

Machine Learning

Neural Networks and Deep Learning

Dr. Edgar Acuna
Department of Mathematical Sciences
Universidad de Puerto Rico- Mayaguez
January 2022

academic.uprm.edu/eacuna

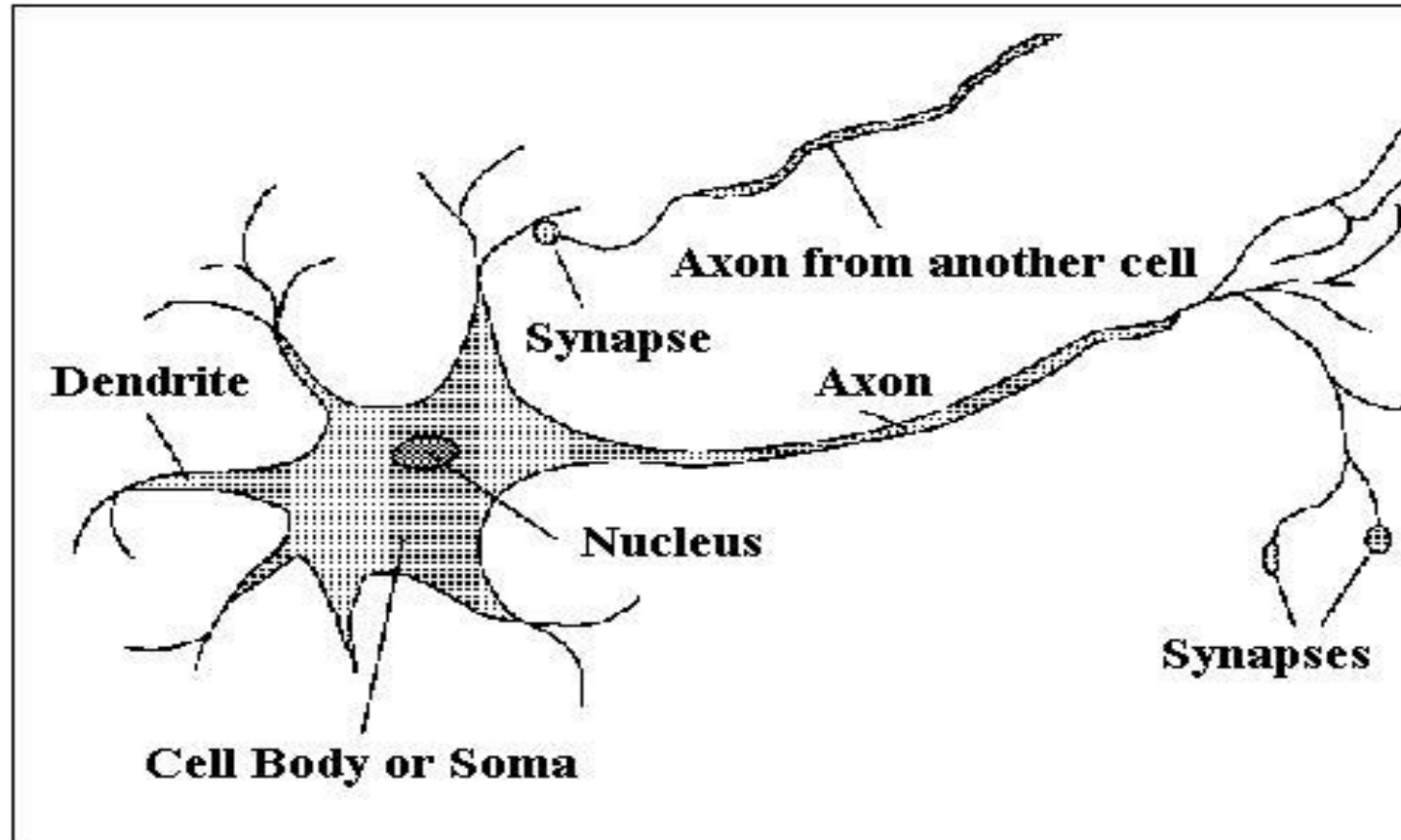
Introduction

The human brain has approximately 1.5×10^{10} neurons with a large number of connections among them called synapses. The number of synapses connecting two neurons varies from 10 to 10^4 . A synapse is a small gap at the end of a neuron that allows a signal to pass from one neuron to the next. The network of neurons form a system of information processing which is massively parallel.

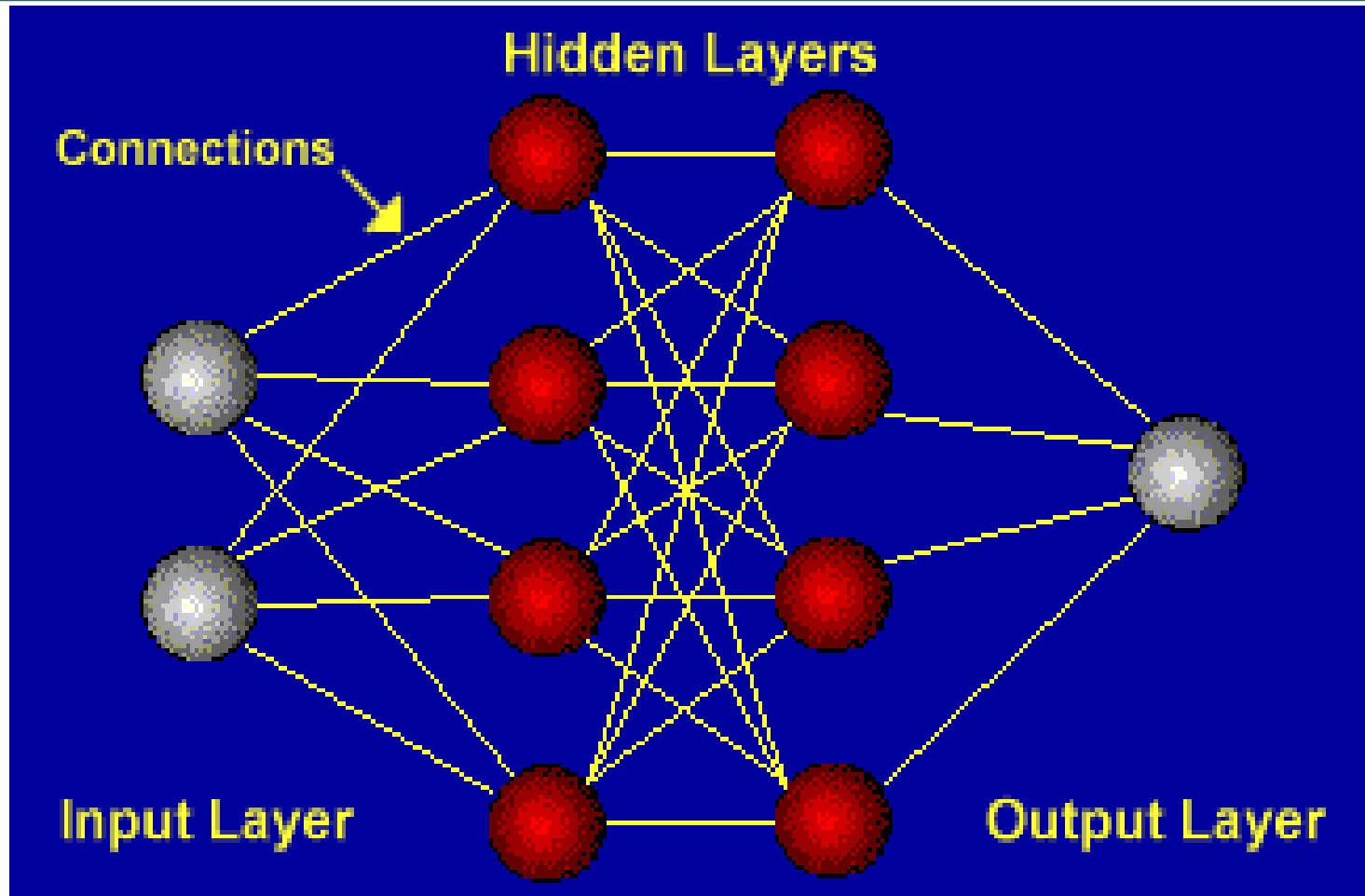
Neural Networks can be considered as an intent to emulate the human brain. In general a neural nets are representations of mathematical models where computational units are connected among them using a mechanism that learns from experience. That is, from the data was recollected.

The first ideas about Neural nets were introduced by McCulloch and Pitts (1943). Later, Rosenblatt (1958, 1962) introduced the concept of “Perceptron” (set of neurons) and he tried to apply them to classification. In 1974, Werbos, published a first gradient descendent algorithm for backpropagation. But, it was until 1986 when, Rumelhart, Hinton and Williams presented the Backpropagation algorithm for learning a neural nets that these started to be used in great demand.

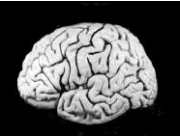

A neuron



A neural network



The Human Brain vs. Computers

	processing elements	element size	energy use	processing speed	style of computation	learns	Intelligent, conscious
	10^{15} synapses	2×10^{-8} m	30 W	1000 Hz	parallel, distributed	yes	usually
	2×10^{10} transistors	10^{-8} m	130 W (Core I7)	3×10^9 Hz	serial, centralized	It is trying	Not(yet)

Estimates are that computers will surpass the capability of human brains around the year 2040 (R Kurzweil, Google).

In Statistics, the research work done by Ripley (1993) and, Chen and Titterigton (1994) were fundamental to attract statisticians into neural network.

Artificial Neural Nets (ANN) have plenty of applications. Among them;

- a. Classification
- b. Regression
- c. Clustering
- d. Outlier Detection
- e. Density function estimation

Comparing Statistics and Neural Nets terms

Statistics	Neural networks
Independent variables	Inputs
Dependent values	Targets
Predicted values	Outputs
Estimacion, Ajuste	Learning, training
Parameters	Weights
Transformations	Functional Links
Outlier Detection	Novelty detection

Type of neural nets

a) For supervised Learning (Linear regression and classification)

- Multilayer Perceptron (MLP)
- Radial basis function Networks (RBF)
- Learning Vector Quantization (LVQ)

b) For unsupervised Learning (Clustering)

- Hopfield Networks
- Kohonen's Self-Organizing Maps
- Adaptive Resonance Theory

c) For time series: Recurrent Neural Networks (RNN) and Long short term memory (LSTM)

The single-layer perceptron

A perceptron is a single layer neural network. A perceptron calculates a linear combination of inputs adding also an intercept called Bias. Then, an **activation function**, which in general is nonlinear is applied to the linear combination to generate a output. Thus, the output y_j is given by

$$y_j = f_j \left(\sum_{inputs:i} w_{ij} x_i \right)$$

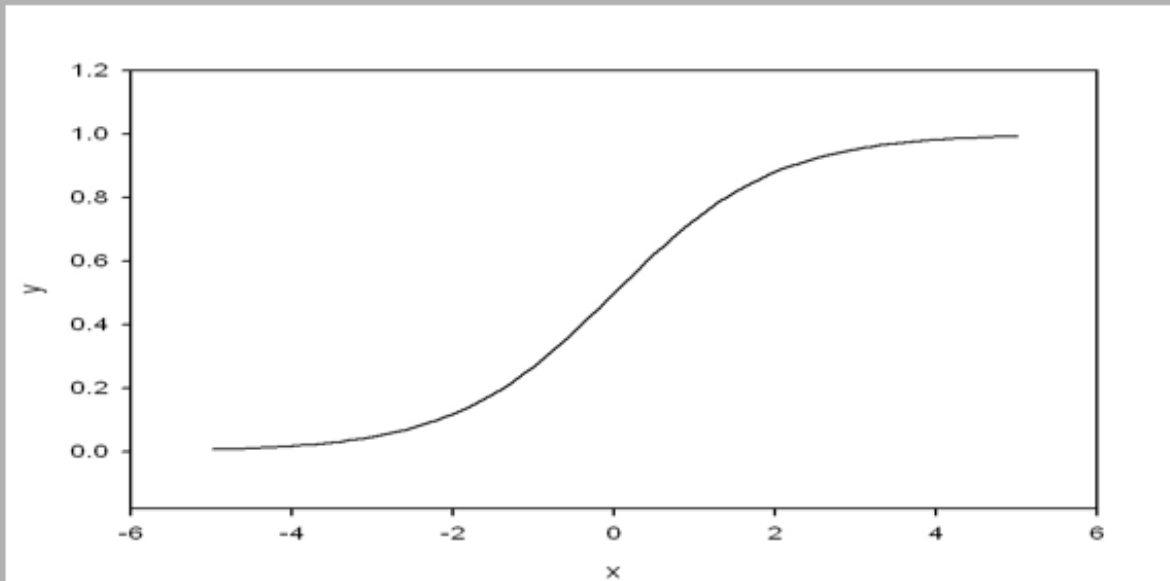
f_j represents the activation function and w_{ij} are the weights. The neural net learns the weights from the recollected data.

Activation Functions

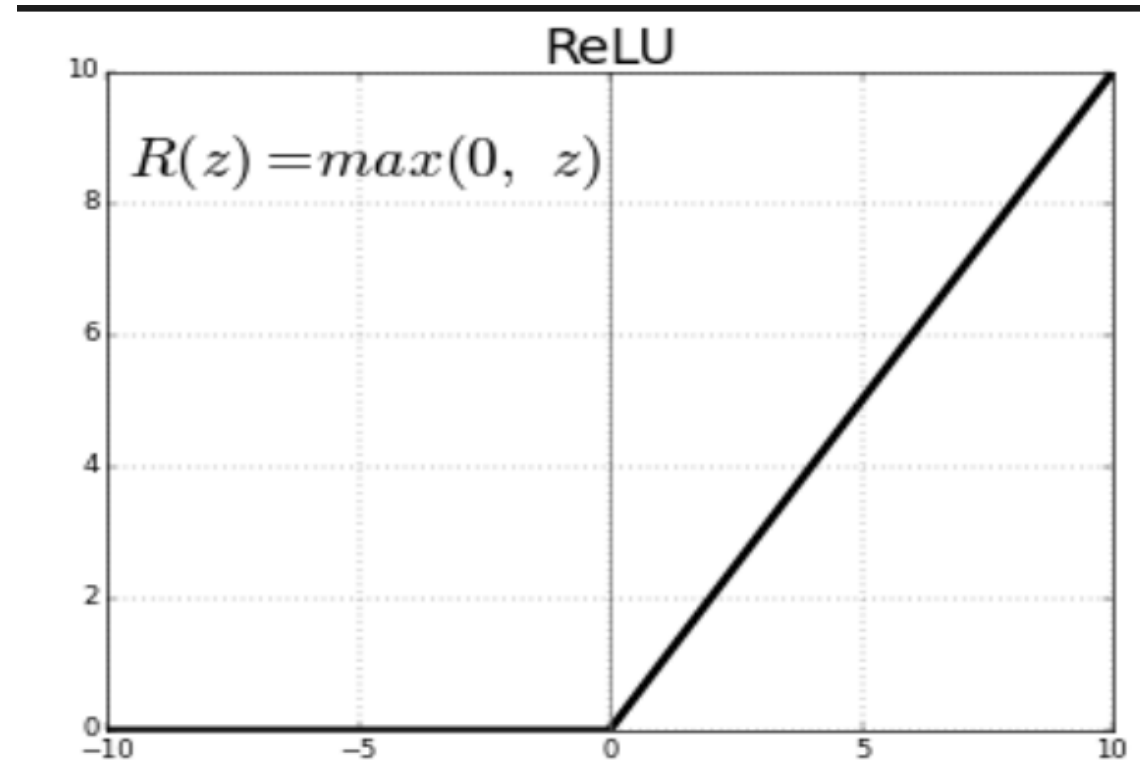
Name	Function
Linear	$A(x)=x$
Logistic (Sigmoid)	$A(x)=(1+e^{-x})^{-1}$
Gaussian	$A(x)=\exp(-x^2/2)$
Relu	$A(x)=0$ si $x<0$, $A(x)=x$ en otro caso
Threshold	$A(x)=0$ si $x<0$, $A(x)=1$ en otro caso
Tanh	$A(x)=2*\text{logistic}(2x)-1$

The Logistic (sigmoid) activation function

$$y = \frac{1}{1 + \exp(-x)}$$

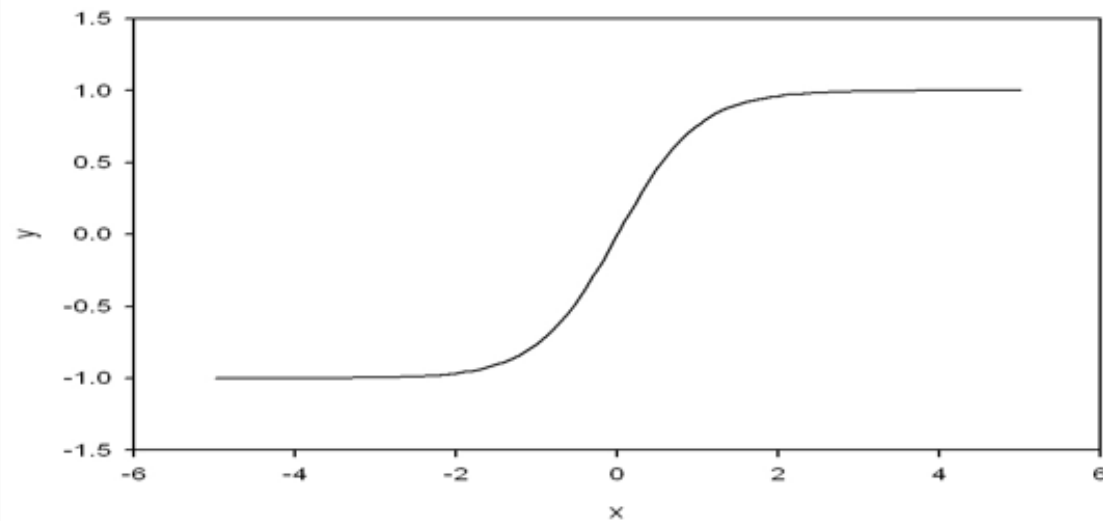


The Relu activation function

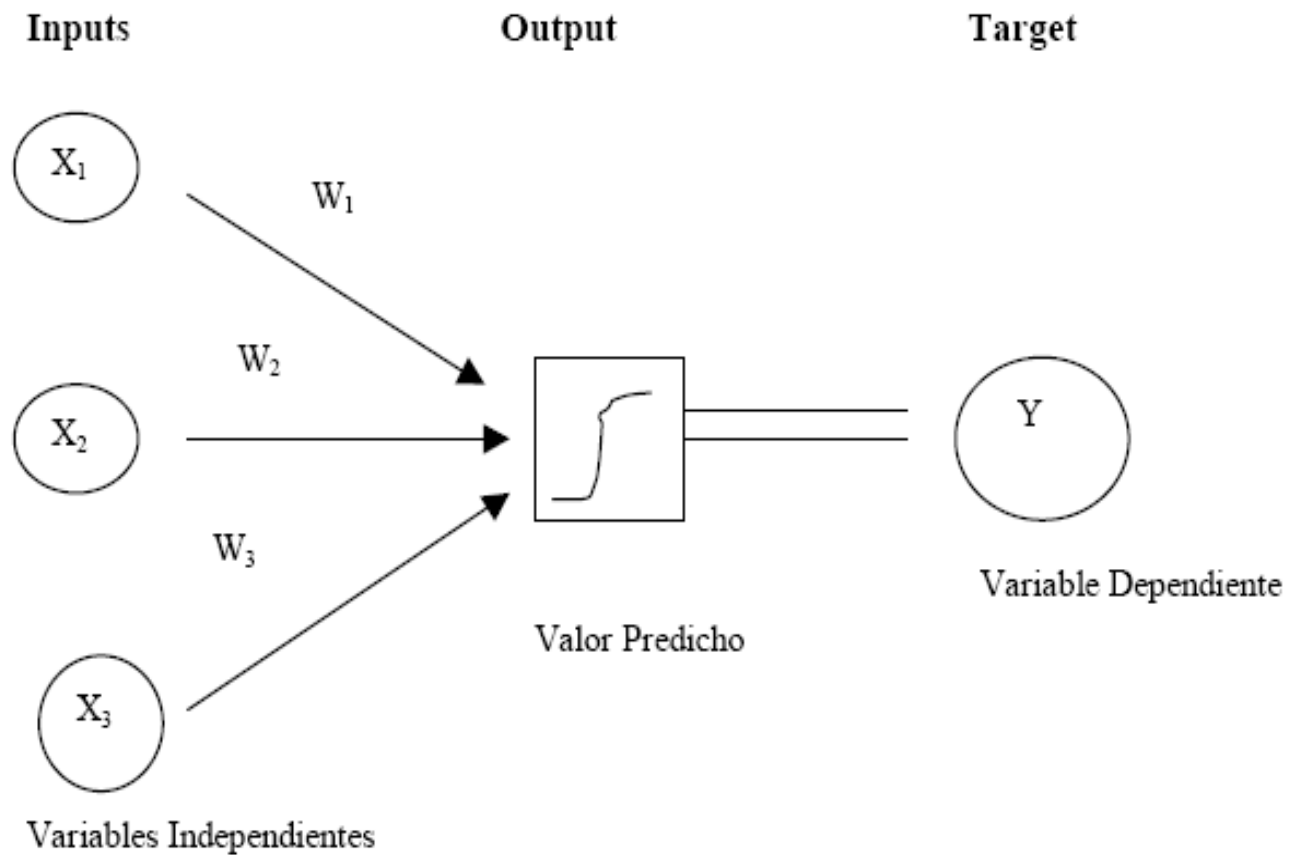


The Tanh activation function

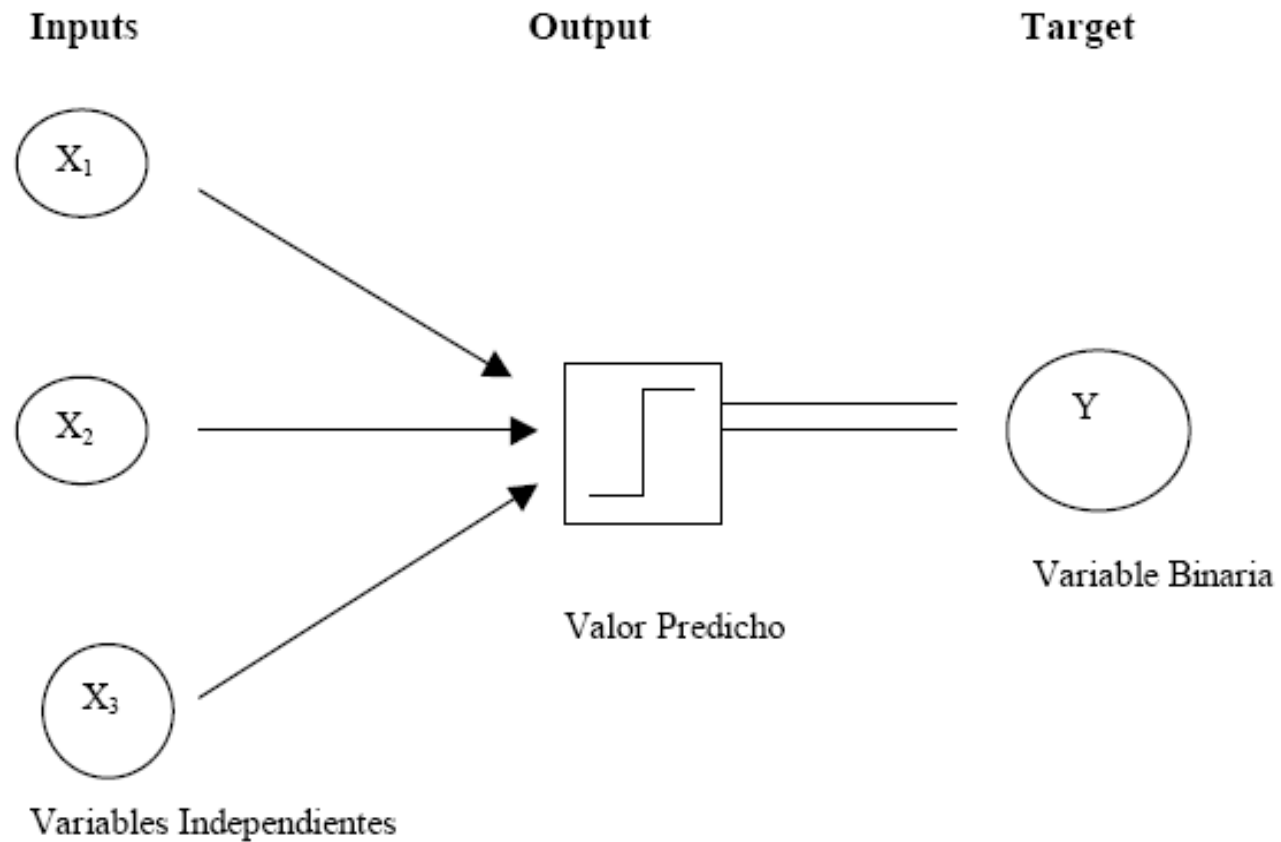
$$y = \tanh(x) = \frac{\sinh(x)}{\cosh(x)} = \frac{\frac{\exp(x) - \exp(-x)}{2}}{\frac{\exp(x) + \exp(-x)}{2}}$$



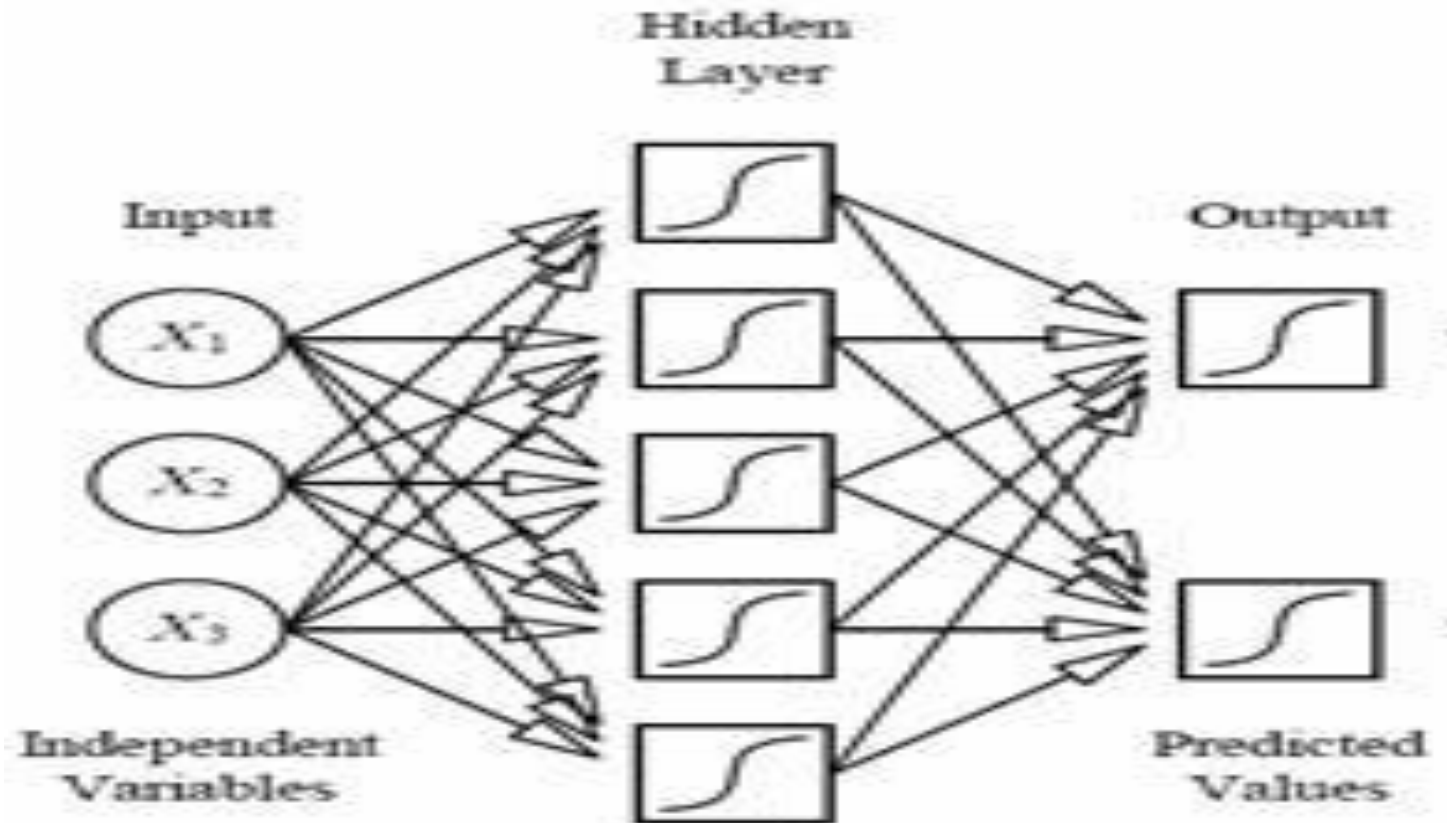
Single Perceptron= Logistic Regression



Perceptron for two-class classification



A MLP Neural Net



The Multilayer Perceptron, MLP

A neural net of two layers can be written as a pair of equations

$$z_j = \phi_h(\alpha_j + \sum_i w_{ji} x_i)$$

$$y_k = \phi_o(\alpha_k + \sum_j w_{jk} z_j)$$

Here z_j are the units in the hidden layer, y_k are the units in the output layer and ϕ_o and ϕ_h are activation functions. Usually ϕ_h is the logistic function, also known as the sigmoidal activation function. On the hand, ϕ_o can be either linear, logistic or threshold function. Replacing the first equation in the second equation, we obtain

$$y_k = \phi_o(\alpha_k + \sum_j w_{jk} \phi_h(\alpha_j + \sum_i w_{ji} x_i))$$

where the weights w_{ij} and the biases α_j must be estimated using the data

The Multilayer Perceptron, MLP(cont)

In classification problems with more than two classes, say C classes, the softmax function is used as an output function. The softmax function, ϕ_o , is defined as

$$\phi_o(z) = (\phi_1(z), \dots, \phi_k(z), \dots, \phi_C(z))$$

where

$$\phi_k(z) = \frac{e^{t_k z}}{\sum_{j=1}^C e^{t_j z}}$$

The MLP also is called a Feed Forward Neural Network (FFNN) or a Backpropagation net

When a neural network has more than one hidden layers is called a “deep” Neural network.

Training of a Neural Net

- This is equivalent to model estimation in the statistics world. The weights w 's of the NN are chosen in such way that some fit measure is minimized.
- **Regression:** The sum squares of errors given by

$$E = \sum_{i=1}^n \sum_{j=1}^J (\hat{y}_j^i - y_j^i)^2$$

Is minimized with respect to the weights $w=(\alpha_j, w_{ij})$

Training of a neural network(cont)

- **Classification:** In two-classes problems, the cross entropy given by

$$C = \sum_{i=1}^n [y_i \log\left(\frac{y_i}{\hat{y}_i}\right) + (1 - y_i) \log\left(\frac{1 - y_i}{1 - \hat{y}_i}\right)]$$

Is minimized with respect the vector of weights w . If the instance \mathbf{x}_i belongs to class 1 then $y_i = 1$ and it is equal to 0 otherwise. Las \hat{y}_i are estimations of the posterior probability to lie in the class 1 when \mathbf{x}_i is observed.

This is equivalent to parameter estimation in a binary logistic regression.

Training of a neural network(cont)

For J classes the entropy function given by

$$C = \sum_{i=1}^n \sum_{j=1}^J y_j^i \log\left(\frac{y_j^i}{\hat{y}_j^i}\right)$$

must be minimized. C can be minimized iteratively using methods from numerical analysis and nonlinear models including:

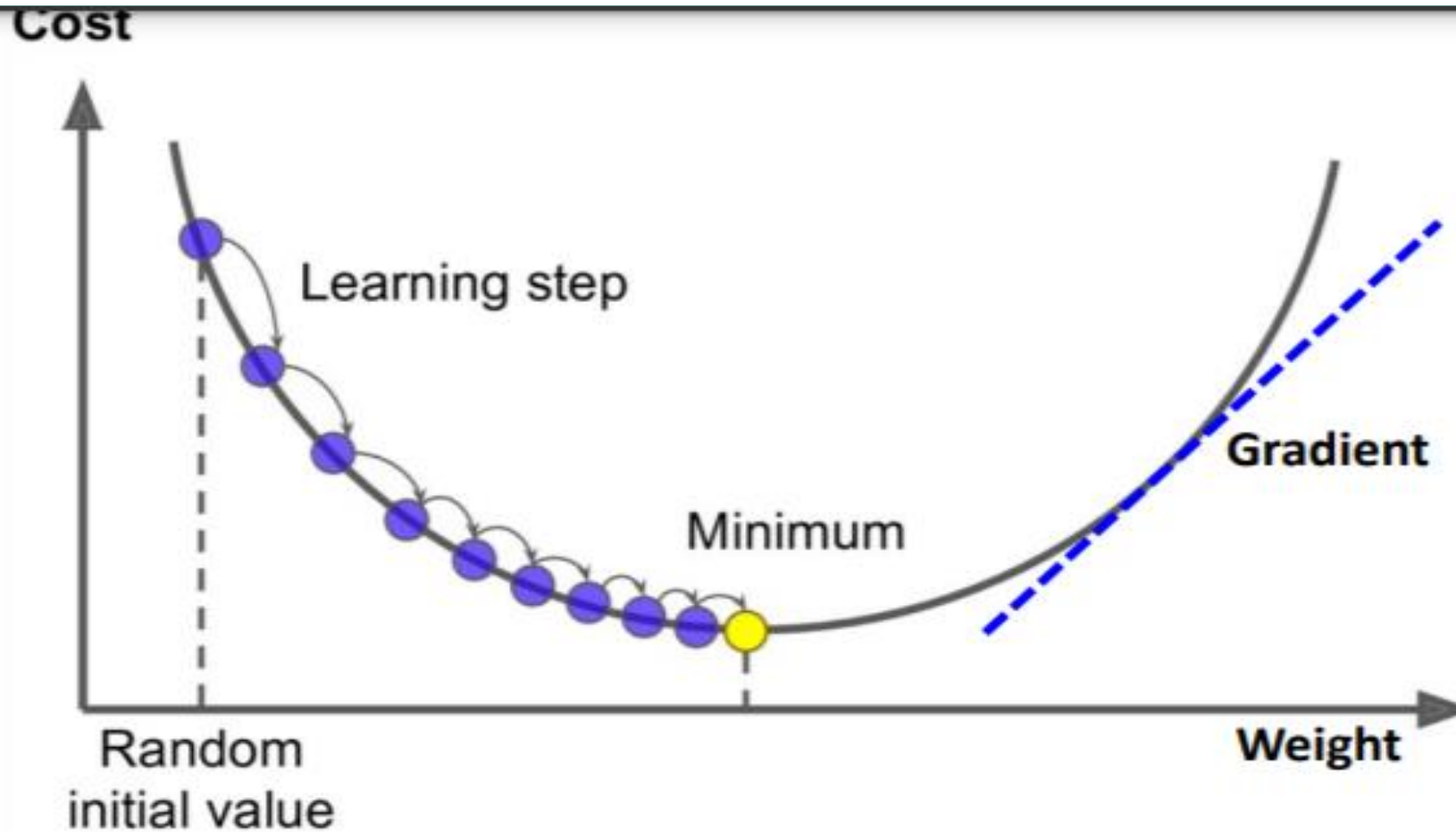
Gradient Descent, Quasi-Newton methods (recommended if the number of weights is less than 1000), Gradient Conjugated method (recommended if there is a large amount of weights to be estimated), Simulated Annealing, Particle Swarm Optimization and Genetic Algorithms.

The greatest challenge to minimize C is the presence of multiple local minimum and there is a risk to choose one that is not the optimum. Frequently, it is necessary to restart the minimization process using different starting values for the iterative process.

The Backpropagation algorithm is used to compute the gradients.

<https://towardsdatascience.com/understanding-backpropagation-algorithm-7bb3aa2f95fd>

The Gradient Descent



The Gradient descent

The updating of the weights at the k-th step is given by:

$$\begin{aligned}w_{k+1} &= w_k - \eta \cdot \nabla_w C \\b_{k+1} &= b_k - \eta \cdot \nabla_b C\end{aligned}$$

For some $\eta > 0$. This tuning parameter is called **Learning Rate**. Choosing the learning rate is challenging as a value too small may result in a long training process that could get stuck, whereas a value too large may result in learning a sub-optimal set of weights too fast or an unstable training process. The gradient for an n -dimensional function $f(x)$ at the point p is

$$\nabla f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_n}(p) \end{bmatrix}$$

The Stochastic Gradient descent (SGD)

Stochastic gradient descent uses this idea to speed up the process of doing gradient descent. Specifically, the input data are randomly partitioned into disjoint groups $M_1, M_2, \dots, M_{n/m}$. We then do the following updates to the weights (biases are done at the same time, but omitted for sake of space):

$$\begin{aligned}w_{k+1} &= w_k - \frac{\eta}{m} \sum_{i \in M_1} \nabla C_i \\w_{k+2} &= w_{k+1} - \frac{\eta}{m} \sum_{i \in M_2} \nabla C_i \\&\vdots \\w_{k+n/m+1} &= w_{k+n/m} - \frac{\eta}{m} \sum_{i \in M_{n/m}} \nabla C_i\end{aligned}$$

Each set M_j is called a **mini-batch** and going through the entire dataset as above is called an **epoch**.

Example 1. Neural Nets applied to the prediction of the final grade based on Ex1 and Ex2

The **MLPClassifier** function from the **scikit-learn** library allow us to perform training of a MLP neural net.

For instance, to predict the final grade, we use the following commands:

```
df=pd.read_csv("http://academic.uprm.edu/eacuna/eje1dis.csv")
y=df['Nota']
X=df.iloc[:,0:2]
#creating a numerical column "pass" to represent the classes
lb_make = LabelEncoder()
df["pass"] = lb_make.fit_transform(df["Nota"])
y2=df['pass']
y1=y2.as_matrix()
X1=X.as_matrix()
```

Example 1 (cont)

```
#Training a neural net with one hidden layer and five units on it
mlp = MLPClassifier(solver='lbfgs',hidden_layer_sizes=(5),max_iter=1000,random_state=99)
mlp.fit(X1, y1)
#Showing the weights
mlp.coefs_
[array([[ -1.93306393e+01, -2.19847626e-02, 4.67280387e+01, -1.21251200e+01, -
 1.67606515e+01], [ 9.31589043e+01, -3.73206703e-01, 4.15846672e+01, -
 1.23368759e+01, -3.21037745e+00]]), array([[ -7.01419931e-04], [ 9.44058482e-01], [
 1.10475826e-02], [ 1.96857243e+00], [ -1.35496680e+00]])]
#Showing the biases
mlp.intercepts_
[array([-7.44811572, 0.45692388, 1.39966161, -0.16317004, 0.57152039]), array([-
 51.31888278])]
There are 21 parameters in the model
```

Example 1 (cont)

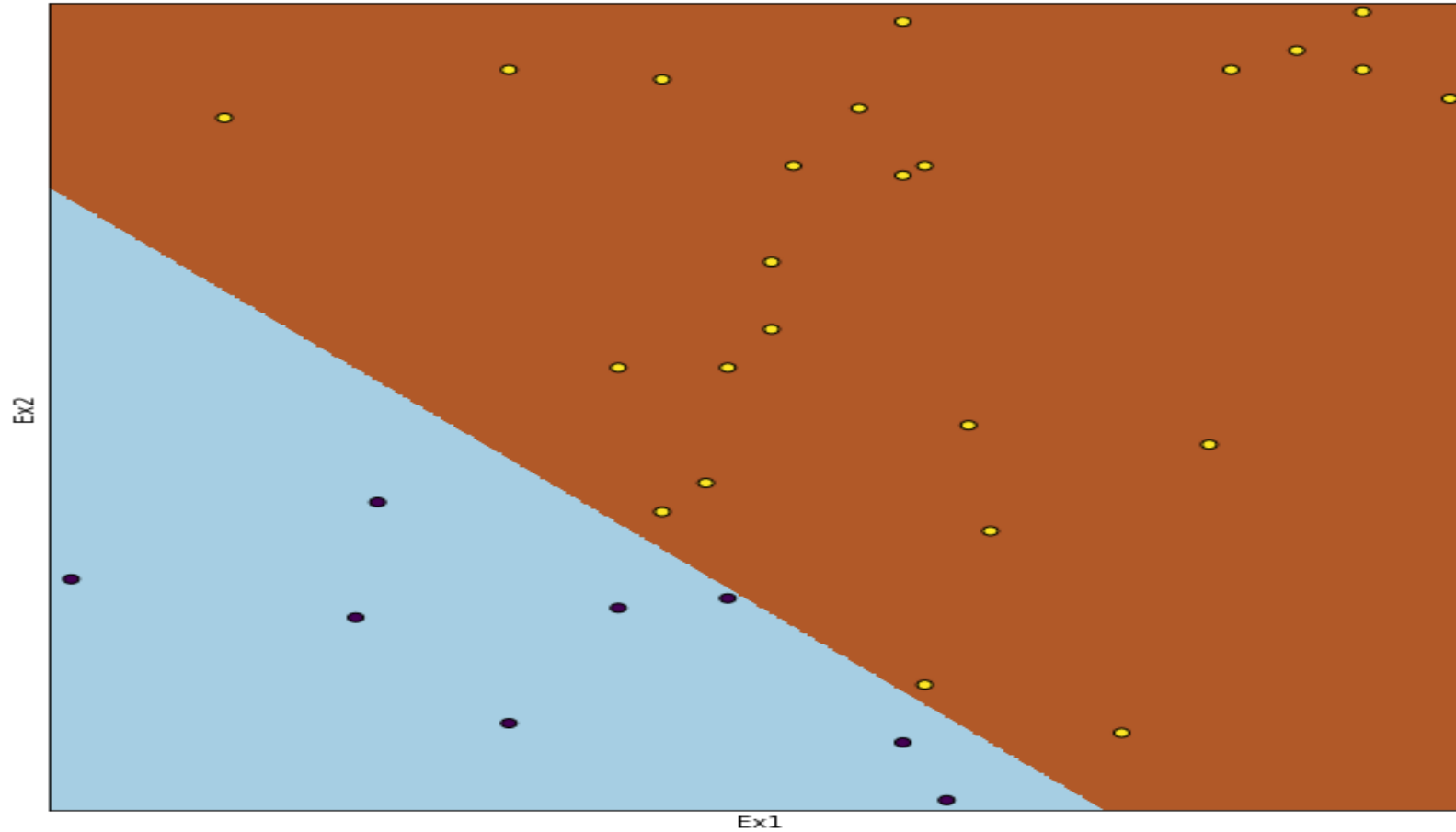
#Calculating the posterior probabilities

```
mlp.predict_proba(X1)
```

```
array([[ 0.00000000e+00,  1.00000000e+00],  
       [ 2.22044605e-16,  1.00000000e+00],  
       [ 0.00000000e+00,  1.00000000e+00],  
       [ 2.22044605e-16,  1.00000000e+00],  
       [ 3.10862447e-15,  1.00000000e+00],  
       [ 6.56805055e-10,  9.99999999e-01],  
       [ 1.19992905e-12,  1.00000000e+00],  
       [ 4.33475522e-09,  9.99999996e-01],  
       [ 2.44449664e-08,  9.99999976e-01],  
       .....])
```

A record of X is assigned to class with the greatest posterior probability

Example 1: Decision Boundary using Neural Nets



Example 2:Diabetes

```
url= "http://academic.uprm.edu/eacuna/diabetes.dat"
names = ['preg', 'plas', 'pres', 'skin', 'test', 'mass', 'pedi', 'age', 'class']
data = pd.read_table(url, names=names)
y=data['class']
X=data.iloc[:,0:8]
y1=y.as_matrix()
X1=X.as_matrix()
#Training a NN with one hidden layer and 20 units
mlp2=MLPClassifier(solver='lbfgs',hidden_layer_sizes=(20),max_iter=1000)
mlp2.fit(X1, y1)
#Estimating the accuracy
mlp2.score(X1, y1)
0.7942708333333337
```

Example 2(cont)

```
#training a NN with two hidden layer and 20 units in each of them
mlp22=MLPClassifier(solver='lbfgs',hidden_layer_sizes=(20,20),max_iter=5000)
mlp22.fit(X1, y1)
mlp22.score(X1, y1)
0.82682291666666663
#Using Training and Test sets
X_train, X_test, y_train, y_test = train_test_split(X, y)
#training a NN with one hidden layer and 20 units
mlp=MLPClassifier(hidden_layer_sizes=(20),max_iter=500)
mlp.fit(X_train, y_train)
pred=mlp.predict(X_test)
mlp.score(X_test, y_test)
0.64583333333333337
```

Accuracy of MLP estimated by CV

```
#Estimating the accuracy using cross validation
from sklearn.model_selection import cross_val_score
scores = cross_val_score(mlp, X1, y1, cv=10)
print 'The accuracy estimated by CV is:', scores.mean()
```

The accuracy estimated by CV is: 0.690105946685

Example of overfitting with nnet

Let us consider the Diabetes dataset

Number of units in the hidden layer	5	20	50	100	200	500
Accuracy by resubstitution	68.48	79.29	85.80	95.57	99.47	100.0
Accuracy by cross-validation	72.4	74.6	73.04	72..00		

There is high variability on the estimation of the accuracy by resubstitution

Remedial measures for overfitting

Neural Nets tend to overfit the data. That is, NN tend to give accuracy of 100%.

The following are some remedial measures to fix this problem:

a) Stop the iterative process to estimate the minimum of E. It is assumed that a validation set is available and the iterative process is stopped when the performance of the neural net on the validation set begins to deteriorate.

b) Regularization: In this case a penalty is added to the function $E(w)$ and then it is minimized. More specifically,

$$\text{Min}_w[E(w) + \lambda\phi(w)].$$

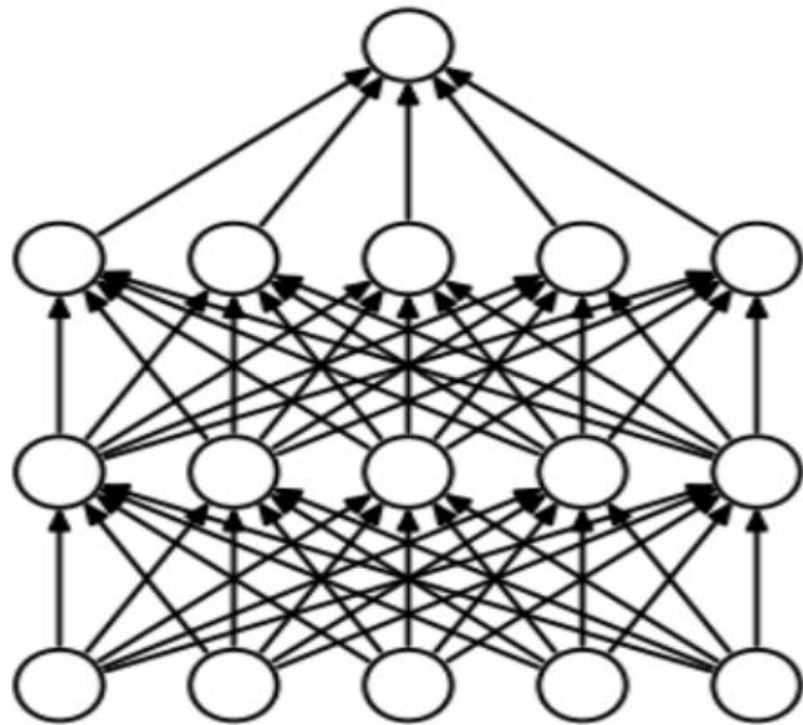
here λ is the regularization constant and ϕ is the penalty function of the model. The most simple regularization method is the one known as weight decay defined by

$$\text{Min}_w[E(w) + \lambda\phi\sum w^2]$$

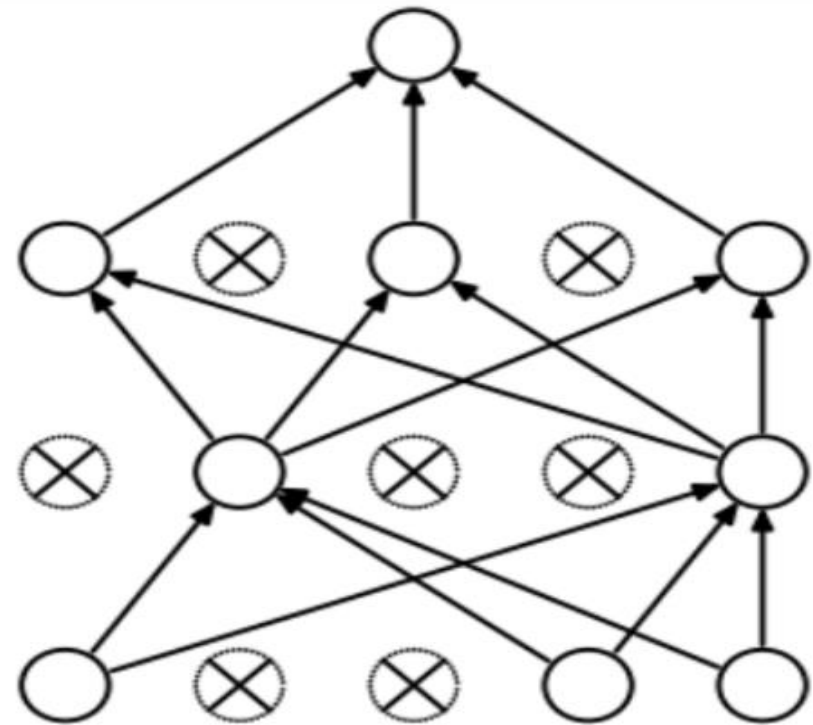
This is very similar to ridge regression in statistics.

Remedial measures for overfitting(cont)

- c) **Average:** In this case several values are chosen as starting values of the iterative process of minimizing $E(w)$ and then the average of the predictions obtained is taken.
- d) **Add noise.** In this case we add noise to each input variable and then a NN is fitted. The process is repeated several times and then the predictions obtained are averaged.
- e) **Bagging** (“Bootstrap aggregating”). Several samples with replacement and of the same size are taken from the training sample. A NN is fitted for each of these samples. Finally, each instance of the training sample is assigned to the most voted class.
- f) **Dropout** is a technique where randomly selected neurons are ignored during training. They are “dropped-out” randomly. This means that their contribution to the activation of downstream neurons is temporally removed on the forward pass and any weight updates are not applied to the neuron on the backward pass.



(a) Standard Neural Net



(b) After applying dropout.

Dropout: Srivastava, et al. JMLR, 2014.

Effect of the use of decay

Let us consider the Diabetes dataset. The parameter alpha of the MLPClassifier determines the decay. Only one hidden layer is used with 20 units.

Decay	0	5	10	.1	.5
Resubstitution Error	80.59	84.11	81.51	78.51	80.46

A large decay gives higher accuracy

Radial basis functions Nets

Initially, the radial basis functions nets (RBF nets) were applied to approximate functions.

Roomhead and Lowe (1988) were the first researchers to use RBF nets in classification.

Mathematically, the RBF nets can be written as a linear combination of nonlinear radially symmetric functions. That is:

$$y = \alpha + \sum_{j=1}^M \beta_j \phi_j(\|x - c_j\|)$$

where c_j are pre-specified centers, α and β_j are weights to be estimated, the ϕ_j 's are the basis functions, usually they are the same for all j . The most used functions are the gaussian $\phi(r)=\exp(-r^2/2\sigma)$, the quadratic $\phi(r)=r^2+ar+b$ and the thin plate spline $\phi(r)=r^2\log(r)$. Here r represents radius.

Radial basis functions Nets(cont)

RBF nets are very similar to density estimation using Gaussian mixtures or kernel methods.

Most of the time, the centers of the model, are the centroids determined for any clustering algorithm, such as k-means.

The number of components, M , is chosen trying to avoid either underfitting or overfitting.

In terms of classification, RBF nets are used to estimate the posterior probability of each class, $P(C_j/x)$ and assign then the object x to the class j with the maximum posterior probability.

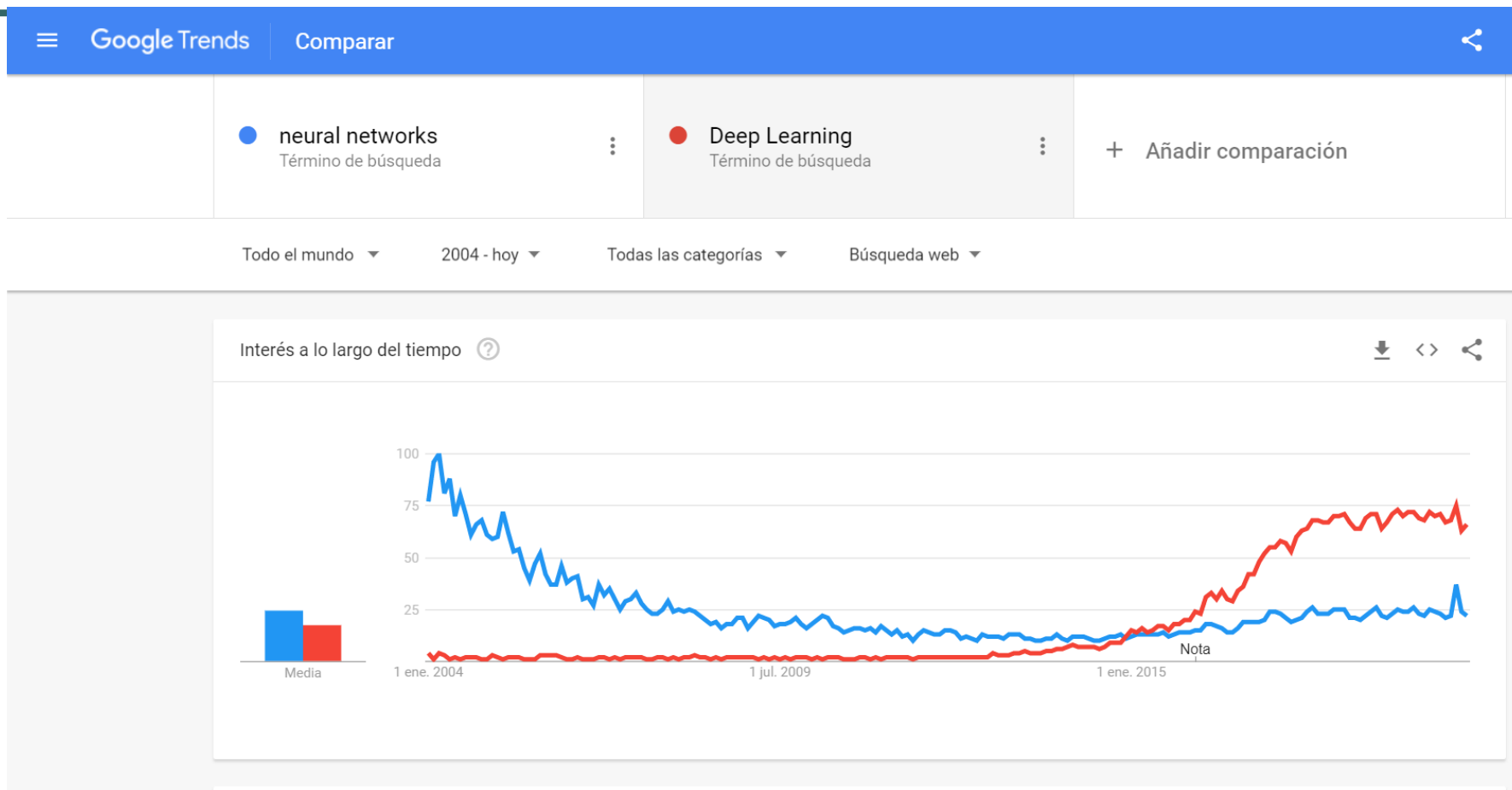
In scikit-learn, the function `sklearn.gaussian_process.kernels.RBF` trains a RBF net.

Radial basis functions Nets(cont)

- The training of RBF net is faster than the training MLP net. Also , an RBF is easier to be interpreted. However the MLP nets is much faster for making predictions.
- RBF nets are more sensitive than MLP nets to the curse of dimensionality problem. Therefore, it does not perform well with a large number of predictors.

Advantages and disadvantages of Neural Nets

- a) Neural nets are good for making predictions but they are hard to understand.
- b) NN are good to analyze large and complex datasets.
- c) Depending on the problem NN can perform better or worst than any other machine learning algorithm.
- d) NN do not have good theory for feature selection or model selection in general.



Deep Learning

In 2006, Geoffrey Hinton (now at Google) showed that a type of neural network called deep belief network, can be efficiently trained using a strategy called greedy layer-wise pretraining.

Hinton and his associates introduced the term “deep learning” to emphasize that now it would be possible to train neural nets more deeply than before.

At the present is claimed that deep neural networks have better performance than any other machine learning methods.

The basic idea of Deep Learning is to simulate the array of a large number of neurons of the Brain's Neocortex in a artificial neural network.

Deep Learning (cont)

Thanks to improvements in the mathematical formulation and the constant increase in computer power, now scientist can model more layers and neurons than before.

In June 2012, a Google's Deep Learning system was capable to identify cats in 10 millions of Youtube videos with the double of efficiency of previous algorithms of image recognition (the team was led for A. Ng and J. Dean). Also, Google has used Deep Learning to reduce the error rate of its algorithm for voice recognition on its androids cell phones.

Also, in September 2012, AlexNet, a CNN consisting of 8 layers wain the ImageNet competition

Deep Learning (cont)

The success of deep learning has led to the appearance of several libraries to train deep neural networks. Most of these libraries are written in Python. For instance, Keras, Theano, H2O and Tensorflow (Google Brain). Now, Tensorflow 2 includes Keras (tf.keras API).

To extend deep learning to applications beyond voice and image recognition it requires more software development (CUDA, 2007) and processing power (GPU).

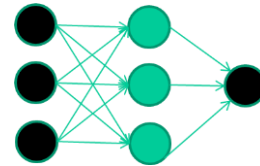
Limitations of Neural Networks

Random initialization + densely connected networks lead to:

- High cost
 - Each neuron in the neural network can be considered as a logistic regression.
 - Training the entire neural network is to train all the interconnected logistic regressions.
- Difficult to train as the number of hidden layers increases
 - Recall that logistic regression is trained by gradient descent.
 - In backpropagation, gradient is progressively getting more dilute. That is, below top layers, the correction signal δ_n is minimal.
- Stuck in local optima
 - The objective function of the neural network is usually not convex.
 - The random initialization does not guarantee starting from the proximity of global optima.
- Solution:
 - Deep Learning/Learning multiple levels of representation

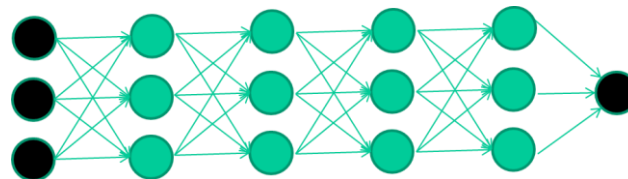
What is new in Deep Learning with respect to NN?

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

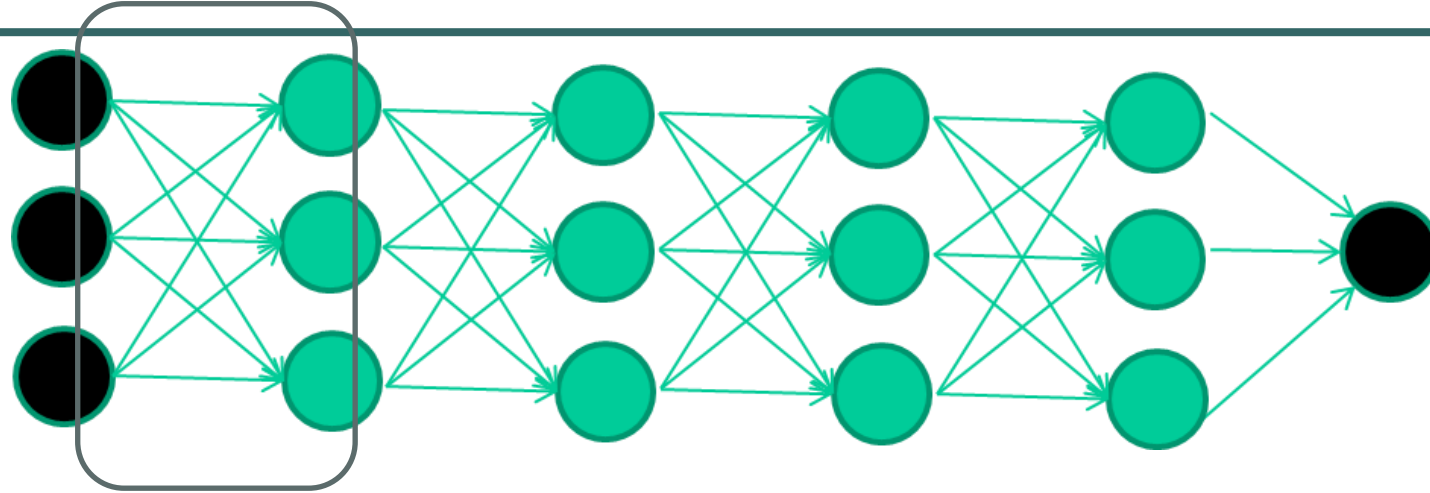


but these algorithms are not good at learning the weights for networks with more hidden layers

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

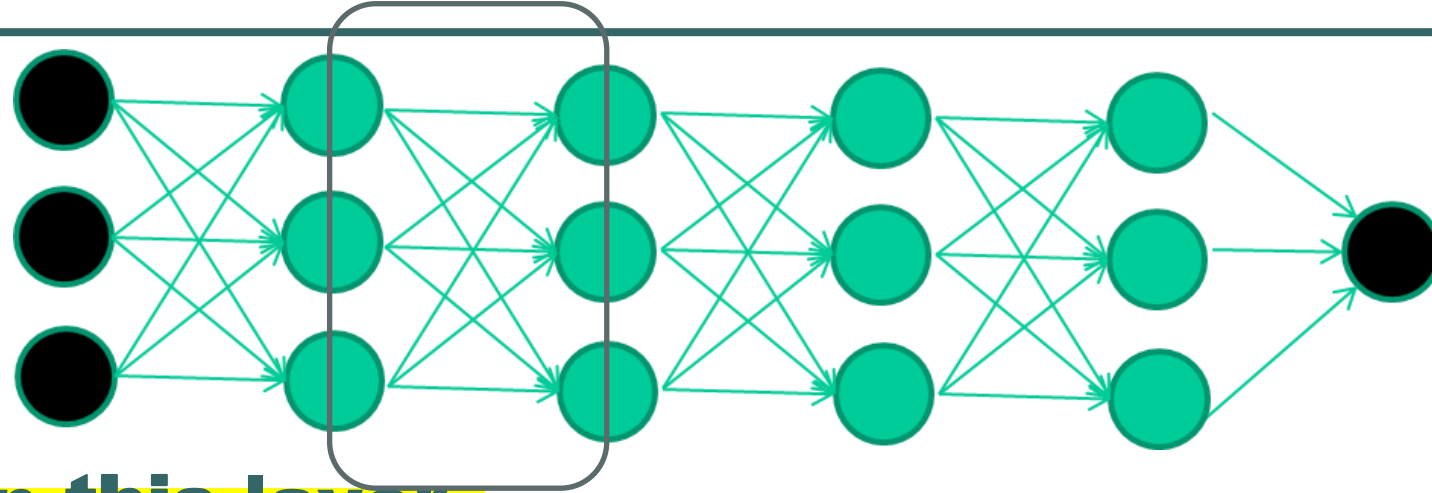


The new way to train multi-layer NNs...



**Train this layer
first**

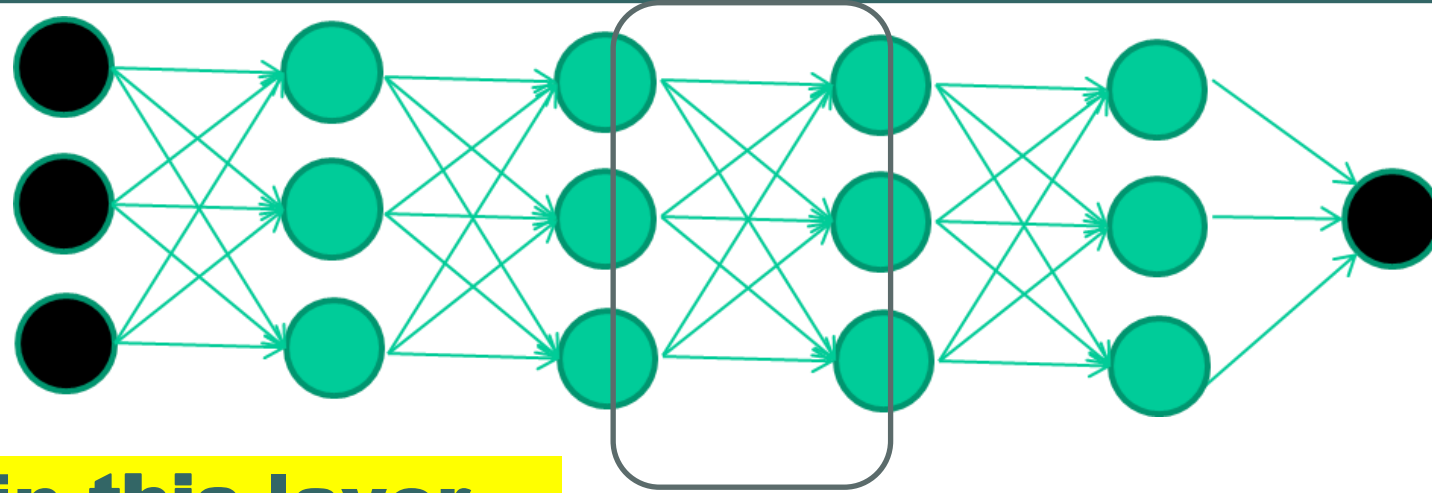
The new way to train multi-layer NNs...



**Train this layer
first**

then this layer

The new way to train multi-layer NNs...

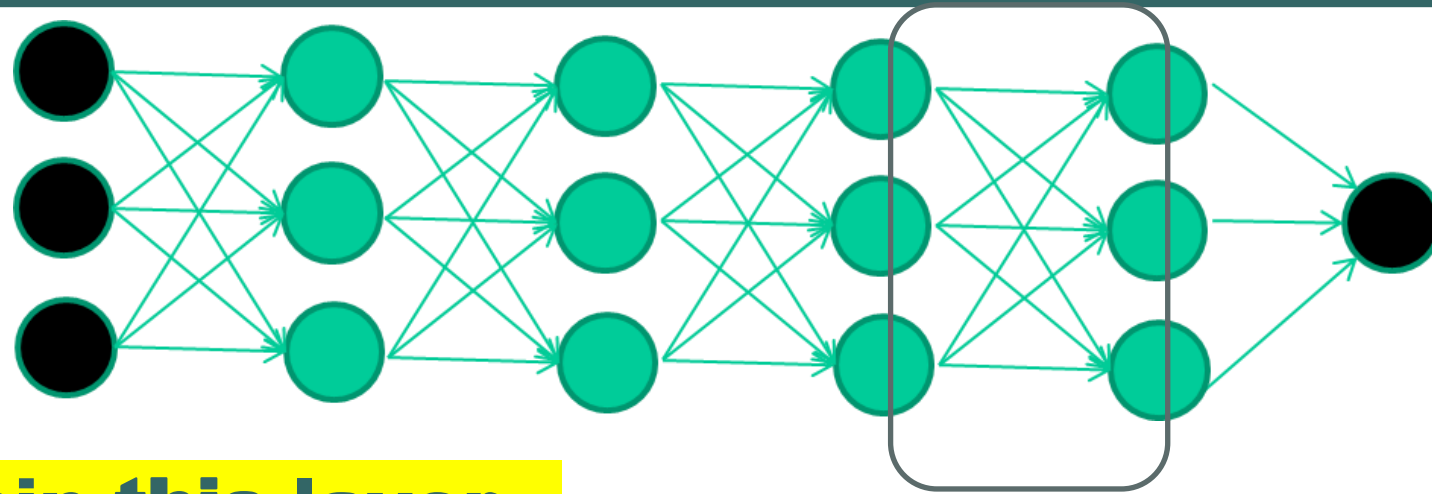


**Train this layer
first**

then this layer

then this layer

The new way to train multi-layer NNs...



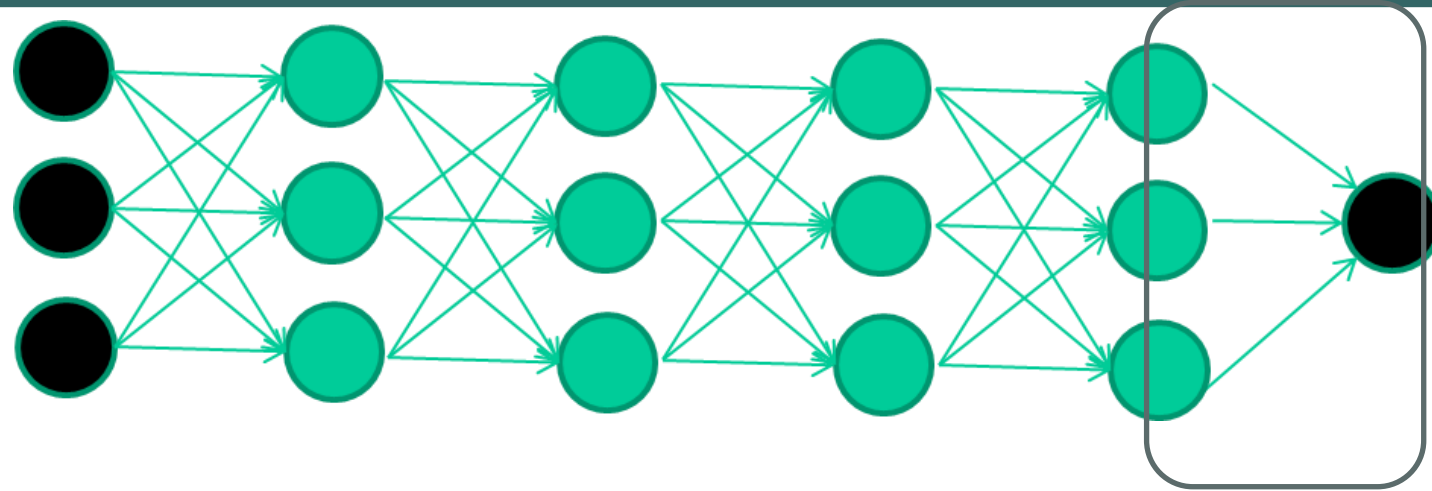
**Train this layer
first**

then this layer

then this layer

then this layer

The new way to train multi-layer NNs...



**Train this layer
first**

then this layer

then this layer

then this layer

finally this layer

Types of Deep Learning architectures

Multilayer Perceptrons (MLPs), useful for numerical data
Convolutional Neural Networks (CNNs), very popular to analyze data from images (Computer Vision).

Recurrent Neural Networks (RNNs), used for learning from sequential data including text, audio and video. Long Short Term Memory networks – usually just called “LSTMs” – are a special kind of RNN, capable of learning long-term dependencies.

Autoencoder is an unsupervised artificial neural network that learns how to efficiently compress and encode data

Deep Learning using h2o

Only performs deep learning using MLP

```
> h2o.init()
```

```
dl_fit3 <- h2o.deeplearning(x = x, y = y, training_frame = data, epochs = 20, hidden=  
c(10,10), seed = 1)
```

```
dl_perf3 <- h2o.performance(model = dl_fit3, newdata = data)
```

For Convolutional Networks and Recurrent Neural Networks, use Deep Water, Theano, Tensorflow, Lasagne, Keras, Mxnet. Caffe,

Deep Learning using h2o

H2OBinomialMetrics: deeplearning

MSE: 0.14129

RMSE: 0.37588

LogLoss: 0.42824

Mean Per-Class Error: 0.21228

AUC: 0.87013

Gini: 0.74025

Confusion Matrix (vertical: actual; across: predicted) for F1-optimal threshold:

	1	2	Error	Rate
1	381	119	0.238000	=119/500
2	50	218	0.186567	=50/268
Totals	431	337	0.220052	=169/768

Deep Learning using h2o

Maximum Metrics: Maximum metrics at their respective thresholds

	metric	threshold	value	idx
1	max f1	0.281453	0.720661	229
2	max f2	0.088843	0.825949	322
3	max f0point5	0.544483	0.735849	137
4	max accuracy	0.488910	0.800781	157
5	max precision	0.986791	1.000000	0
6	max recall	0.031608	1.000000	364
7	max specificity	0.986791	1.000000	0
8	max absolute_mcc	0.281453	0.552705	229
9	max min_per_class_accuracy	0.309897	0.782000	218
10	max mean_per_class_accuracy	0.281453	0.787716	229

Deep Learning using h2o

```
library(h2o)
> h2o.init(nthreads = -1, #Number of threads -1 means use all cores on your machine
+         max_mem_size = "8G") #max mem size is the maximum memory to allocate to
H2O
dl_fit2 <- h2o.deeplearning(x = x, y = y, training_frame = train,model_id = "dl_fit2",
epochs = 20, hidden= c(10,10),seed = 1)
dl_perf2 <- h2o.performance(model = dl_fit2,newdata = test)
```


Convolutional Neural Networks

This material on CNN has been adapted from:

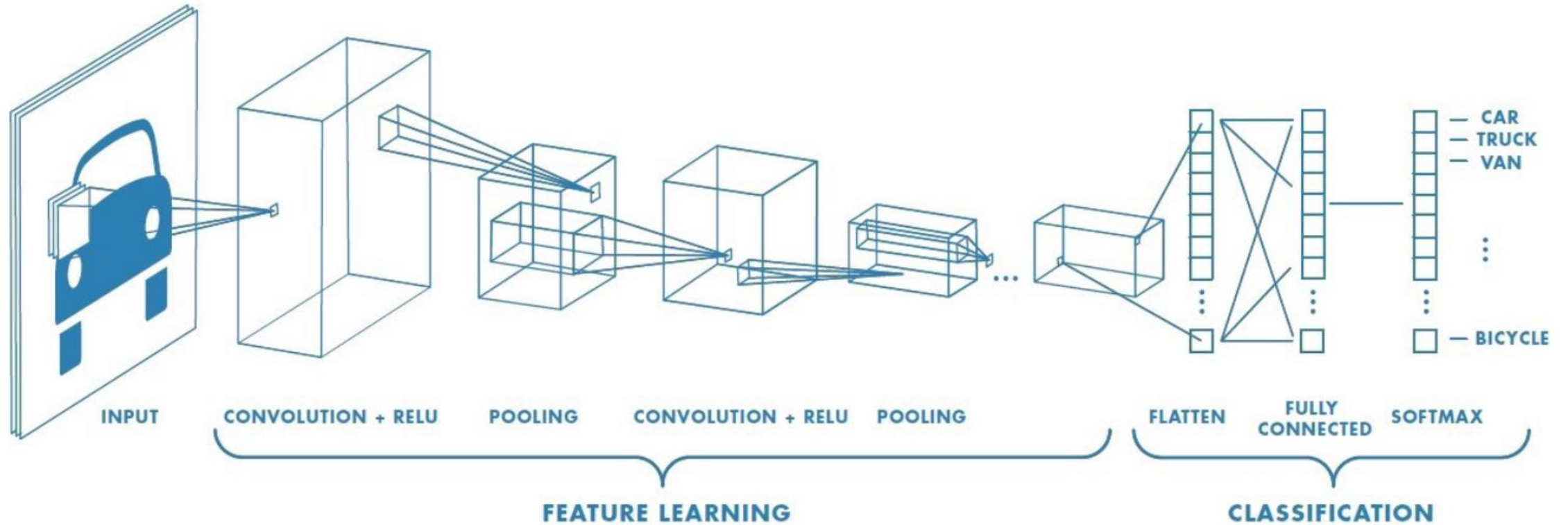
<https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-eli5-way-3bd2b1164a53>.

and slides from Prof Ming-Li website: : <http://cs.uwaterloo.ca/~mli/>

A **Convolutional Neural Network (CNN, Le Cun (1998))** is a Deep Learning algorithm which can take in an input image, assign importance (learnable weights and biases) to various aspects/objects in the image and be able to differentiate one from the other.

A CNN is able to **successfully capture the Spatial and Temporal dependencies** in an image through the application of relevant filters. The architecture performs a better fitting to the image dataset due to the reduction in the number of parameters involved and reusability of weights.

A Typical CNN



CNN (cont)

A CNN has several layers:

The convolutional layer: Its goal is to **extract the high-level features** such as eyes, nose, mouth, etc from the input image. The first CNN Layer is responsible for capturing the Low-Level features such as edges, color, gradient orientation, etc. The following layers capture the high-level features.

There are two types of results to the convolution operation — one in which the convolved feature is reduced in dimensionality as compared to the input (Valid Padding) , and the other in which the dimensionality remains the same (Same Padding). Same Padding results in padding with zeros evenly to the left/right or up/down of the input.

The Pooling layer: It is used for reducing the spatial size of the Convolved Feature. This is to **decrease the computational power required to process the data** through dimensionality reduction. Furthermore, it is useful for **extracting dominant features** which are rotational and positional invariant. There are two types of Pooling: Max Pooling and Average Pooling. **Max Pooling** returns the **maximum value** from the portion of the image covered by the Kernel. On the other hand, **Average Pooling** returns the **average of all the values** from the portion of the image covered by the Kernel.

CNN (cont)

Classification — Fully Connected Layer (FC Layer)

It is used for learning non-linear combinations of the high-level features as represented by the output of the convolutional layer. The Fully-Connected layer is learning a possibly non-linear function in that space.

Now that we have converted our input image into a suitable form for our Multi-Level Perceptron, we have to flatten the image into a column vector. The flattened output is fed to a feed-forward neural network and backpropagation is applied to every iteration of training. Over a series of epochs, the model is able to distinguish between dominating and certain low-level features in images and classify them using the **Softmax Classification** technique.

A CNN compresses a fully connected network in two ways:

- Reducing number of connections
- Shared weights on the edges
- Max pooling further reduces the complexity

Convolution

These are the network parameters to be learned.

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

⋮ ⋮

Each filter detects a small pattern (3 x 3).

Convolution

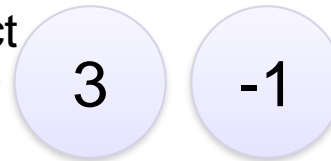
stride=1

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

Dot
product



6 x 6 image

Convolution

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

If stride=2

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

3 -3

Convolution

-1	1	-1
-1	1	-1
-1	1	-1

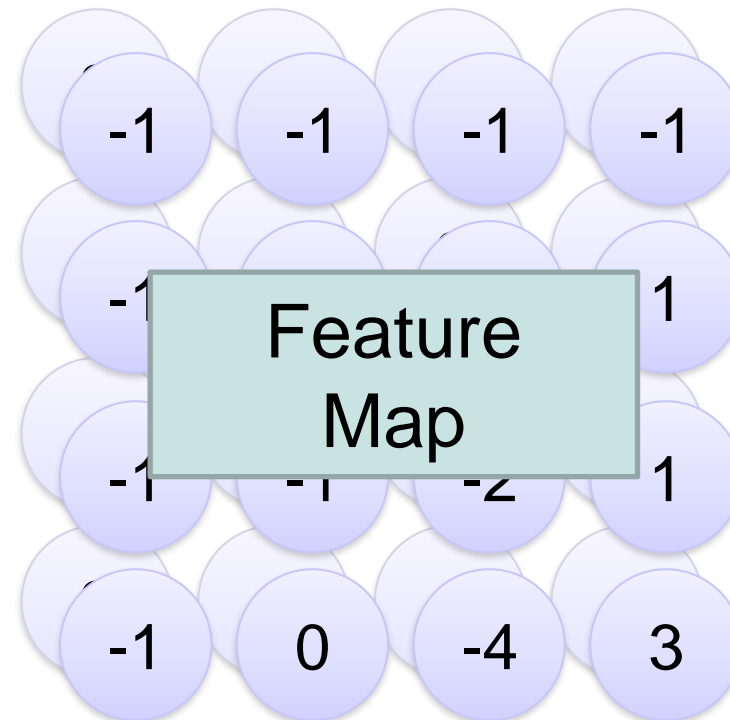
Filter 2

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

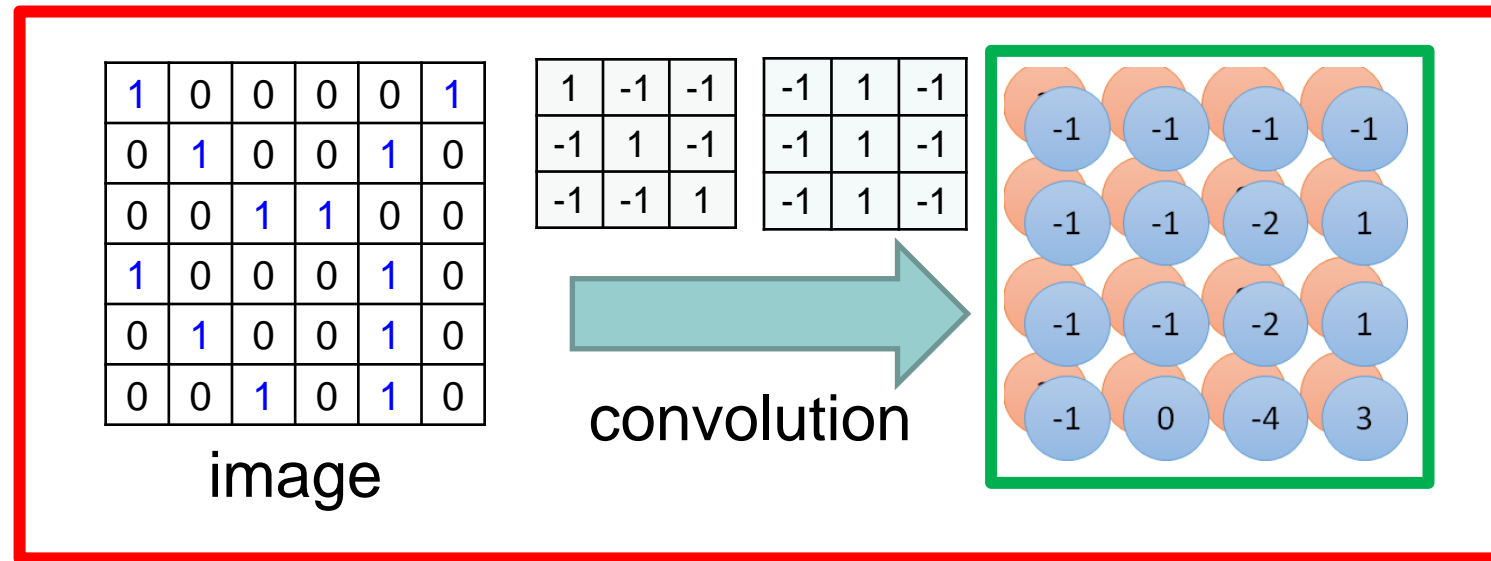
6 x 6 image

Repeat this for each filter

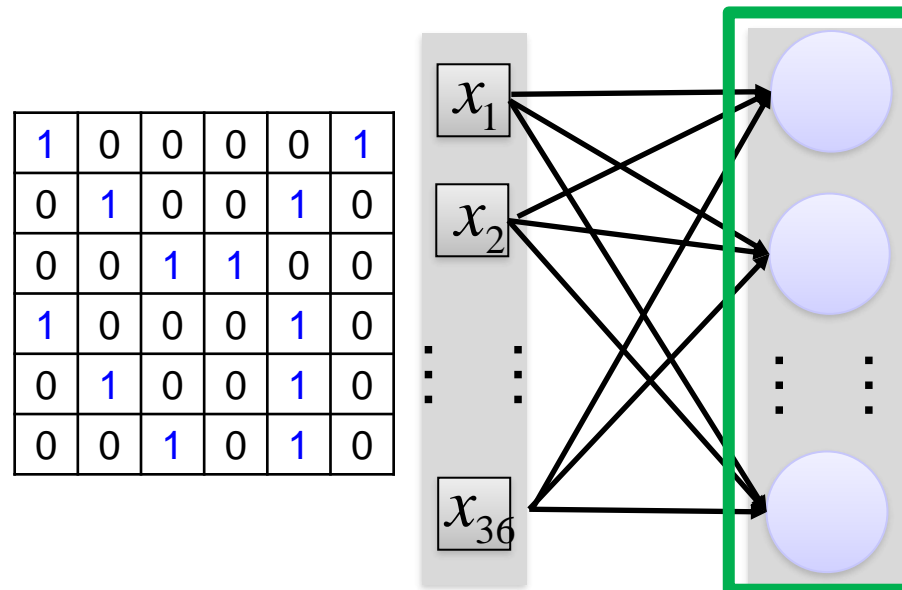


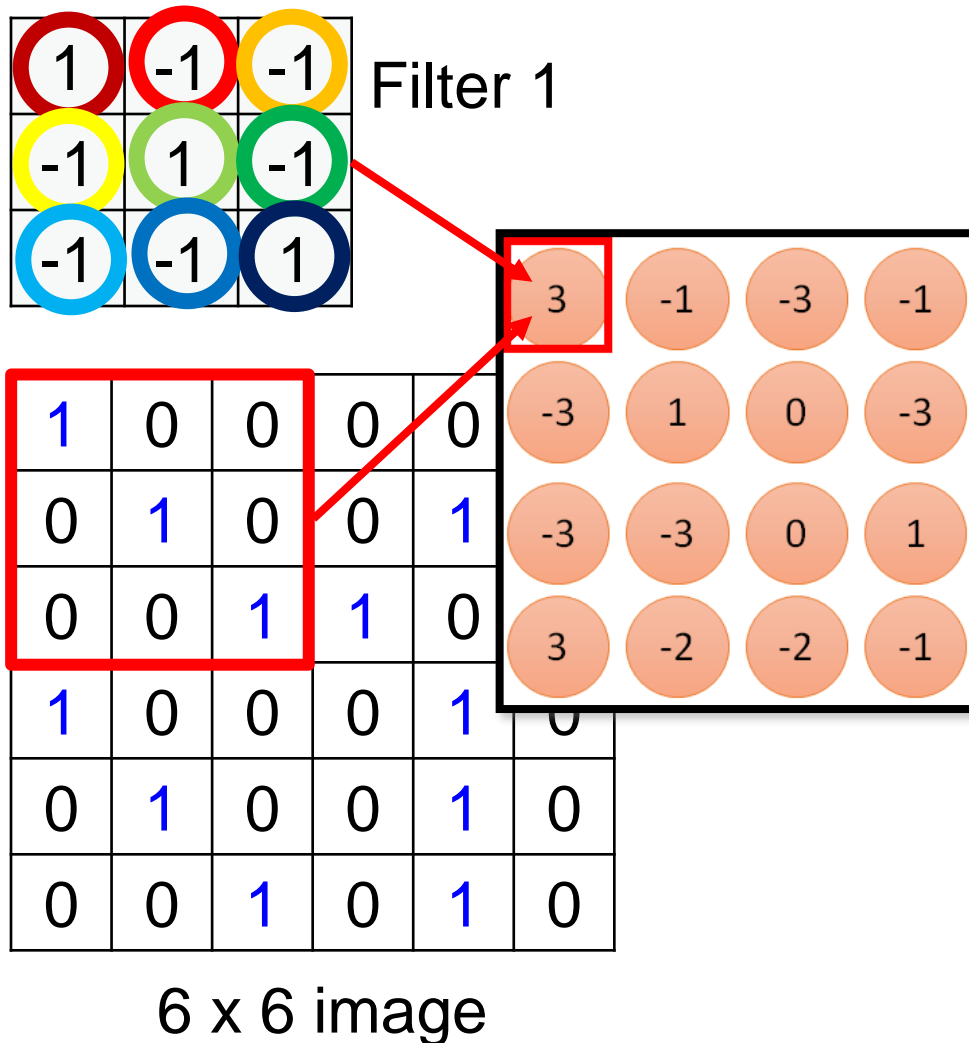
Two 4 x 4 images
Forming 2 x 4 x 4 matrix

Convolution v.s. Fully Connected

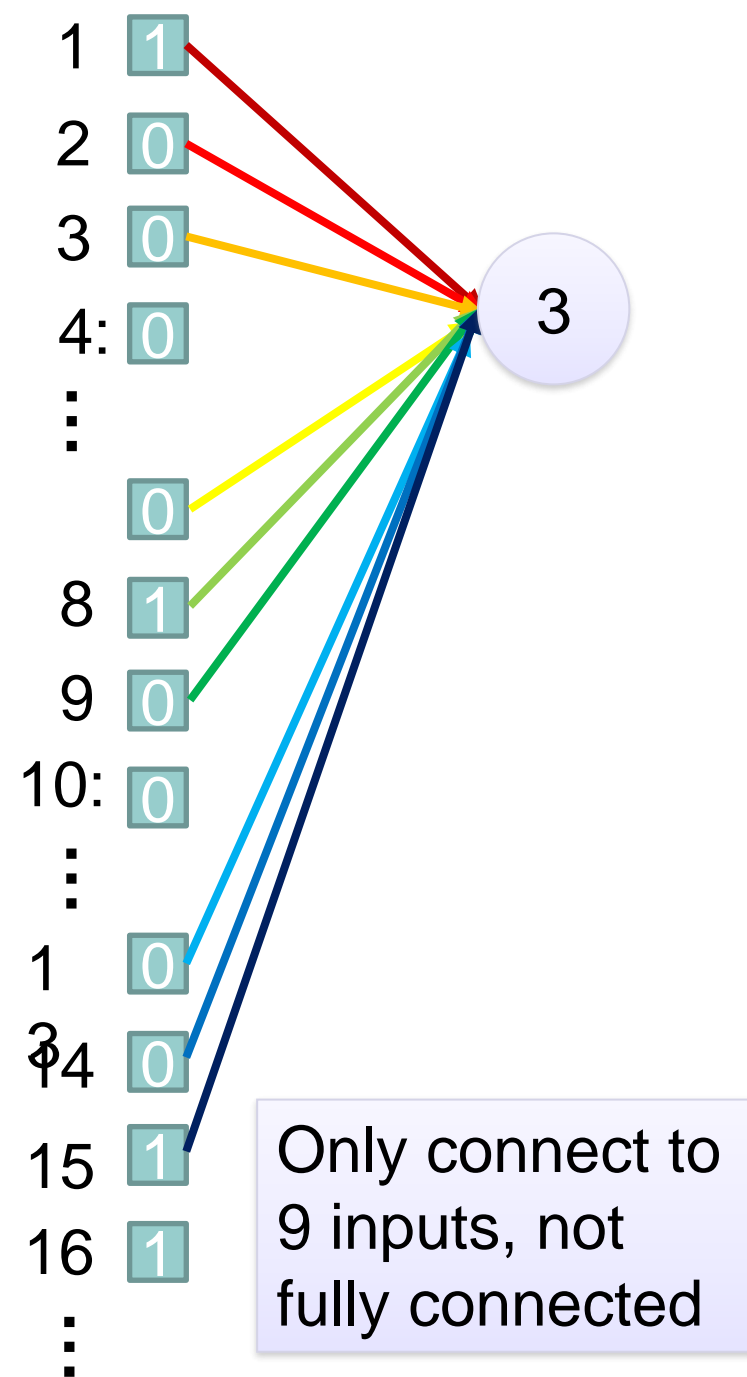


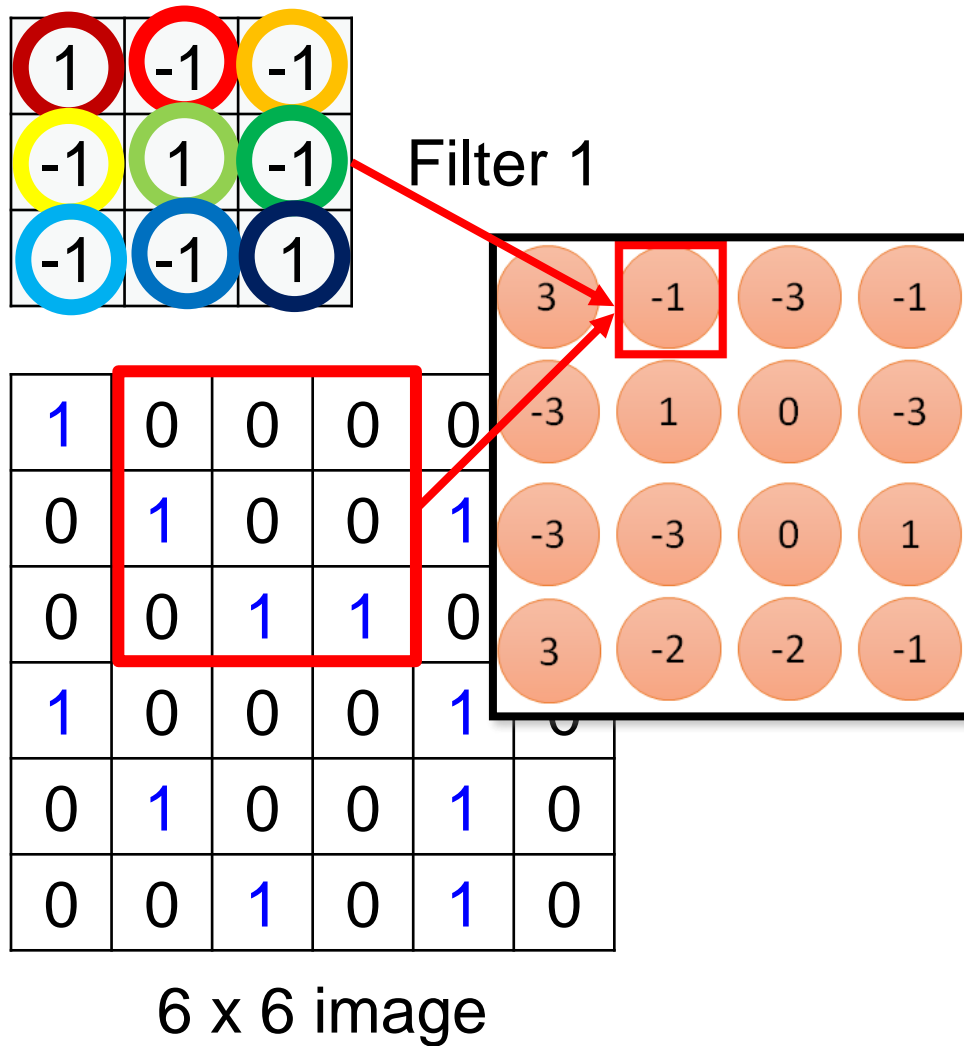
Fully-
connected





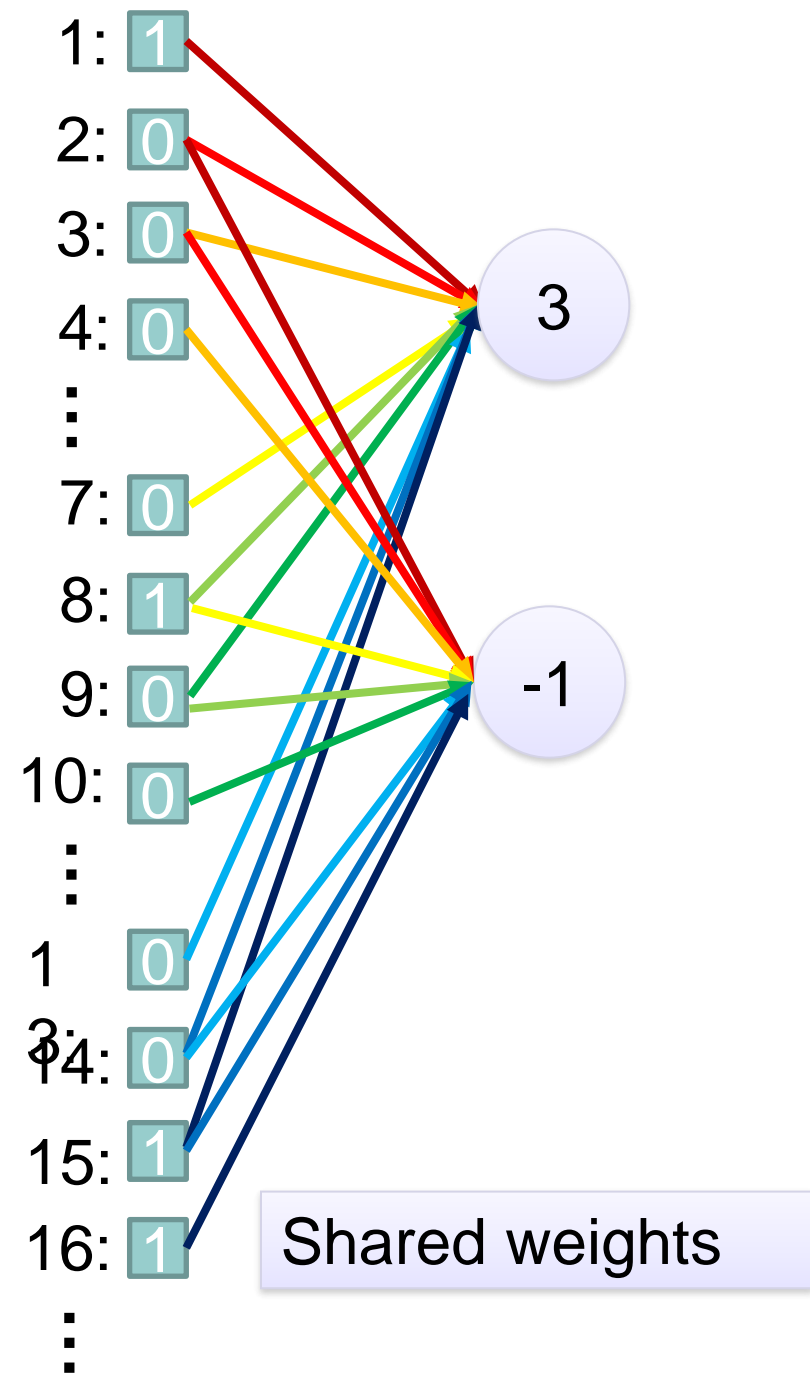
fewer parameters!





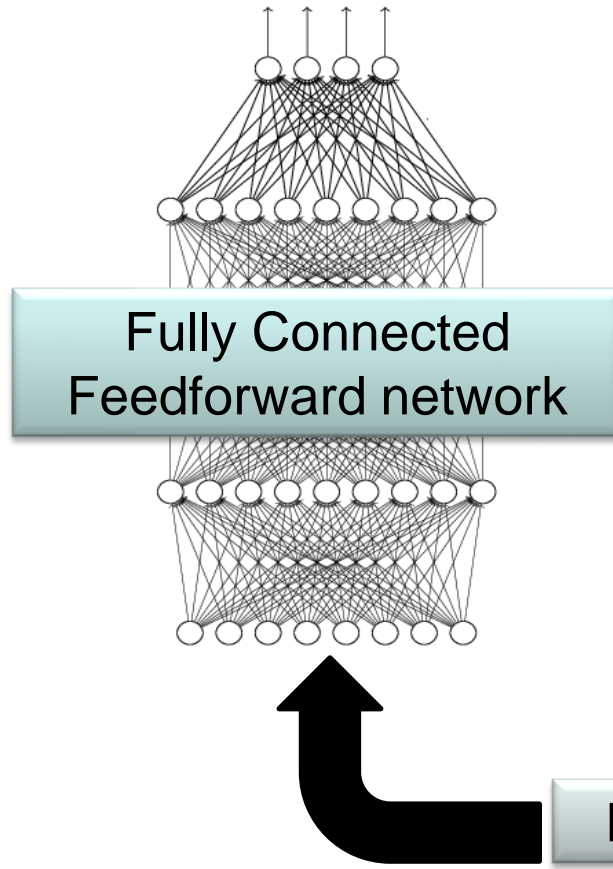
Fewer parameters

Even fewer parameters



The whole CNN

cat dog



Convolution

Max Pooling

Convolution

Max Pooling

Can
repeat
many
times

Flattened

Max Pooling

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

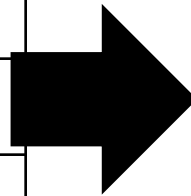
3	-1	-3	-1
-3	1	0	-3
-3	-3	0	1
3	-2	-2	-1

-1	-1	-1	-1
-1	-1	-2	1
-1	-1	-2	1
-1	0	-4	3

Max Pooling

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

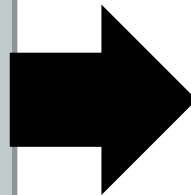
6 x 6 image



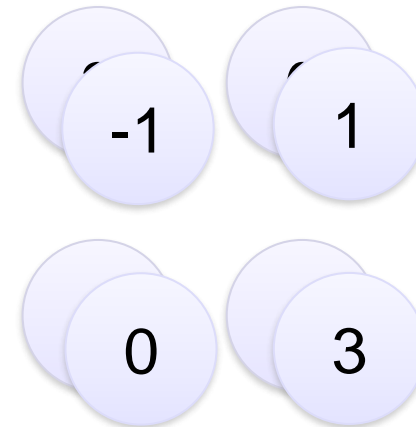
Conv



Max
Pooling



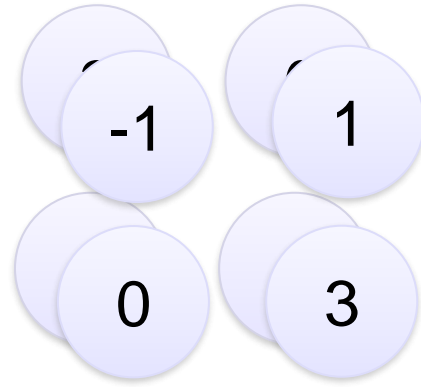
New image
but smaller



2 x 2 image

Each filter
is a channel

The whole CNN



A new image

Smaller than the original image

The number of channels is the number of filters



Convolution

Max Pooling

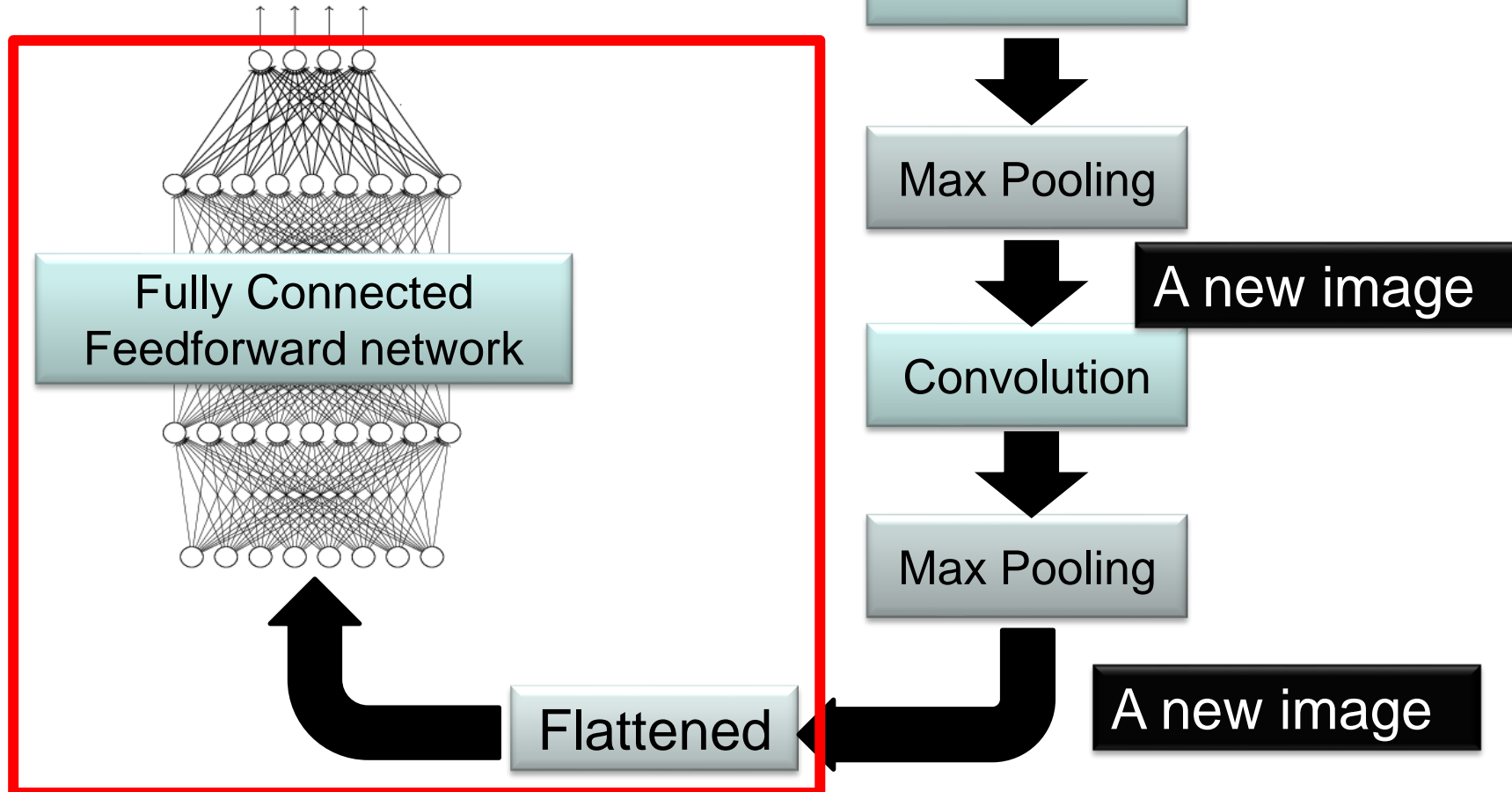
Convolution

Max Pooling

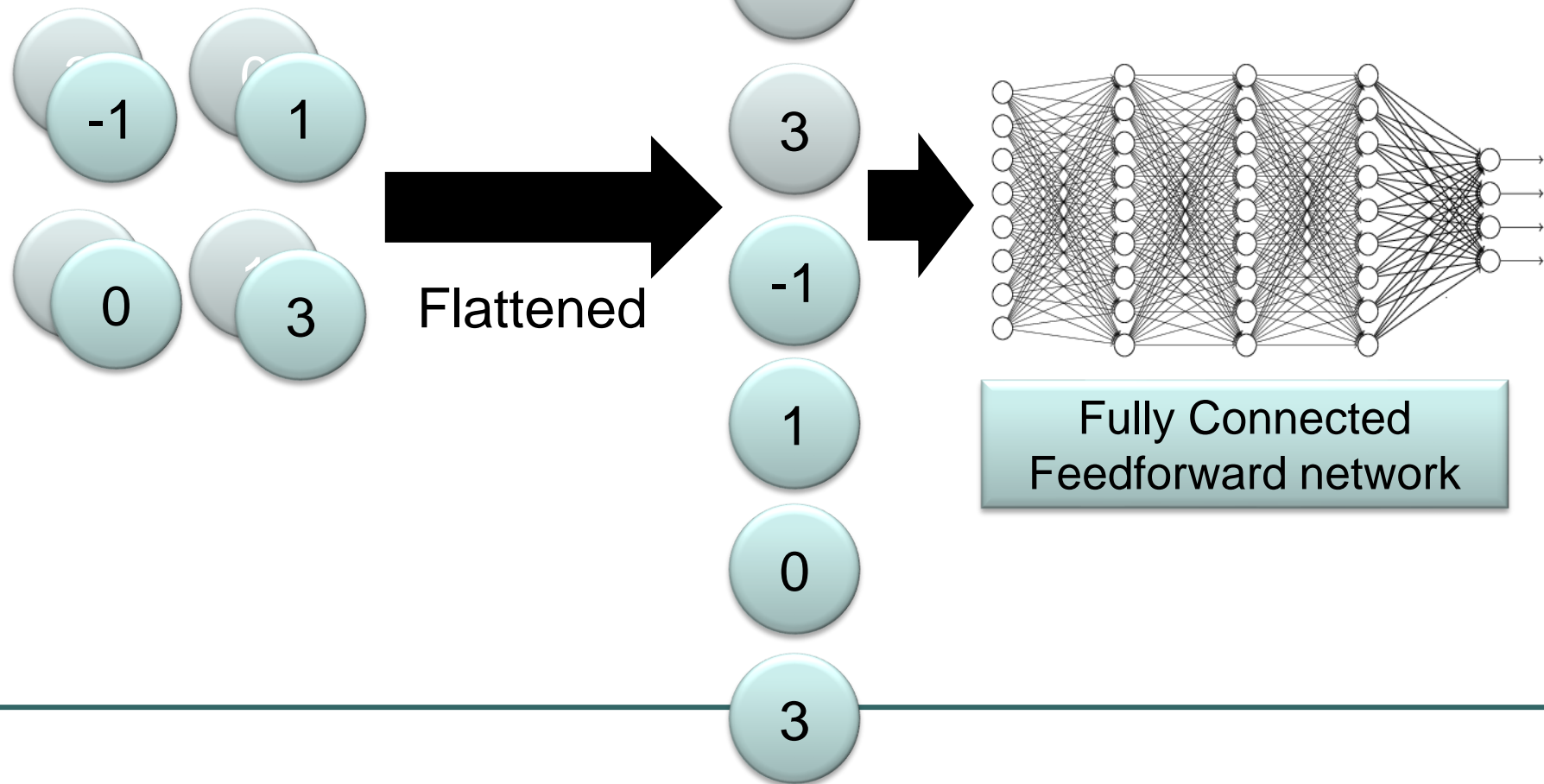
Can repeat many times

The whole CNN

cat dog



Flattening



CNN in Keras

Only modified the *network structure* and *input format* (vector -> 3-D tensor)

```
model2.add( Convolution2D( 25, 3, 3,  
                           input_shape=(28, 28, 1)) )
```

1	-1	-1	1	-1
-1	1	-1	1	-1
-1	-1	-1	1	-1
		-1	1	-1

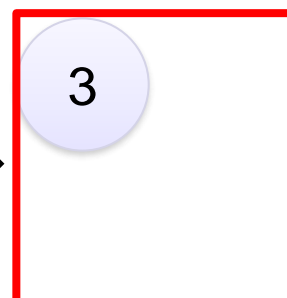
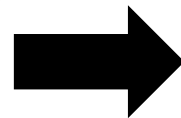
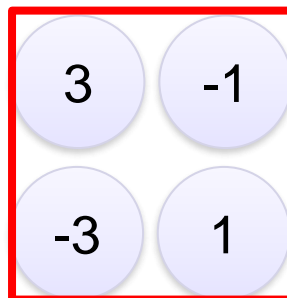
There are
25 3x3
filters.

Input_shape = (28 , 28 , 1)

28 x 28 pixels

1: black/white, 3: RGB

```
model2.add(MaxPooling2D( (2, 2) ))
```



input



Convolution



Max Pooling



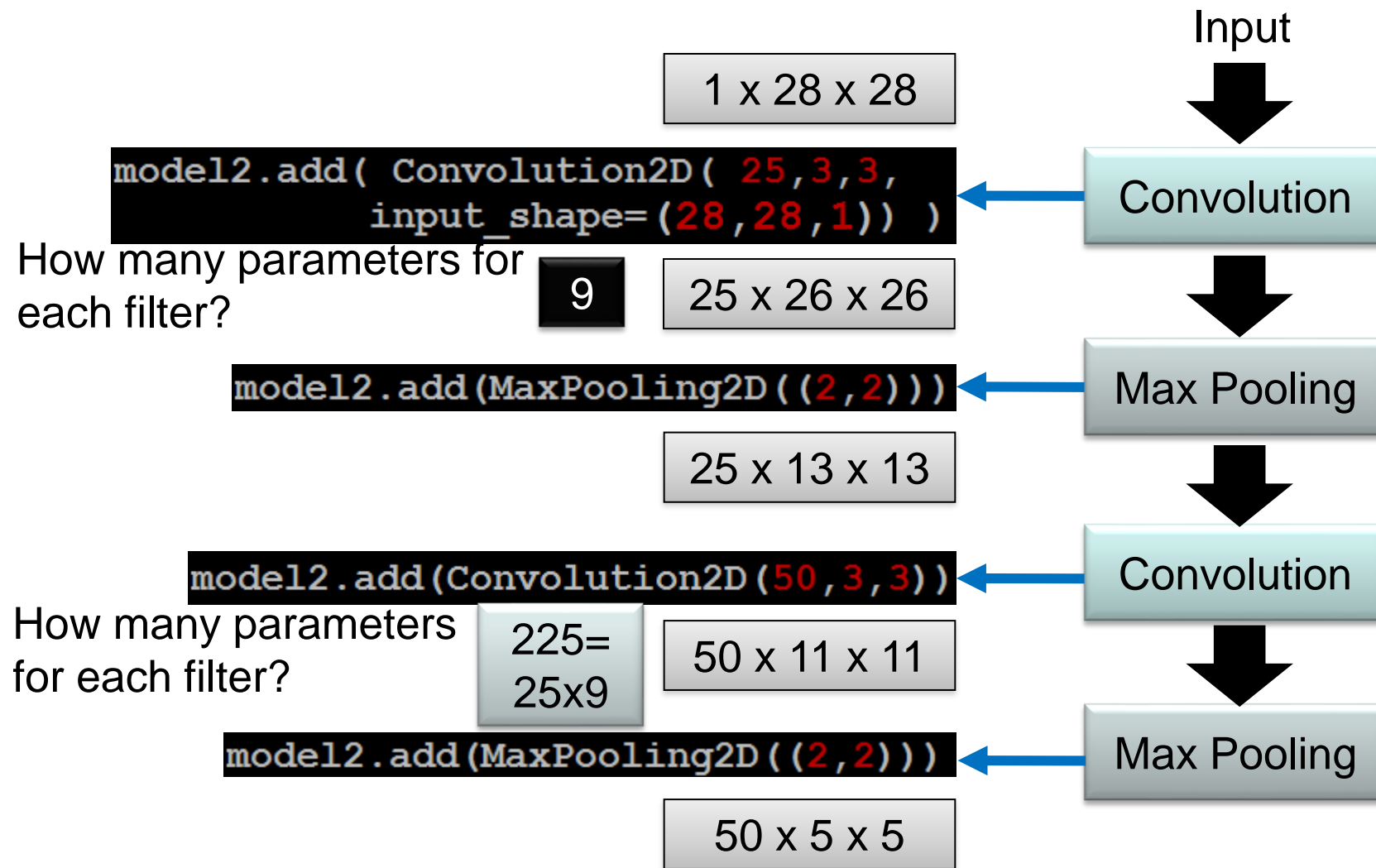
Convolution



Max Pooling

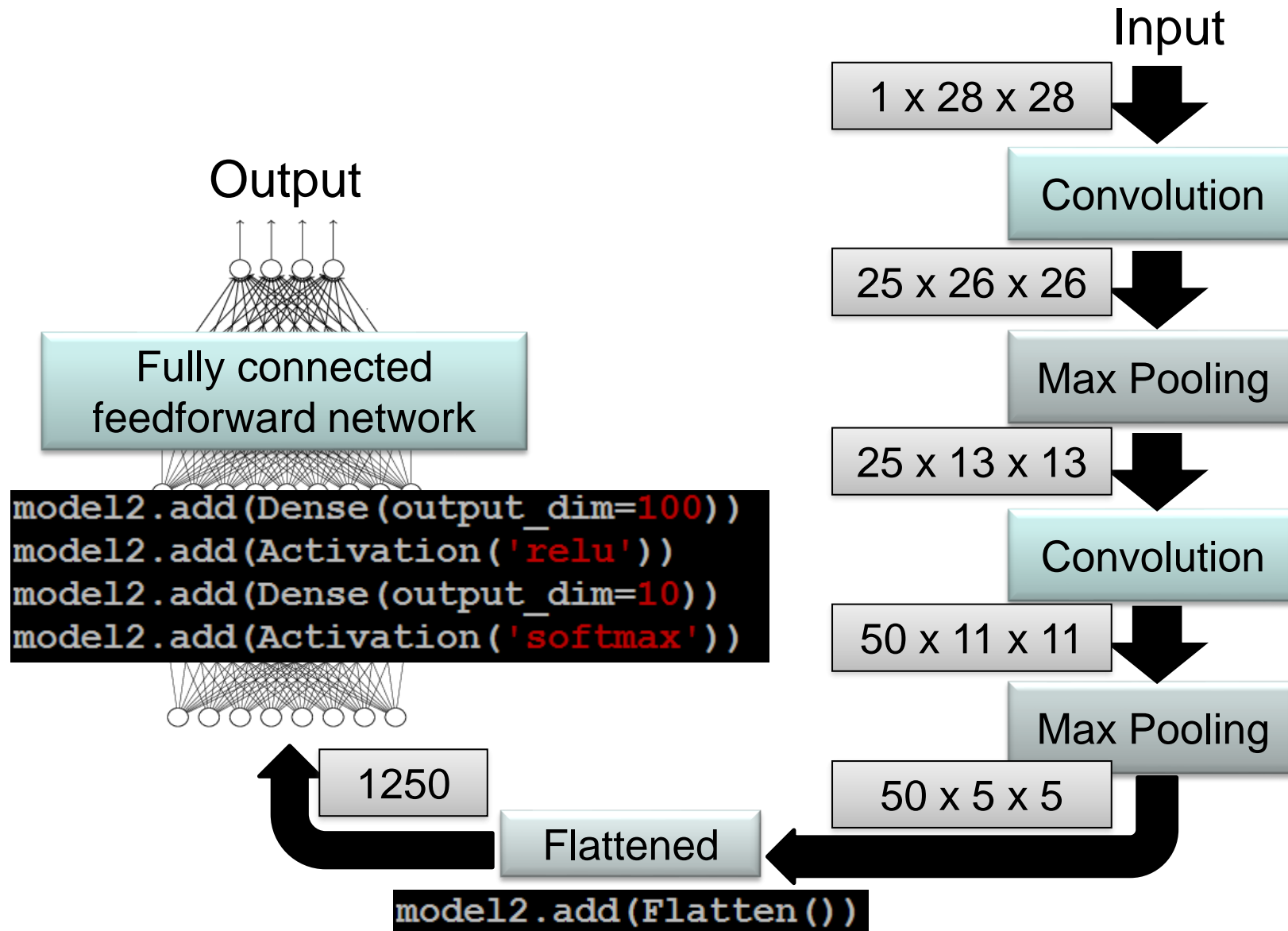
CNN in Keras

Only modified the *network structure* and *input format (vector -> 3-D array)*



CNN in Keras

Only modified the *network structure* and *input format (vector -> 3-D array)*



Data Mining and Machine Learning

Outlier Detection (Novelty Detection)
using autoencoders

Dr. Edgar Acuna
Departamento de Matematicas

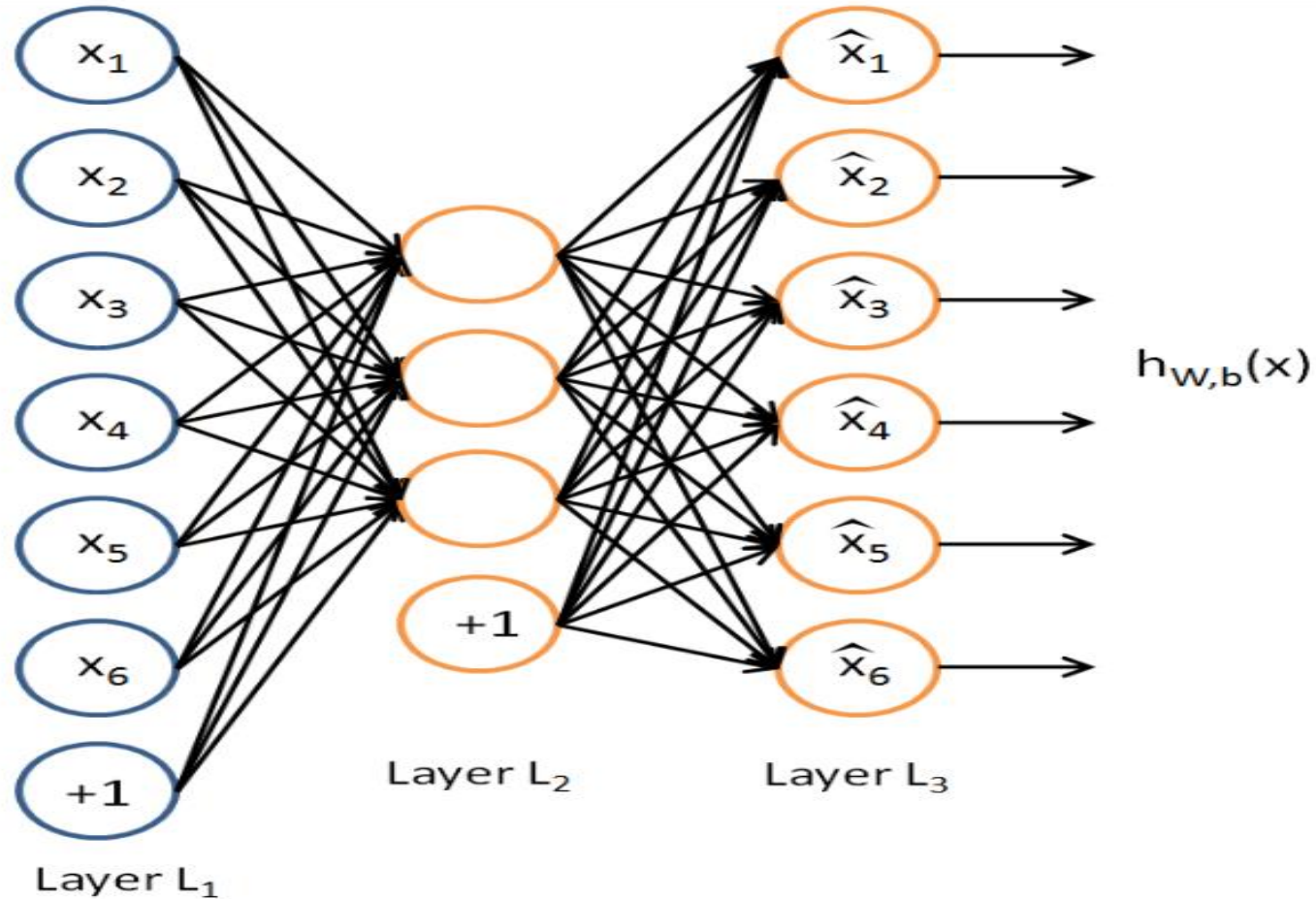
Universidad de Puerto Rico- Mayaguez

academic.uprm.edu/eacuna

Autoencoders

- Autoencoders are:
- feed Forward Neural Network models for unsupervised tasks (No Labels).
- Applies backpropagation, setting the target values to be equal to the inputs.
- simple to understand!

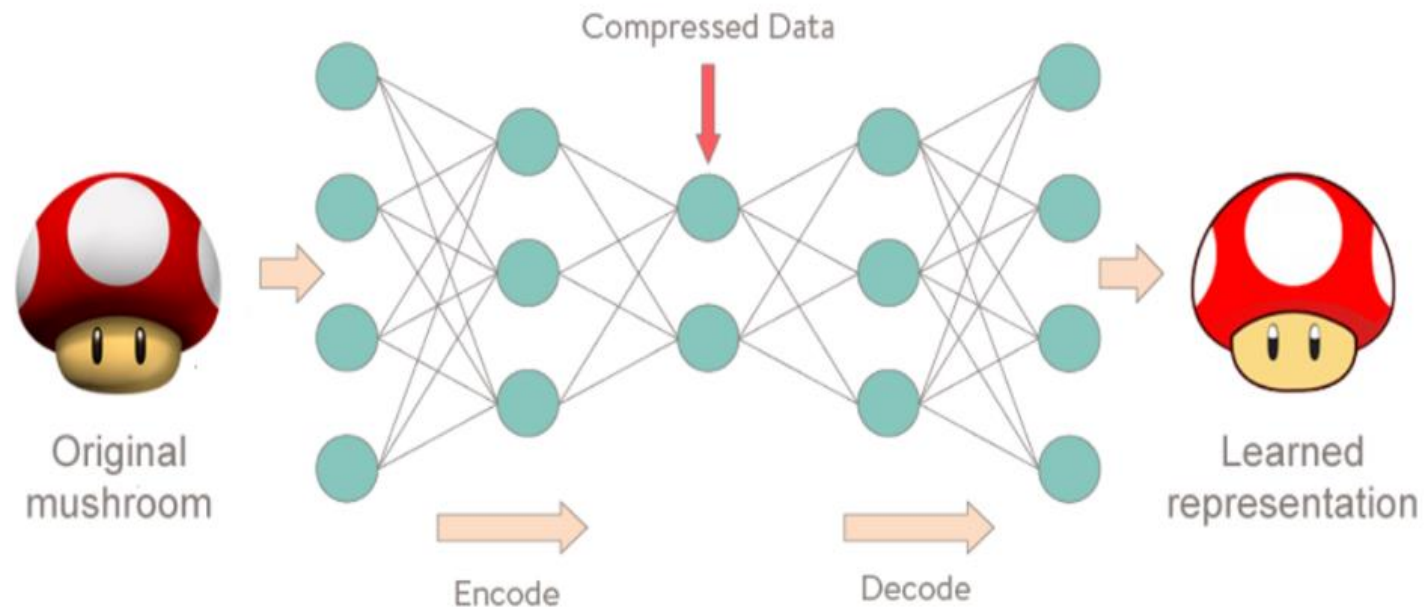
Example of autoencoder



Autoencoders are used for:

- Compressing data and learn some features, especially in non-structured data(e.g images).
- Anomaly (outlier) detections (Fraud detection is one area of application).
- Data Augmentation using the reconstructed data.

Autoencoder with 5 layers



The autoencoder model tries to minimize the reconstruction error (RE), which is the mean squared distance between input and output.

$$L(x, x') \approx ||x - x'||^2$$

Reconstruction Error

The model will train on the normal dataset by minimizing the RE. It is expected that the RE to be relatively high when it is tested on a new abnormal datapoint (outlier)

However, one has to do tuning on the RE's threshold to detect outliers

Detecting outliers for Diabetes using autoencoders

The neural network model is built by Keras. The optimized Autoencoder model has 3 fully connected hidden layers, which has 4,2,4 neurons respectively. The input and out layers have 8 neurons each for the 8 features.

The tangent hyperbolic (tanh) activation function is used in the first and third layer and the relu activation function in the second and fourth layer. The model is trained for 20 epochs and used a batch size of 50, and 10 % of the training data is used for validation during the training process.

Outliers detected on Diabetes using autoencoders

Using all the data

4, 12, 13, 39, 43, 45, 86, 88, 159, 177, 186, 193, 215,
228, 247, 259, 270, 298, 357, 362, 370, 375, 445, 453,
458, 459, 487, 542, 558, 579, 590, 674, 691, 740, 744, 763

Using train and test datasets

445, 428, 123, 228, 274, 660, 517, 323, 339, 100, 588, 708,
691, 582, 614, 84, 24, 459, 509, 319, 12, 612, 672, 298

Deep Learning for Sequential Data

Examples of sequential Data are time series and spectral data

We can use:

One dimensional Convolution: Conv1D followed by one dimensional MaxPooling

Recurrent Neural Networks (RNN),

Long- Short Term Memory (LSTM) models, and

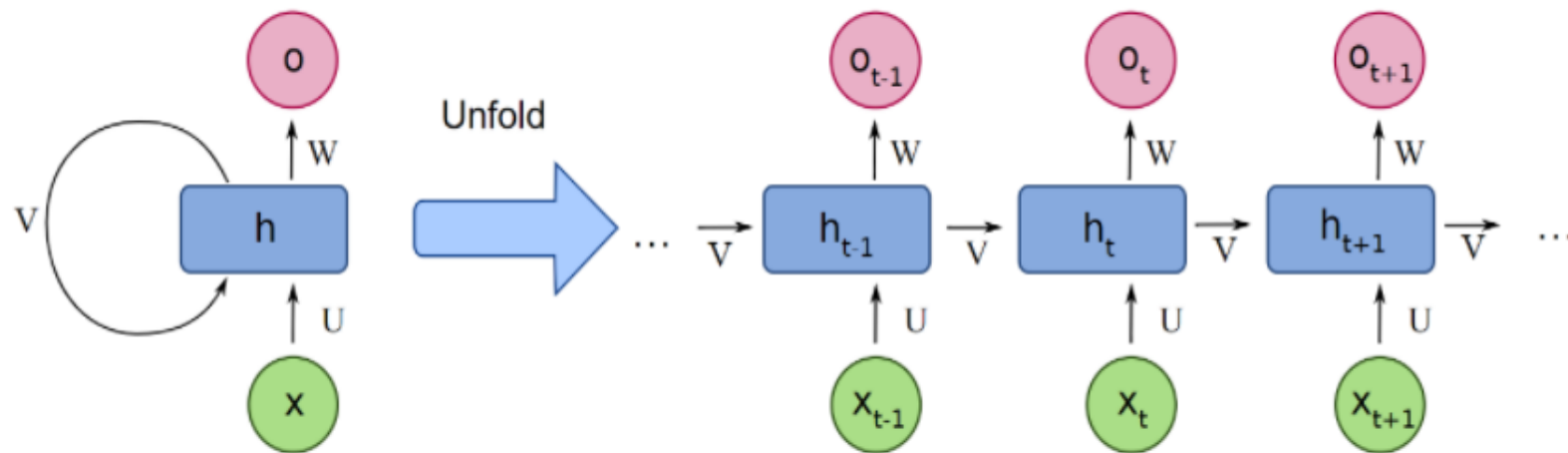
Hybrids models such as CNN-LSTM

RNN and LSTM Networks

In the MLP model, the pre-layer and post-layer are fully connected, and there is no connections between the neurons in the same layer.

The recurrent Neural Network (RNN, Rumelhart 1986) adds a weighted sum of the hidden layer with the previous input when calculating the output of the hidden layer. Therefore, the input of the hidden layer includes not only the output of the pre-layer, but also all the previous outputs. This introduces a feedback mechanism in the hidden layer to learn the context-related information, which can effectively process the sequence data (e.g. time series, spectral data speech recognition data). There are several type of RNN according to the number o inputs and outputs.

RNNs use a slightly different approach to traditional neural networks in the form of backpropagation through time also known as BPTT. BPTT is more geared towards sequential data and the main difference is that it also sums the errors at each time step. At every timestep the RNNs are fed one input and it predicts an output using a copy of the network and one output. BPTT calculates the error accumulated in each individual timestep and when finished it rolls back to update the weights as the normal backpropagation algorithm does.



Representation of RNN both in folded and unfolded forms

$$h_t = f(Vh_{t-1} + Ux_t)$$

$$o_t = Wh_t$$

U, V and W are parameters to be estimated

h is called the RNN cell or hidden state

h_{t-1}, h_t, h_{t+1} are called hidden units of the hidden state

RNN and LSTM (continued)

Recurrent Neural Networks suffer from short-term memory. If a sequence is long enough, they'll have a hard time carrying information from earlier time steps to later ones. RNN's may leave out important information from the beginning.

As iterations of the RNN occur, the gradient continues to shrink, and the weights of the model tend to 0 making the algorithm unable to learn anything else. This is called the vanishing gradient problem resulting that the continued training does not produce any further improvement.

RNN cannot process very long sequences if using *tanh* or *relu* as an activation function. Its computation is very slow. Also, RNN cannot consider any future input for the current state.

RNN and LSTM (continued)

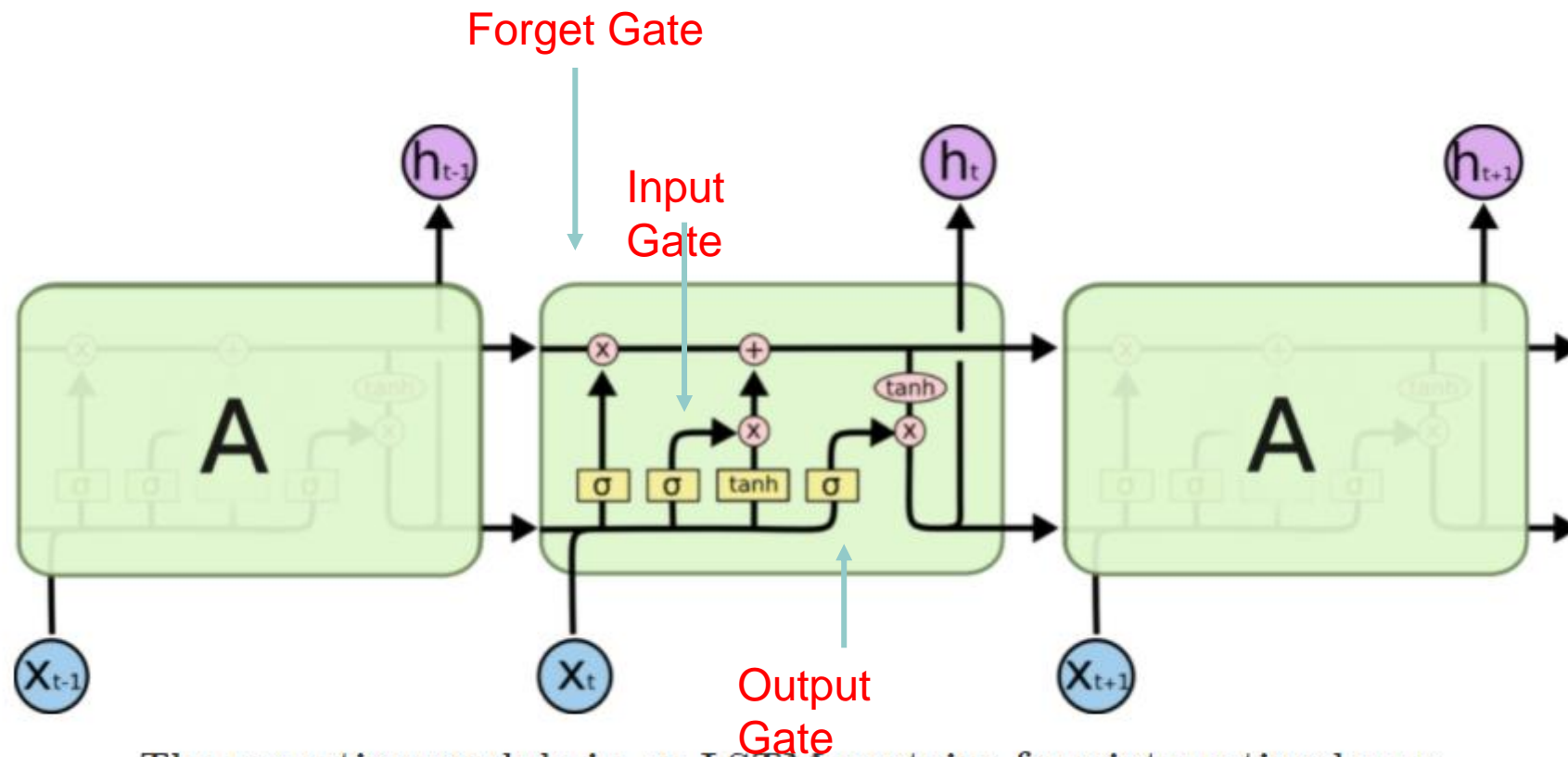
There are two improved version of the RNN model. Gated Recurrent Unit (GRU) and Long Short term Memory (LSTM). Both of them have internal mechanisms called gates that can regulate the flow of information.

GRU has two gates: reset and update.

Long short-term memory or LSTM (Hochreiter and Schmidhuber, 1997) is a recurrent neural network that is able to deal with long-term dependencies between inputs and outputs.

LSTM has three gates: Forget gate, Input Gate and Output Gate. LSTM can learn which data in a sequence is important to keep or throw away. By doing that, it can pass relevant information down the long chain of sequences to make predictions.

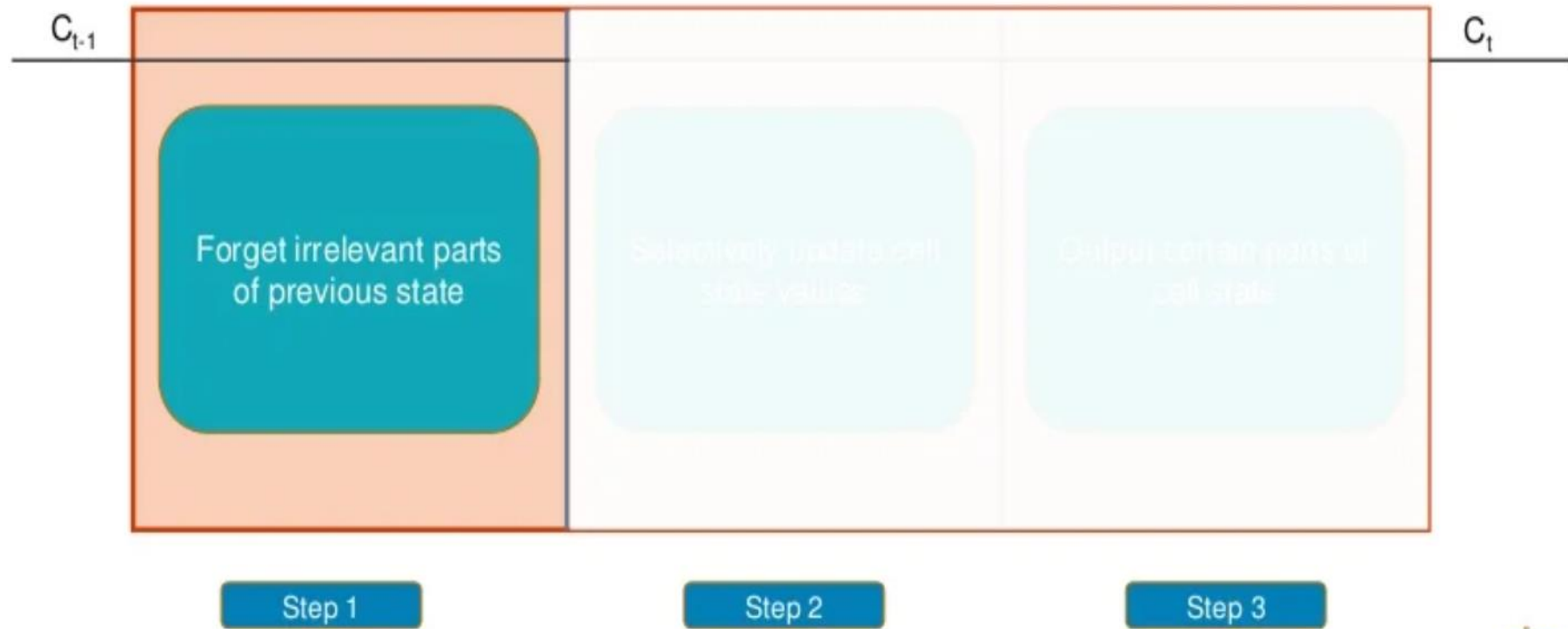
But still sequences can not have more than 400 timesteps.



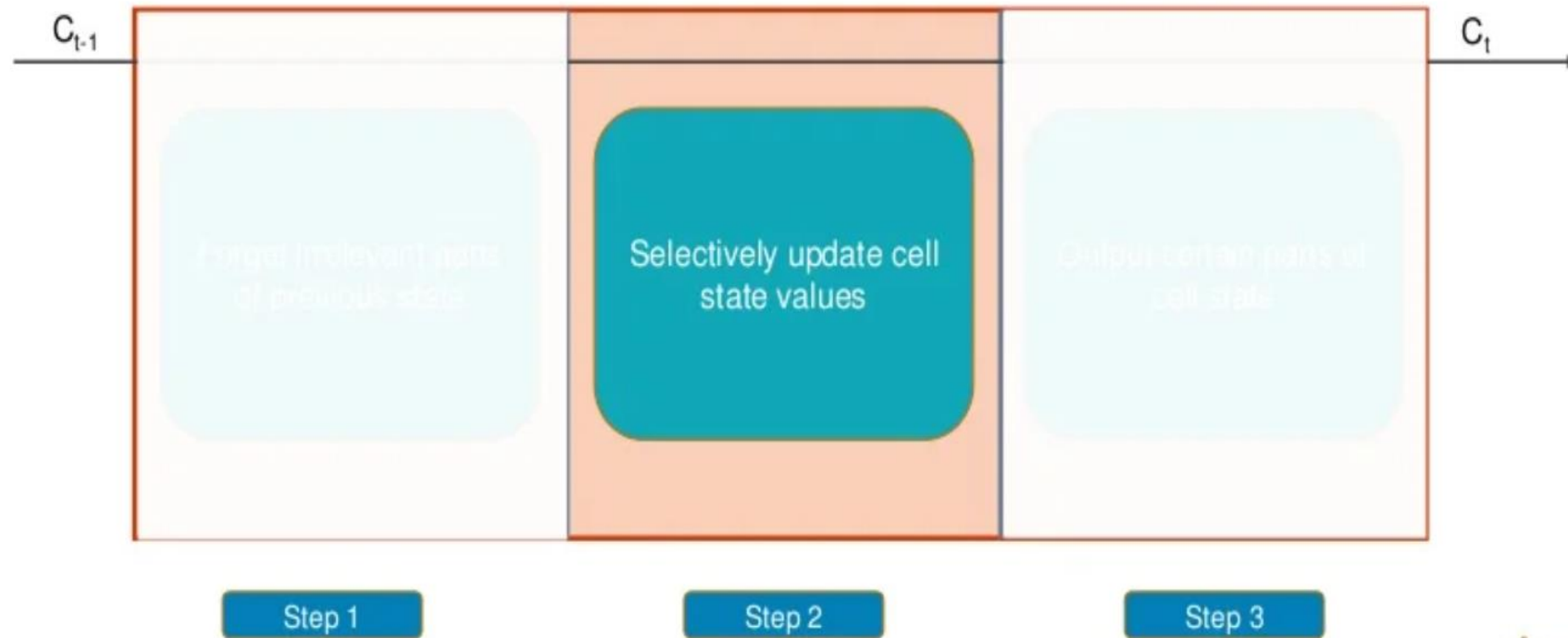
The repeating module in an LSTM contains four interacting layers.

Source: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

3 step process of LSTMs



3 step process of LSTMs



3 step process of LSTMs



CNN-LSTM

CNN performs feature extraction in subsequences of the main sequence. The extracted features are then flattened and provided to the LSTM model to read, extracting its own features before a final mapping to an analyte is made (Sainath et.al, 2015).

We can define a CNN-LSTM model by first defining the CNN layers, wrapping them in a Time Distributed layer and then defining the LSTM and output layers.