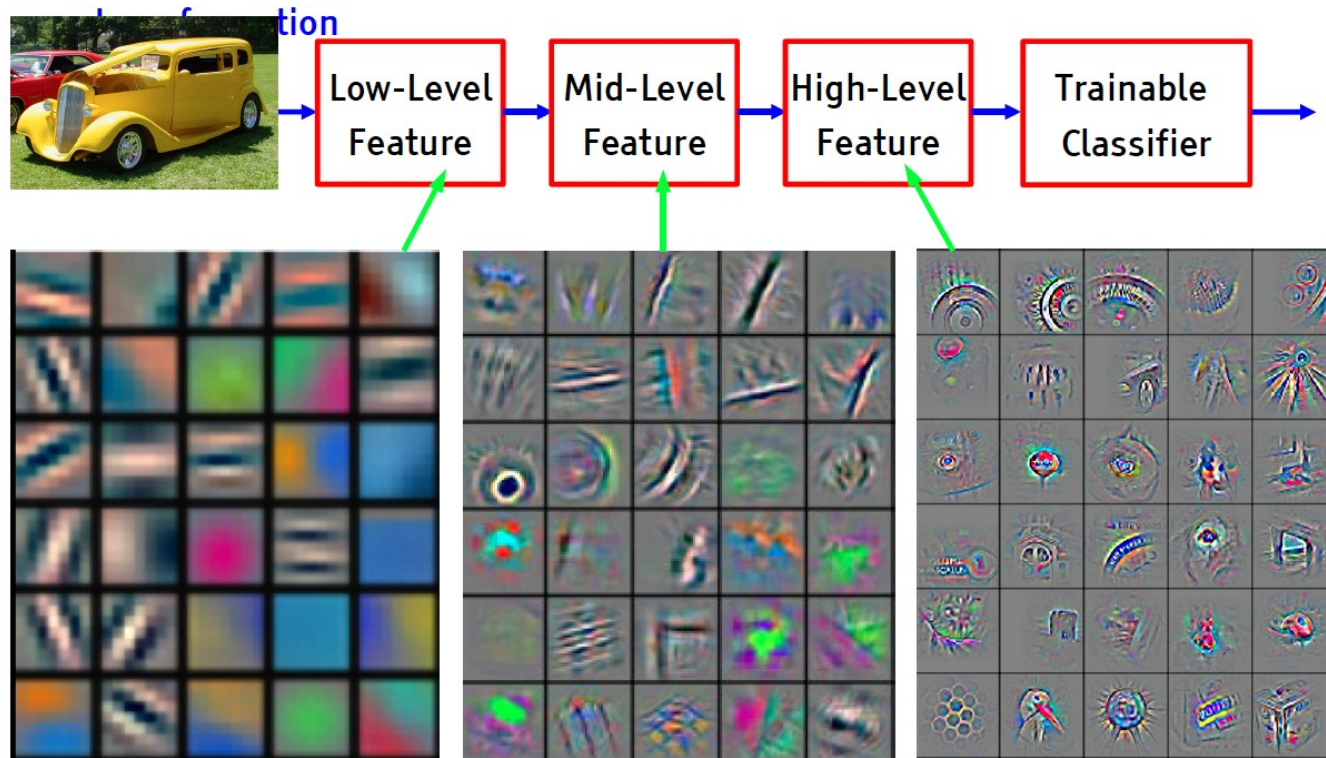
○ Complex units
○ Simple units



Hierarchical models

Riesenhuber & Poggio. Nature Neurosci 1999

Deep Learning = Learning Hierarchical Representations

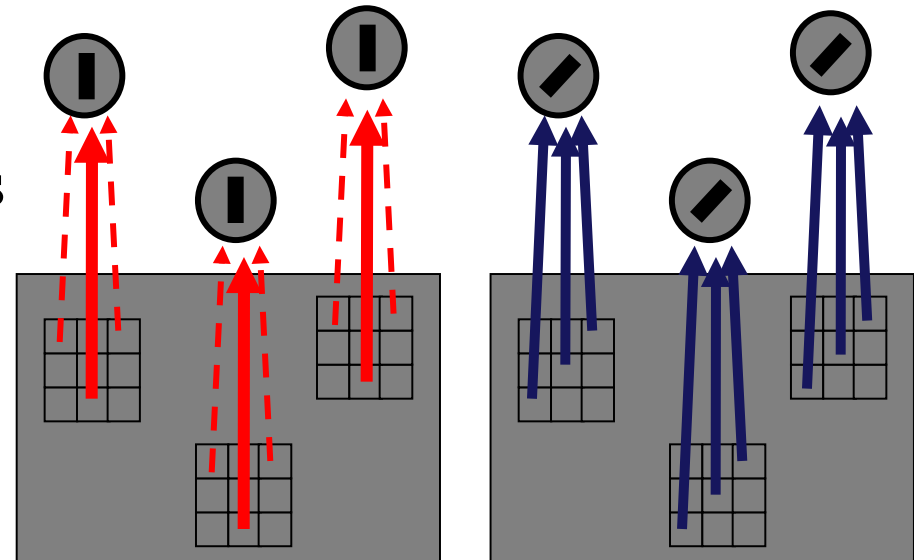


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

Convolutional Networks (ConvNet or CNN)

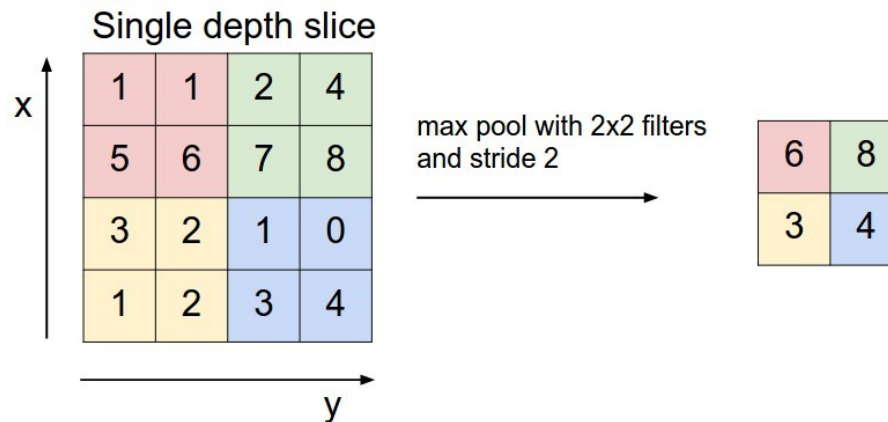
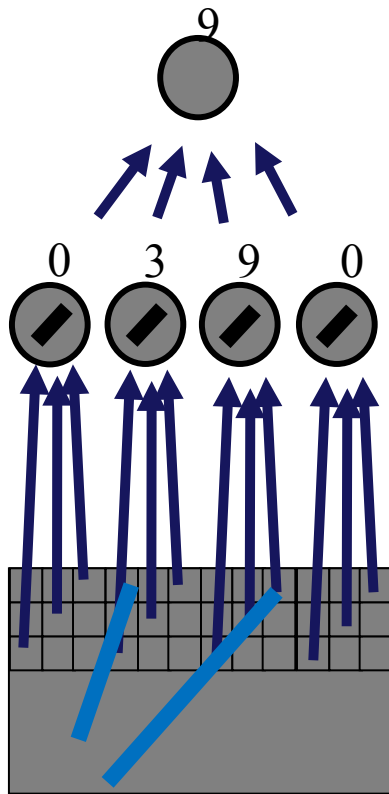
(currently the dominant approach for neural networks)

- Use many different copies of the same feature detector with different positions.
 - Replication greatly reduces the number of free parameters to be learned.
- Use several different feature types, each with its own map of replicated detectors.
 - Allows each patch of image to be represented in several ways.

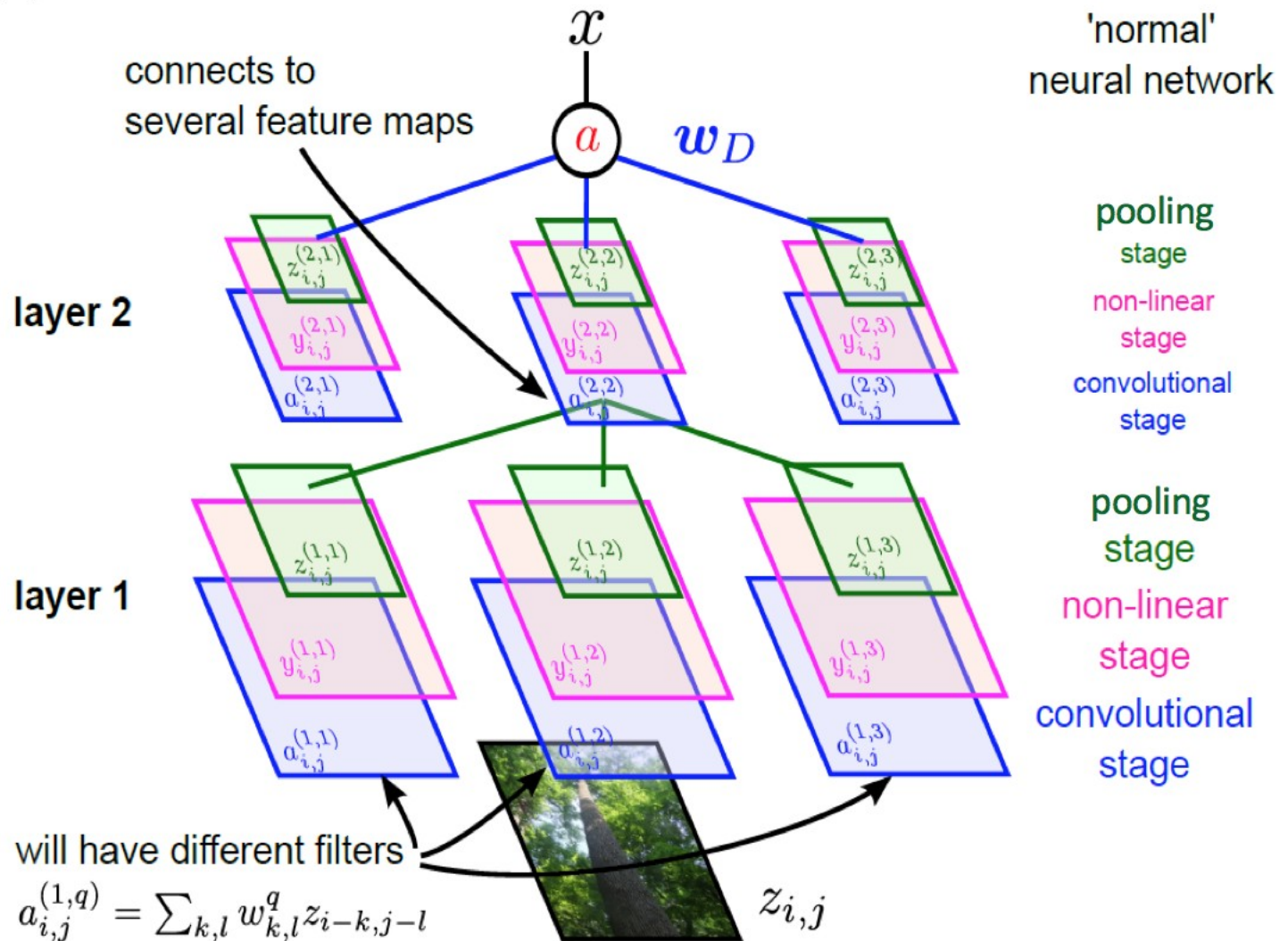


CNN Architecture: Pooling Layer

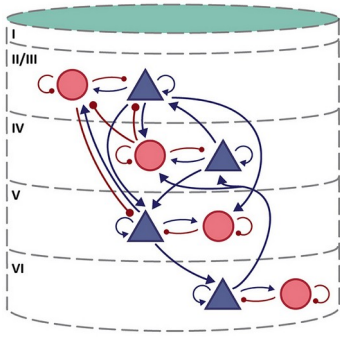
- Pooling partitions the input image into a set of non-overlapping rectangles and, for each such sub-region, outputs the maximum value of the features in that region.
- Intuition: to progressively **reduce the spatial size** of the representation to **reduce the amount of parameters and computation** in the network, and hence to also **control overfitting**



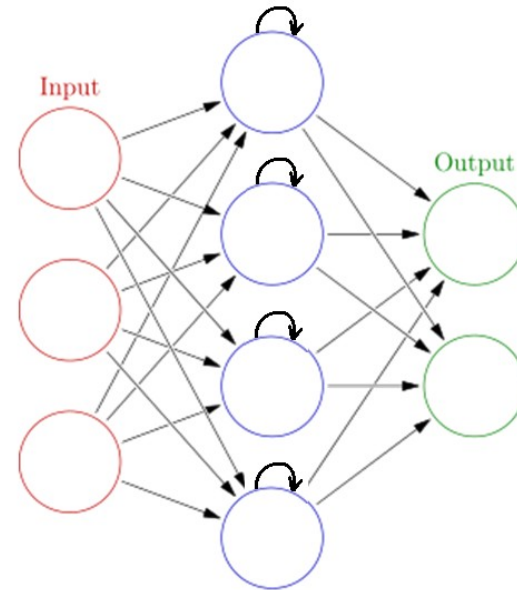
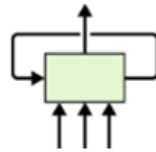
Full CNN



Recurrent Neural Networks and LSTM neurons

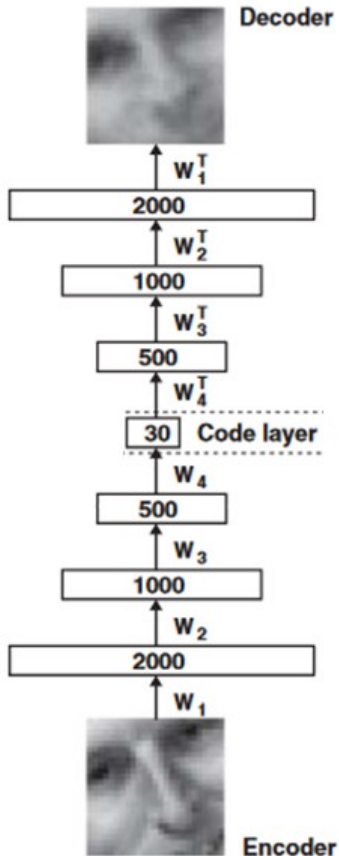


Potjans and Diesmann (2014)



Note: No top-down feedback connections from top layers

Autoencoder

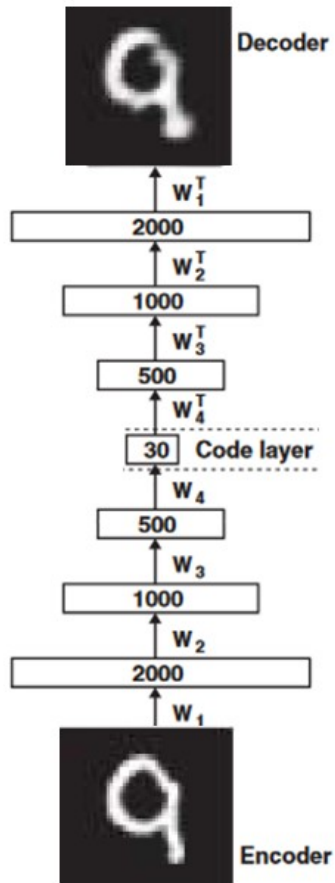


- Train the neural network to reproduce its input vector as its output
- This forces it to compress as much information as possible into few numbers in the central bottleneck.
- These few (here 30) numbers are then a good way to represent data.

Autoencoder

Reducing the Dimensionality of Data with Neural Networks

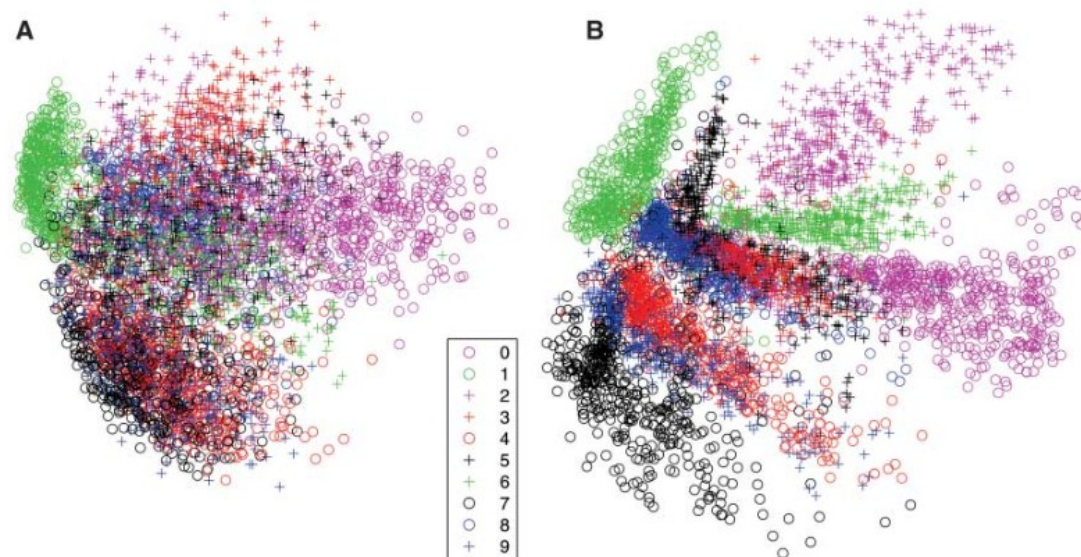
G. E. Hinton* and R. R. Salakhutdinov 28 JULY 2006 VOL 313 SCIENCE



```

0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9
    
```

Fig. 3. (A) The two-dimensional codes for 500 digits of each class produced by taking the first two principal components of all 60,000 training images. (B) The two-dimensional codes found by a 784-1000-500-250-2 autoencoder. For an alternative visualization, see (8).



Convolutional Autoencoder

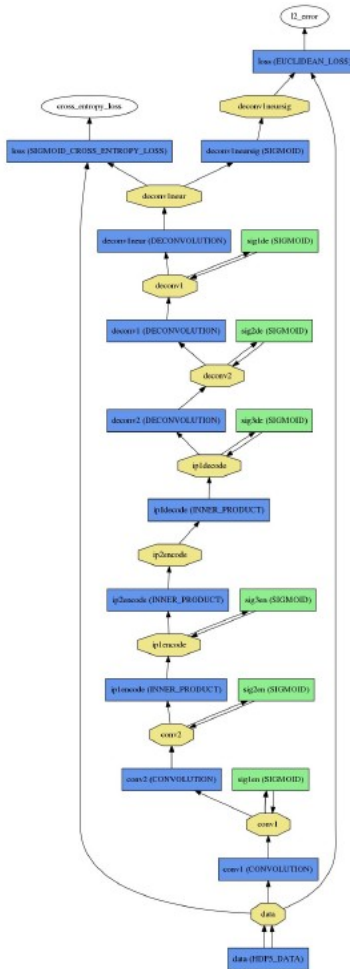


Figure 3. CAE model in Caffe

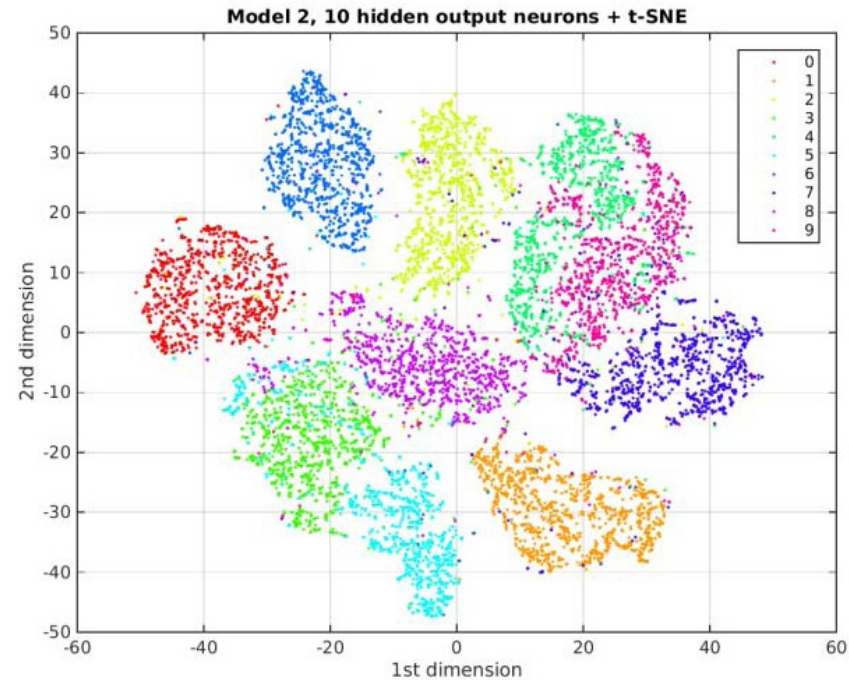
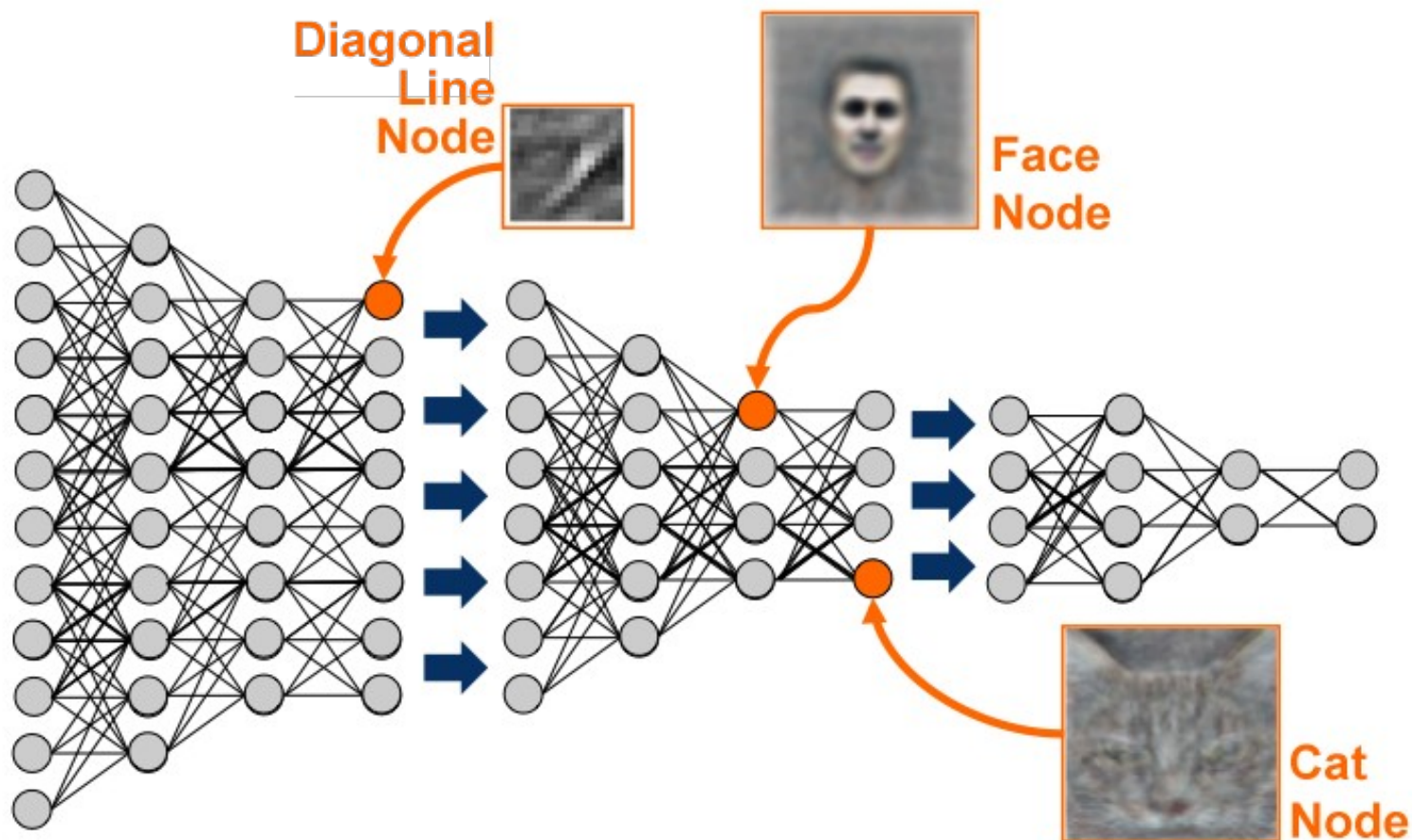


Figure 5. Visualization of MNIST test set in a 2D space by 10-dimensional CAE Model 2 + t-SNE

Turchenko & Luczak, IEEE IDAACS 2017

Deep Neuronal Networks

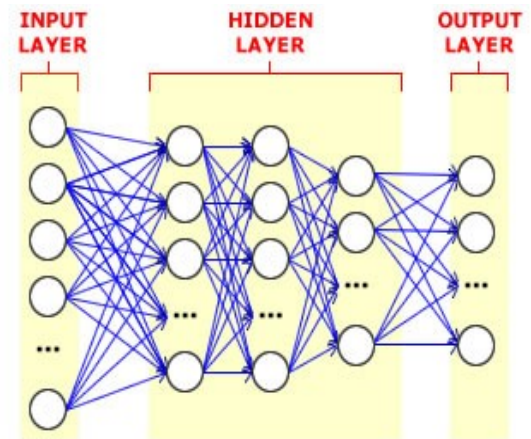


Le et al. (2013) *ICASSP, IEEE International Conference*

<http://theanalyticsstore.ie/deep-learning/>

And that's that

- That's the basic idea
- There are many many types of deep learning,
- different kinds of autoencoder, variations on architectures and training algorithms, etc...
- Very fast growing area ...



Thanks

- Shangsong Liang
- Sun Yat-sen University
- liangshangsong@gmail.com