

Machine Learning in HEP

A. Sznajder

UERJ
Instituto de Fisica

September - 2019

Outline

- 1 What is Machine Learning ?
- 2 Neural Networks
- 3 Learning as a Minimization Problem (Gradient Descent and Backpropagation)
- 4 Deep Learning Revolution
- 5 Deep Architectures and Applications

What is Machine Learning ?

Machine Learning (ML)

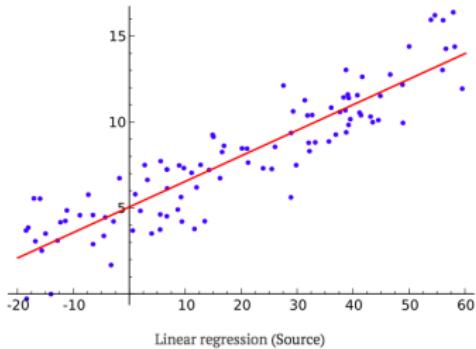
Machine learning (ML) is the study of computer algorithms capable of building a mathematical model out of a data sample, by learning from examples. The algorithm builds a predictive model without being explicitly programmed (sequence of instructions) to do so



Centuries Old Machine Learning¹

Take some points on a 2D graph, and fit a function to them. What you have just done is generalized from a few (x, y) pairs (examples) , to a general function that can map any input x to an output y

The Centuries Old Machine Learning Algorithm

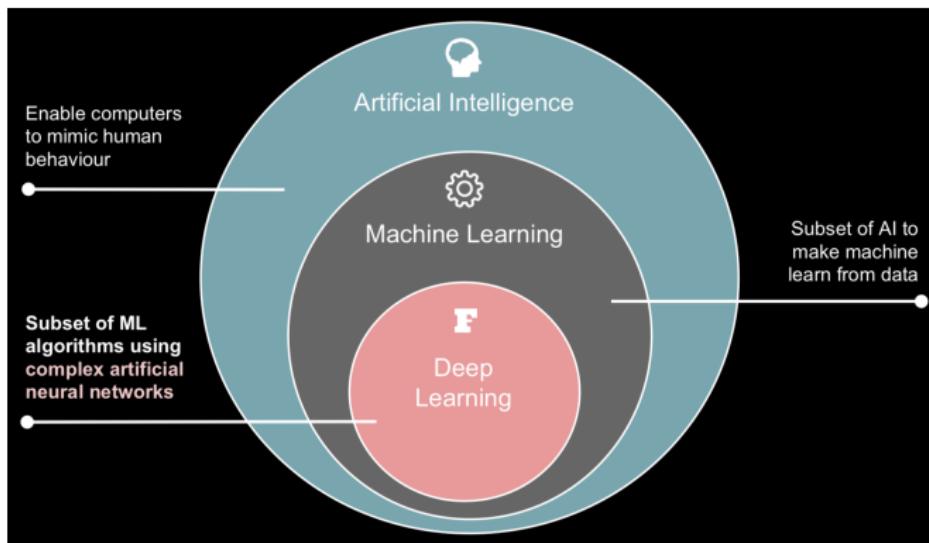


Linear regression is a bit too wimpy a technique to solve the problems of image , speech or text recognition, but what it does is essentially what supervised ML is all about: learning a function from a data sample

¹<http://www.andreykurenkov.com/writing/ai/a-brief-history-of-neural-nets-and-deep-learning/>

Artificial Intelligence

Intelligence is the ability to process current information to inform future decisions



Artificial Intelligence

None of the systems we have nowadays are real AI. The brain learns so efficiently that no ML method can match it

AI versus Brain

- Brain has 10^{14} parameters and we live only 10^9 s (lot more parameters than data)
- So, it must do lots of unsupervised learning and must predict what we observe !
- Supervised (reinforcement) learning requires millions of examples(trials)
- Still missing a learning paradigm that builds predictive models through observation and action

Yann LeCun on the Epistemology of Deep Learning:

<https://www.youtube.com/watch?v=gG5NCkMerHU&t=944s>

<https://www.youtube.com/watch?v=cWzi38-vDbE&t=768s>

Introduction to Machine Learning

Machine learning can be implemented by many different algorithms (ex: SVM, BDT, Bayesian Net , Genetic Algo ...) , but we will discuss only Neural Networks(NN)

1) Neural Network Topologies

- Feed Forward NN
- Recurrent NN

2) Learning Paradigms

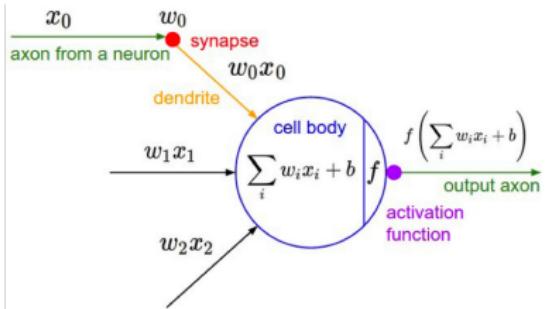
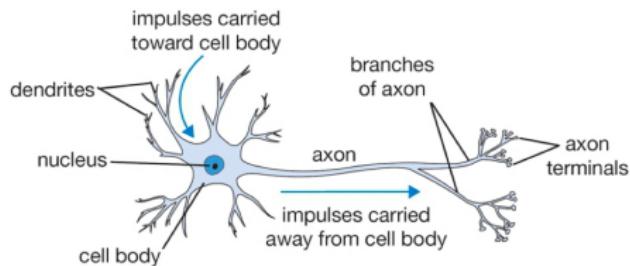
- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning

3) Neural Networks Architectures

- Multilayer Perceptron(MLP)
- Convolutional Neural Network (CNN)
- Recurrent Neural Network (RNN)
- Long Short Time Memory(LSTM)
- Autoencoder(AE)
- Variational Autoencoder(VAE)
- Generative Adversarial Network(GAN)

Neural Networks

Artificial Neural Networks (NN) are computational models vaguely inspired² by biological neural networks. A Neural Network (NN) is formed by a network of basic elements called neurons, which receive an input, change their state according to the input and produce an output.



Original goal of NN approach was to solve problems like a human brain. However, focus moved to performing specific tasks, deviating from biology. Nowadays NN are used on a variety of tasks: image and speech recognition, translation, filtering, playing games, medical diagnosis, autonomous vehicles, ...

²Design of airplanes was inspired by birds, but they don't flap wings to fly !

Artificial Neuron

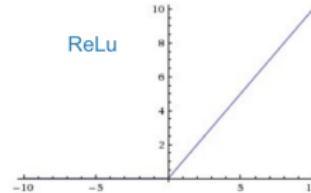
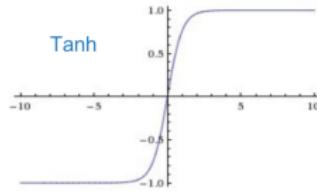
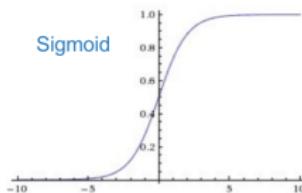
Artificial Neuron (Node)

Each node of a NN receives inputs $\vec{x} = \{x_1, \dots, x_n\}$ from other nodes or an external source and computes an output y according to the expression

$$y = F \left(\sum_{i=1}^n W_i x_i + B \right) = F(\vec{W} \cdot \vec{x} + B) \quad (1)$$

, where W_i are connection weights, B is the threshold and F the activation function ³

There are a variety of possible activation function and the most common ones are



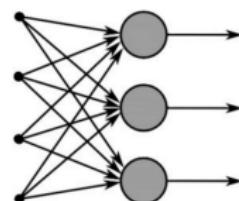
³ Nonlinear activation is fundamental for nonlinear decision boundaries

Neural Network Topologies

Neural Networks can be classified according to the type of neuron interconnections and the flow of information

Feed Forward Networks

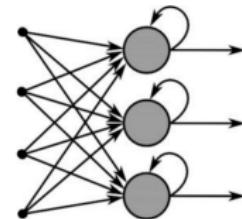
A feedforward NN is a neural network wherein connections between the nodes do not form a cycle. The information always moves one direction, from input to output, and it never goes backwards.



Feed-Forward Neural Network

Recurrent Neural Network

A Recurrent Neural Network (RNN) is a neural network that allows connections between nodes in the same layer, with themselves or with previous layers. RNNs can use their internal state (memory) to process sequential data.



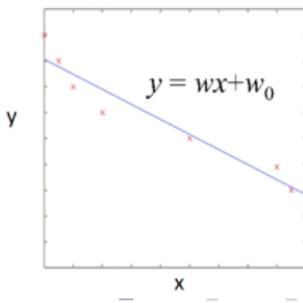
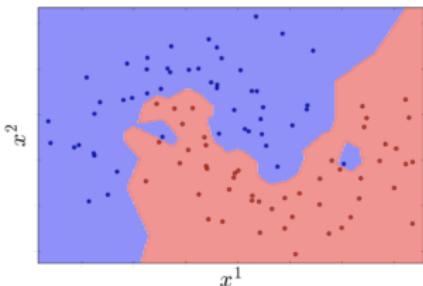
Recurrent Neural Network

A NN layer is called a dense layer to indicate that it's fully connected.

Supervised Learning

Supervised Learning

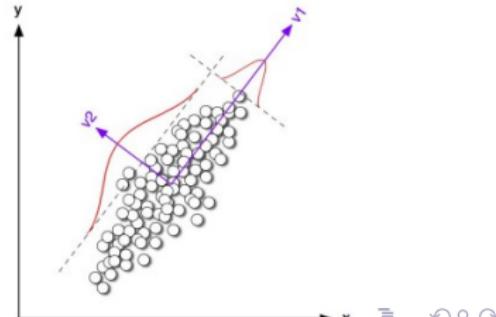
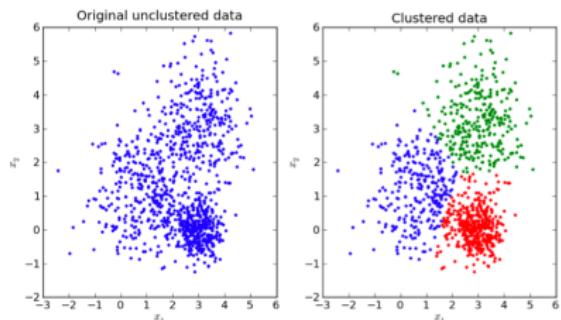
- During training a learning algorithm adjust the network's weights to values that allows the NN to map the input to the correct output.
- It calculates the error between the target output and a given NN output and use error to correct the weights.
- Given some labeled data $D = \{(\vec{x}_1, \vec{t}_1), \dots, (\vec{x}_n, \vec{t}_n)\}$ with features $\{\vec{x}_i\}$ and targets $\{\vec{t}_i\}$, the algorithm finds a mapping $\vec{t}_i = F(\vec{x}_i)$
- Classification:** $\{\vec{t}_1, \dots, \vec{t}_n\}$ (finite set of labels)
- Regression:** $\vec{t}_i \in \mathbb{R}^n$



Unsupervised Learning

Unsupervised Learning

- No labeled data is used at training and the NN finds patterns within input data
- Given some data $D = \{\vec{x}_1, \dots, \vec{x}_n\}$, but no labels, find structures in the data
 - Clustering:** partition the data into sub-groups $D = \{D_1 \cup D_2 \cup \dots \cup D_k\}$
 - Dimensional Reduction:** find a low dimensional representation of the data with a mapping $\vec{y} = F(\vec{x})$, where $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with $n \gg m$

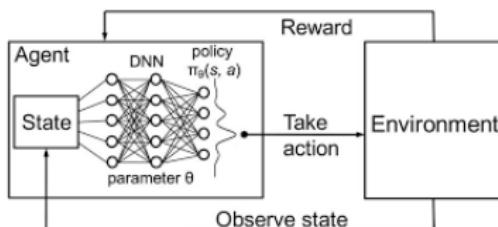


Reinforcement Learning

Reinforcement Learning

Similar to supervised learning but instead of a target output, a reward is given based on how well the system performed. The algorithm takes actions to maximize some notion of cumulative reward.

- Inspired by behaviorist psychology and strongly related with how learning works in nature
- Maximize the reward the system receives through trial-and-error
- Algorithm learns to make the best sequence of decisions to achieve goal.



Requires a lot of data, so applicable in domains where simulated data is readily available: robotics, self-driving vehicles, gameplay ...

Reinforcement Learning

Google Deepmind neural network playing Atari game ⁴

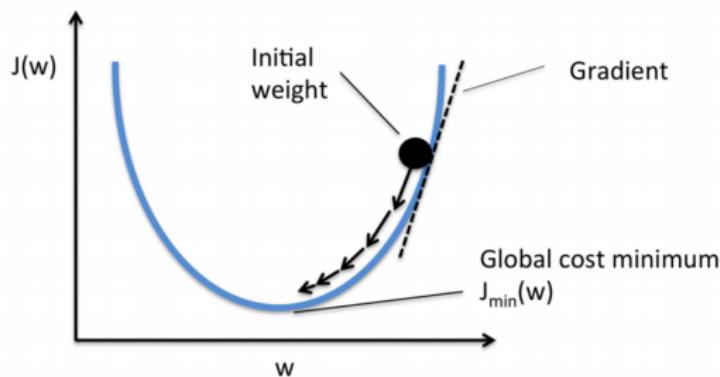


⁴<https://www.youtube.com/watch?v=TmPfTpjtdgg>

Supervised Learning - Training Process

Learning as an Error Minimization Problem

- ① Random NN parameters (weights and bias) initialisation
- ② Choose a Loss Function , differentiable with respect to model parameters
- ③ Use training data to adjust parameters (gradient descent + back propagation) that minimize the loss
- ④ Repeat until parameters values stabilize or loss is below a chosen threshold



Supervised Learning - Loss Function

The Loss function quantifies the error between the NN output $\vec{y} = F(\vec{x})$ and the desired target output \vec{t} .

Squared Error Loss (Regression)

If we have a target $t \in \mathcal{R}$ and a real NN output y

$$L = (y - \vec{t})^2$$

Cross Entropy Loss⁵ (Classification)

If we have m classes with binary targets $t \in \{0, 1\}$ and output probabilities y :

$$L = - \sum_{i=1}^m t_i \log(y_i)$$

Cost (Objective) Function

The Cost function is the mean Loss over a data sample $\{\vec{x}_i\}$

$$\bar{L} = \frac{1}{n} \sum_{i=1}^n L(\vec{x}_i)$$

The activation function type used in the output layer is directly related to loss used for the problem !

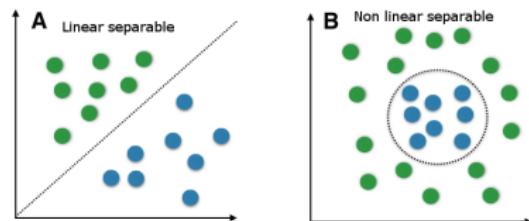
⁵For $m = 2$ we have the Binary Cross Entropy $L = -t \log y - (1 - t) \log(1 - y)$

The Perceptron

The perceptron algorithm is a binary linear classifier invented in 1957 by F.Rosenblatt. It's formed by a single neuron that takes input $\vec{x} = (x_1, \dots, x_n)$ and outputs $y = 0, 1$ according to

Perceptron Model⁶

$$y = \begin{cases} 1, & \text{if } (\vec{W} \cdot \vec{x} + B) > 0 \\ 0, & \text{otherwise} \end{cases}$$



To simplify notation define $W_0 = B$, $\vec{x} = (1, x_1, \dots, x_n)$ and call θ the Heaviside step function

Perceptron Learning Algorithm

Initialize the weights and for each example (\vec{x}_j, t_j) in training dataset D , perform the following steps:

- ① Calculate the output error: $y_j = \theta(\vec{W} \cdot \vec{x}_j)$ and $Error = 1/m \sum_{j=1}^m |y_j - t_j|$
- ② Modify(update) the weights to minimize the error: $\delta W_i = r \cdot (y_j - t_j) \cdot X_i$, where r is the learning rate
- ③ Return to step 1 until output error is acceptable

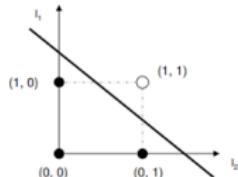
⁶Equation of a plane in \mathbb{R}^n is $\vec{W} \cdot \vec{x} + B = 0$

Perceptron and XOR Problem

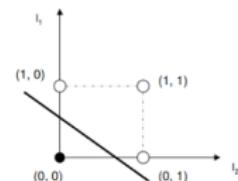
Problem:

Perceptrons are limited to linearly separable problems => unable to learn the *XOR* boolean function

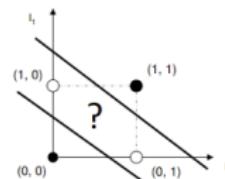
AND		
I_1	I_2	out
0	0	0
0	1	0
1	0	0
1	1	1



OR		
I_1	I_2	out
0	0	0
0	1	1
1	0	1
1	1	1



XOR		
I_1	I_2	out
0	0	0
0	1	1
1	0	1
1	1	0



Solution:

Need two neurons (layer) to solve the *XOR* problem in two-stages

Multilayer Perceptron (MLP)

The Multilayer Perceptron(MLP) is a fully connected NN with at least 1 hidden layer and nonlinear activation function F ⁷. It is the simplest feed forward NN.⁸

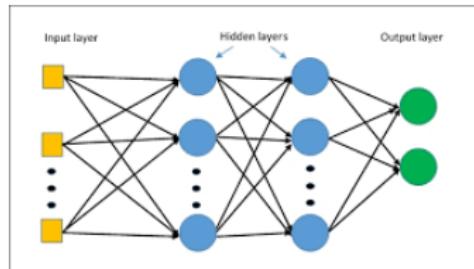
Multilayer Perceptron Model

For a MLP with inputs nodes $\vec{x}^{(0)}$, one hidden layer of nodes $\vec{x}^{(1)}$ and output layer of nodes $\vec{x}^{(2)}$, we have

$$\begin{cases} \vec{x}^{(1)} = \vec{F}^{(1)} (\vec{W}^{(1)} \cdot \vec{x}^{(0)}) \\ \vec{x}^{(2)} = \vec{F}^{(2)} (\vec{W}^{(2)} \cdot \vec{x}^{(1)}) \end{cases}$$

Eliminating the hidden layer variables \vec{H} we get

$$\Rightarrow \vec{x}^{(2)} = \vec{F}^{(2)} (\vec{W}^{(2)} \cdot \vec{F}^{(1)} (\vec{W}^{(1)} \cdot \vec{x}^{(0)})) \quad (2)$$



A MLP can be seen as a parametrization of a mapping $F_{w,b} : \mathbb{R}^n \rightarrow \mathbb{R}^m$

⁷ A MLP with m layers using linear activation functions can be reduced to a single layer !

⁸ The thresholds \vec{B} are represented as weights by redefining $\vec{W} = (B, W_1, \dots, W_n)$ and $\vec{x} = (1, x_1, \dots, x_n)$ (bias is equivalent to a weight on an extra input of activation=1)

Multilayer Perceptron as a Universal Approximator

Universal Approximation Theorem

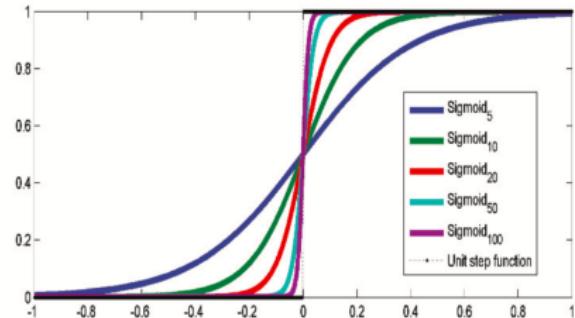
A single hidden layer feed forward neural network with a linear output unit can approximate any continuous function arbitrarily well, given enough hidden neurons⁹

The theorem doesn't tell us how many neurons or how much data is needed !

Sigmoid → Step Function

For large weight W the sigmoid turns into a step function, while B gives its offset

$$y = \frac{1}{1 + e^{-(w.x+b)}} \quad (3)$$



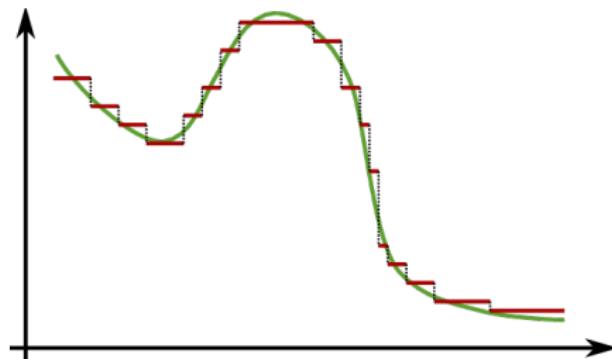
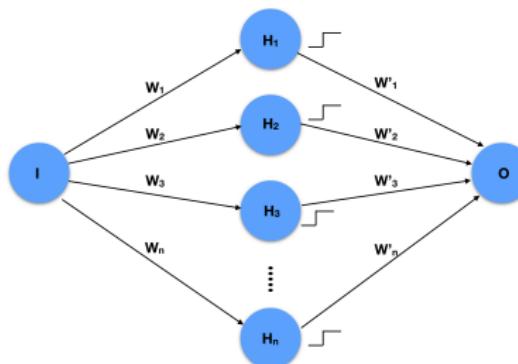
⁹Cybenko,G.(1989) Approximations by superpositions of sigmoidal functions, Math.ofCtrl.,Sig.,andSyst.,2(4),303
Hornik,K.(1991) Approximation Capabilities of Multilayer Feedforward Networks, Neural Networks, 4(2), 251

Multilayer Perceptron as a Universal Approximator

Approximating $F(x)$ by Sum of Steps

A continuous function can be approximated by a finite sum of step functions. The larger the number of steps(nodes), the better the approximation

Consider a network composed of a single input , n hidden nodes and a single output. Tune the weights such that the activations approximate steps functions with appropriate thresholds and add them together !



The same works for any activation function $f(x)$, limited when $x \rightarrow \pm\infty$. One can always tune weights W and thresholds B such that it behaves like a step function !

Learning as a Minimization Problem (Gradient Descent and Backpropagation)

Gradient Descent Method (Loss Minimization)

The learning process is a loss $L(\vec{W})$ minimization problem, where weights are adjusted to achieve a minimum output error. This minimization is usually implemented by the Gradient Descent method

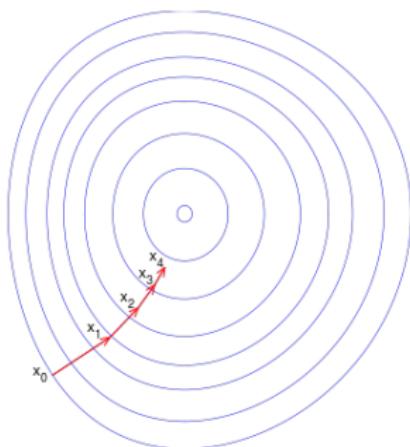
Gradient Descent

A multi-variable function $F(\vec{x})$ decreases fastest in the direction of its negative gradient $-\nabla F(\vec{x})$ ¹⁰.

From an initial point \vec{x}_0 , a recursion relation gives a sequence of points $\{\vec{x}_1, \dots, \vec{x}_n\}$ leading to a minimum

$$\vec{x}_{n+1} = \vec{x}_n - \lambda \nabla F(\vec{x}_n) \text{ , where } \lambda \text{ is the step}$$

The monotonic sequence $F(\vec{x}_0) \geq F(\vec{x}_1) \geq \dots \geq F(\vec{x}_n)$ converges to a local minimum !

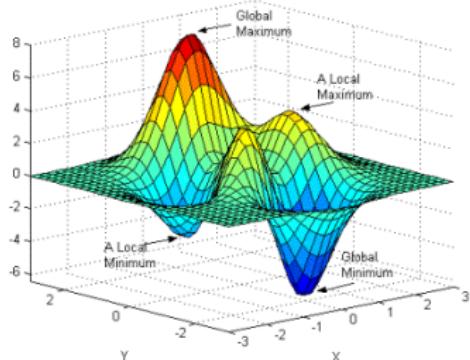


Gradient Descent uses 1^o derivatives, which is efficiently and simply calculated by backpropagation. Methods like Newton uses 2^o derivative(Hessian), which is computationally costly and memory inefficient.

¹⁰The directional derivative of $F(\vec{x})$ in the \vec{u} direction is $D_{\vec{u}} = \hat{u} \cdot \nabla F$

Stochastic Gradient Descent (SGD)

NN usually have a non-convex loss function, with a large number of local minima (**permutations of neurons in a layer leads to same loss !**)



- GD can get stuck in local minima
- Convergence issues
- Should use SGD and adaptive variants

Stochastic Gradient Descent(SGD)

SGD¹¹ selects data points randomly, instead of the order they appear in the training set. This allows the algorithm to try different "minimization paths" at each training epoch

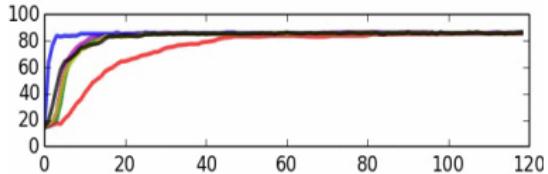
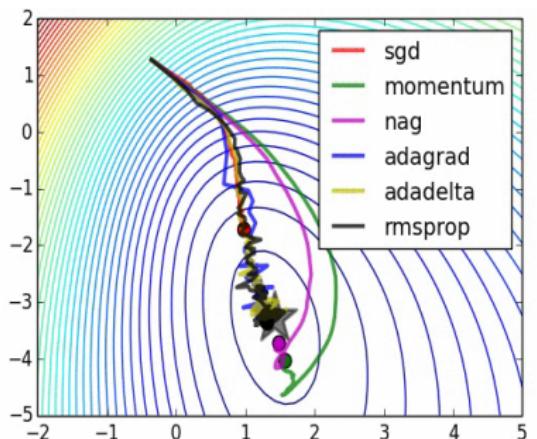
- It can also average gradient with respect to a set of events (minibatch)
- Noisy estimates average out and allows "jumping" out of bad critical points
- Scales well with dataset and model size

¹¹ <https://deeplearn.csail.mit.edu/6.S095/Fall2018/lec02.html>

SGD Algorithm Improvements¹³

SGD Variants¹² :

- **Vanilla SGD**
- **Momentum SGD** : uses update Δw of last iteration for next update in linear combination with the gradient
- **Annealing SGD** : step, exponential or $1/t$ decay
- **Adagrad** : adapts learning rate to updates of parameters depending on importance
- **Adadelta** : robust extension of Adagrad that adapts learning rates based on a moving window of gradient update
- **Rmsprop** : rescale gradient by a running average of its recent magnitude
- **Adam** : rescale gradient averages of both the gradients and the second moments of the gradients



¹²<http://danielnouri.org/notes/category/deep-learning>

¹³<http://ruder.io/optimizing-gradient-descent>

Backpropagation

Backpropagation is a technique to apply gradient descent to multilayer networks. An error at the output is propagated backwards through the layers using the chain rule¹⁴

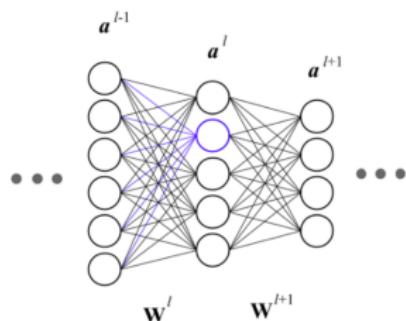
MLP Loss Function

Consider a MLP with n layers (not counting input layer) and a quadratic loss.

Defining the activation $a_i^{(l)} = F(W_{ij}^{(l)} a_j^{(l-1)})$ as the neuron output we have

$$L(\vec{W}) = \frac{1}{2} [a_i^{(n)} - t_i]^2 = \frac{1}{2} [F(W_{ij}^{(n)} \cdots F(W_{rs}^{(1)} a_s^{(0)}) \cdots) - t_i]^2$$

The loss is a n -composite function of the weights, where $W_{ij}^{(l)}$ connects the neuron j in layer $(l-1)$, to neuron i in layer l .



¹⁴<http://neuralnetworksanddeeplearning.com/chap2.html>

Backpropagation

Let's define $z_i^{(l)} = W_{ij}^{(l)} a_j^{(l-1)}$, such that $a_i^{(l)} = F(z_i^{(l)})$. The loss gradients in l -layer are

$$\frac{\partial L}{\partial W_{kj}^{(l)}} = \left[\frac{\partial L}{\partial z_k^{(l)}} \right] \underbrace{\frac{\partial z_k^{(l)}}{\partial W_{kj}^{(l)}}}_{a_j^{(l-1)}} = \left[\left(\sum_m \underbrace{\frac{\partial L}{\partial z_m^{(l+1)}}}_{W_{mk}^{(l+1)}} \underbrace{\frac{\partial z_m^{(l+1)}}{\partial a_k^{(l)}}}_{F'} \right) \underbrace{\frac{\partial a_k^{(l)}}{\partial z_k^{(l)}}}_{\frac{\partial z_k^{(l)}}{\partial W_{kj}^{(l)}}} \right] \underbrace{\frac{\partial z_k^{(l)}}{\partial W_{kj}^{(l)}}}_{a_j^{(l-1)}}$$

$$\Rightarrow \frac{\partial L}{\partial z_k^{(l)}} = \left(\sum_m \frac{\partial L}{\partial z_m^{(l+1)}} W_{mk}^{(l+1)} \right) F'(z_k^{(l)})$$

Backpropagation Formulas

The backpropagation master formulas for the gradients of the loss function in layer- l are given by

$$\boxed{\frac{\partial L}{\partial W_{kj}^{(l)}} = \delta_k^{(l)} a_k^{(l-1)}}$$

and

$$\boxed{\delta_k^{(l)} = \left(\sum_m \delta_m^{(l+1)} W_{mk}^{(l+1)} \right) F'(z_k^{(l)})}$$

, where the 'errors' in each layer are defined as $\delta_k^{(l)} = \frac{\partial L}{\partial z_k^{(l)}}$

→ The only derivative one needs to calculate is F' !

Backpropagation Algorithm

Given a training dataset $D = \{(x_i, t_i)\}$ we first run a *forward pass* to compute all the activations throughout the network, up to the output layer. Then, one computes the network output “error” and backpropagates it to determine each neuron error contribution $\delta_k^{(l)}$

Backpropagation Algorithm

- ➊ Initialize the weights $W_{kj}^{(l)}$ randomly
- ➋ Perform a feedforward pass, computing the arguments $z_k^{(l)}$ and activations $a_k^{(l)}$ for all layers
- ➌ Determine the network output error: $\delta_i^{(n)} = [a_i^{(n)} - t_i] F'(z_i^{(n)})$
- ➍ Backpropagate the output error: $\delta_k^{(l)} = (\sum_m \delta_m^{(l+1)} W_{mk}^{(l+1)}) F'(z_k^{(l)})$
- ➎ Compute the loss gradients using the neurons error and activation: $\frac{\partial L}{\partial W_{kj}^{(l)}} = \delta_k^{(l)} a_k^{(l-1)}$
- ➏ Update the weights according to the gradient descent: $\Delta W_{kj}^{(l)} = -\lambda \frac{\partial L}{\partial W_{kj}^{(l)}}$

This algorithm can be applied to other NN architectures (different connectivity patterns)

Evaluation of Learning Process

Split dataset into 3 independent parts , one for each learning phase

Training

Train(fit) the NN model by iterating through the training dataset (an epoch)

- High learning rate will quickly decay the loss faster but can get stuck or bounce around chaotically
- Low learning rate gives very low convergence speed

Validation

Check performance on independent validation dataset and tune hyper-parameters

- Evaluate the loss over the validation dataset after each epoch
- Examine for overtraining(overfitting), and determine when to stop training

Test

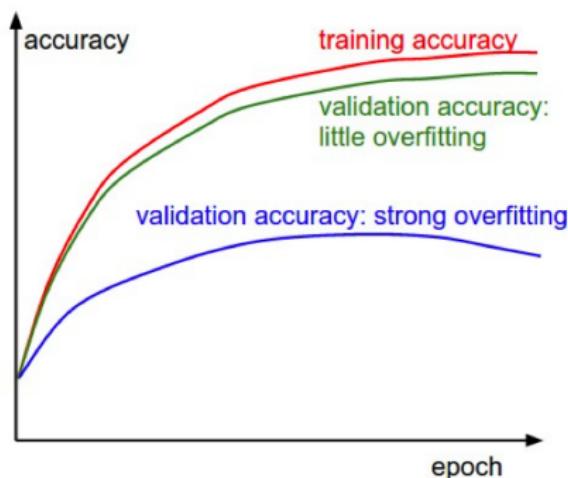
Final performance evaluation after finished training and hyper-parameters are fixed. Use the test dataset for an independent evaluation of performance obtaining a ROC curve

Overtrainging(Overfitting)

Overtraining(Overfitting)

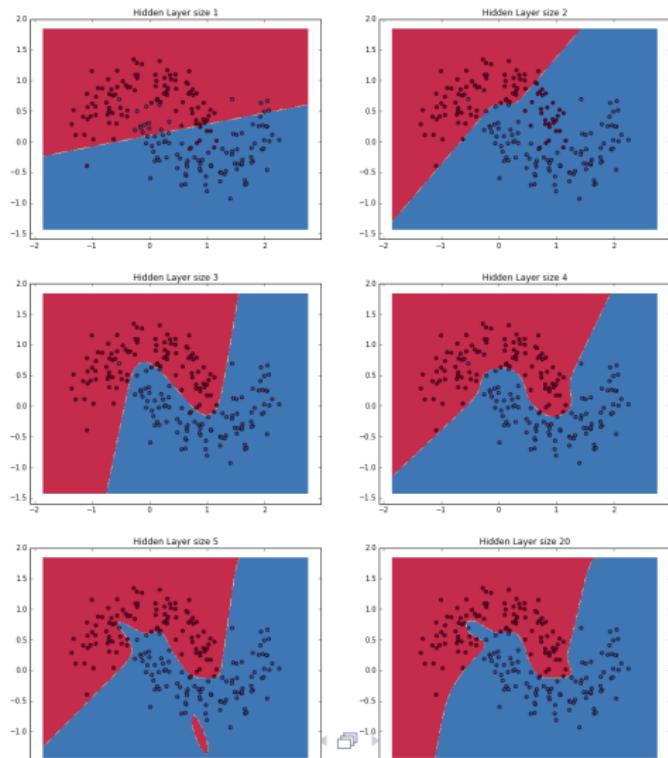
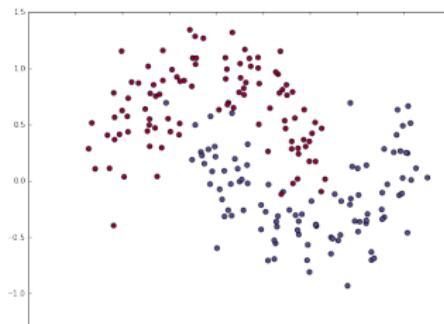
Gap between training and validation accuracy indicates the amount of overfitting

- If validation error curve shows small accuracy compared to training it indicates overfitting \Rightarrow add regularization or use more data.
- If validation accuracy tracks the training accuracy well, the model capacity is not high enough \Rightarrow use larger model



Overfitting - Hidden Layer Size X Decision Boundary ¹⁵

- Hidden layer of low dimensionality nicely captures the general trend of data.
- Higher dimensionalities are prone to overfitting (“memorizing” data) as opposed to fitting the general shape
- If evaluated on independent dataset (and you should !), the smaller hidden layer generalizes better
- Can counteract overfitting with regularization, but picking correct size for hidden layer is much simpler



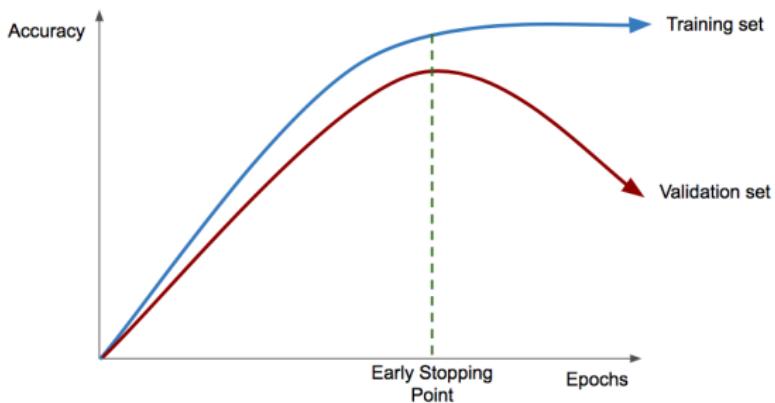
Regularization

Regularization techniques prevent the neural network from overfitting

Early Stopping

Early stopping can be viewed as regularization in time. Gradient descent will tend to learn more and more the dataset complexities as the number of iterations increases.

Early stopping is implemented by training just until performance on the validation set no longer improves or attained a satisfactory level. Improving the model fit to the training data comes at the expense of increased generalization error.

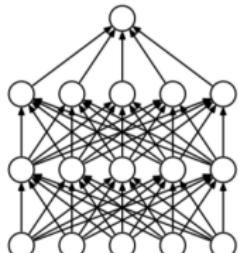


Dropout Regularization

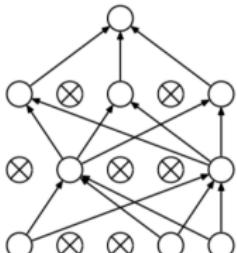
Dropout Regularization

Regularization inside network that remove nodes randomly during training.

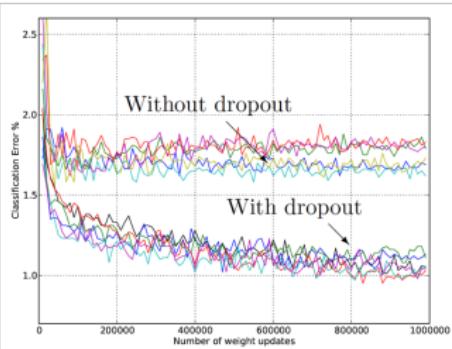
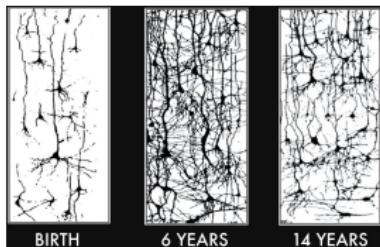
- It's ON in training and OFF in validation and testing
- Avoid co-adaptation on training data
- Essentially a large model averaging procedure



(a) Standard Neural Net



(b) After applying dropout.



Usually worsens the results during training, but improves validation and testing results !

L1 & L2 Regularization

Between two models with the same predictive power, the 'simpler' one is to be preferred (NN Occam's razor)

L1 and L2 regularizations add a term to the loss function that tames overfitting

$$L'(\vec{w}) = L(\vec{w}) + \alpha\Omega(\vec{w})$$

L1 Regularization

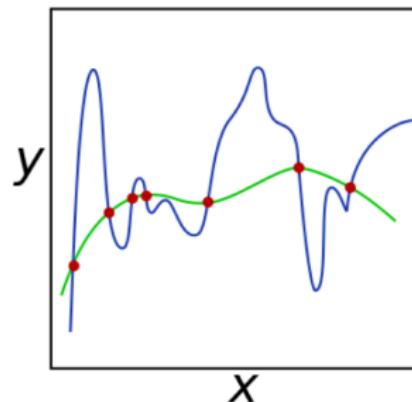
Works by keeping weights sparse

$$\Omega(\vec{w}) = |\vec{w}|$$

L2 Regularization

Works by penalising large weights

$$\Omega(\vec{w}) = |\vec{w}|^2$$

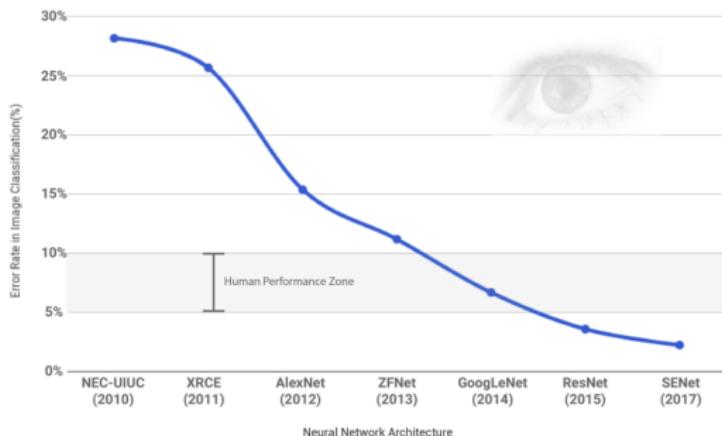


The combined $L1 + L2$ regularization is called Elastic

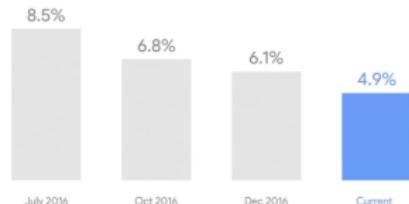
Why Deep Learning ? and Why Now ?

Why Deep Learning ?

Image and Speech Recognition performance (DNN versus Humans)



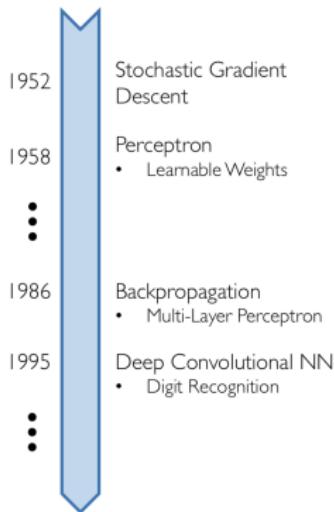
Speech Recognition
Word Error Rate



<https://arxiv.org/pdf/1409.0575.pdf>

Why Now ?

Neural networks date back decades , so why the current resurgence ?



The main catalysts for the current Deep Learning revolution have been:

- **Software:**
TensorFlow, PyTorch, Keras and Scikit-Learn
- **Hardware:**
GPU, TPU and FPGA
- **Large Datasets:**
MNIST

Training Datasets

Large and new open source datasets for machine learning research¹⁶



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

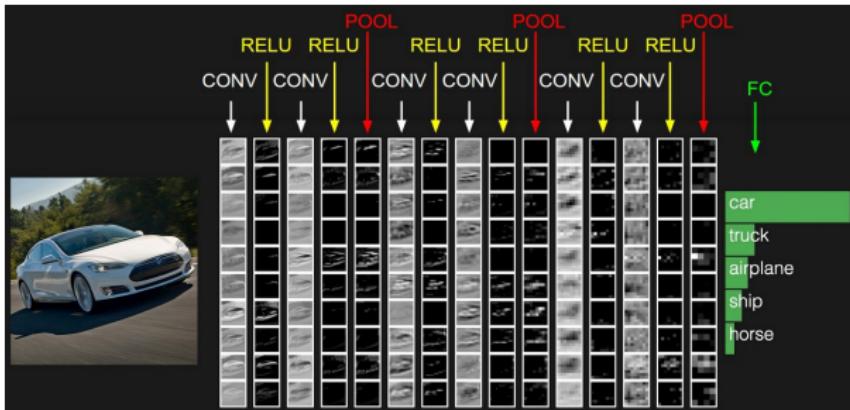


¹⁶https://en.wikipedia.org/wiki/List_of_datasets_for_machine_learning_research

Deep Learning - Need for Depth

Deep Neural Networks(DNN)

Depth allows the NN to factorize the data features, distributing its representation across the layers, exploring the compositional character of nature¹⁷



⇒ DNN allows a hierarchical representation of data features !

¹⁷ Deep Learning , Y.LeCunn, J.Bengio, G.Hinton , Nature , vol. 521, pg. 436 , May 2015

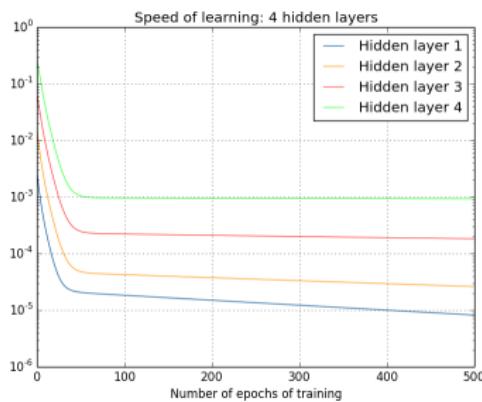
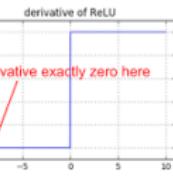
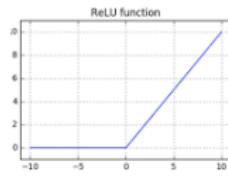
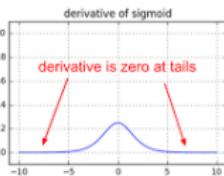
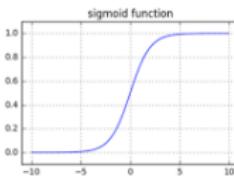
Deep Learning - Vanishing Gradient Problem

Vanishing Gradient Problem ¹⁸

Backpropagation computes gradients iteratively by multiplying the activation function derivate F' through n layers.

$$\delta_k^{(l)} = \left(\sum_m \delta_m^{(l+1)} W_{mk}^{(l+1)} \right) F'(z_k^{(l)})$$

For $\sigma(x)$ and $Tanh(x)$ the derivate F' is asymptotically zero, so weights updates gets vanishing small when backpropagated \Rightarrow Then, earlier layers learns much slower than later layers !!!

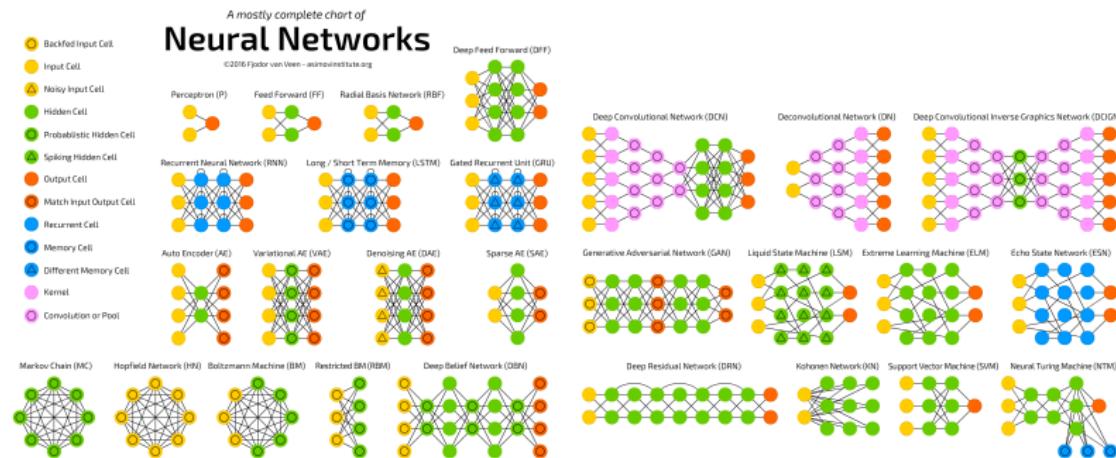


¹⁸ <http://neuralnetworksanddeeplearning.com/chap5.html>

Deep Architectures and Applications

Neural Network Zoo

NN architecture and nodes connectivity can be adapted for the problem at hand¹⁹



¹⁹ <http://www.asimovinstitute.org/neural-network-zoo>

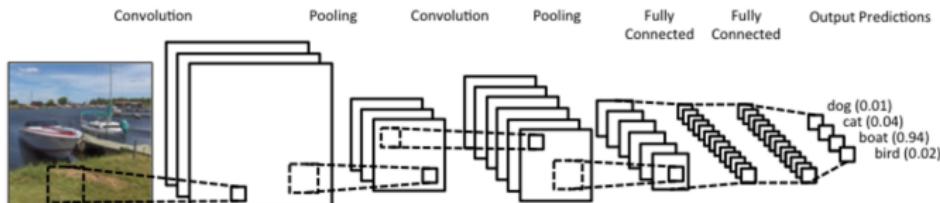
Convolutional Neural Network

For high-dimensional inputs as images, fully connected NN is computationally too expensive (ex: 1000x1000 image has $O(10^6)$ pixels). Convolutional neural network²⁰ emulate the visual cortex, where neurons respond to stimuli only in a restricted region of the visual field.

Convolutional Neural Network (CNN or ConvNet)

CNN mitigates the challenges of high dimensional inputs by restricting the connections between the input and hidden neurons. It connects only a small contiguous region of input nodes, exploiting local correlation. Its architecture is a stack of distinct and specialized layers:

- ① Convolutional
- ② Pooling (Downsampling)
- ③ Fully connected (MLP)



²⁰<http://ufldl.stanford.edu/tutorial/supervised/FeatureExtractionUsingConvolution>

CNN - Convolutional Layer

Convolutional Layer

CNN convolutional layer²¹ is its core building block and its parameters are learnable filters (kernels), that extracts image features corresponding to a small receptive field.

A convolution is like a sliding window transformation that applies a filter to extract local image features

Discrete Convolution

$$y_{ij}^{(l+1)} = \sum_{a=0}^{n-1} \sum_{b=0}^{m-1} w_{ab} x_{(i+a)(j+b)}^{(l)}$$

1 x ₁	1 x ₀	1 x ₁	0	0
0 x ₀	1 x ₁	1 x ₀	1	0
0 x ₁	0 x ₀	1 x ₁	1	1
0	0	1	1	0
0	1	1	0	0

Image

4		

Convolved Feature

²¹ <http://ufldl.stanford.edu/tutorial/supervised/FeatureExtractionUsingConvolution>

<https://machinelearningmastery.com/padding-and-stride-for-convolutional-neural-networks>

Convolution Filters

An image convolution (filter) can apply an effect (sharpen, blurr), as well as extract features (edges, texture)²²

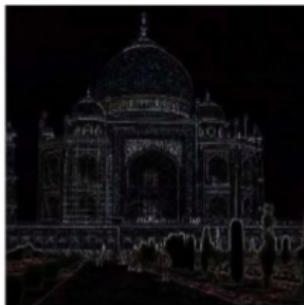
$$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & -1 & 5 & -1 & 0 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$



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



$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$



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



²²<https://docs.gimp.org/2.6/en/plug-in-convmatrix.html>

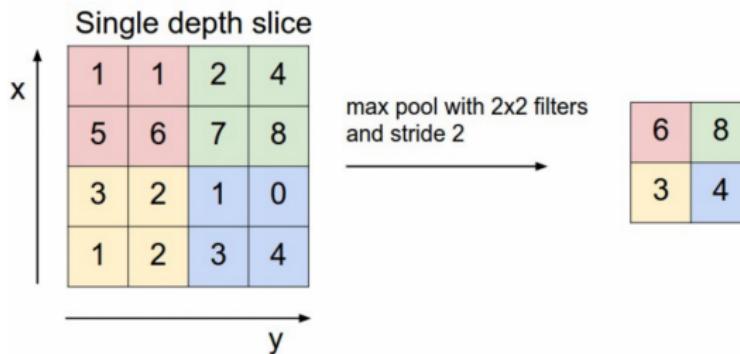
CNN - Pooling Layer

Pooling(Downsampling) Layer

The pooling (downsampling) layer²³ has no learning capabilities and serves a dual purpose:

- Decrease the representation size \Rightarrow reduce computation
- Make the representation approximately invariant to small input translations and rotations

Pooling layer partitions the input in non-overlapping regions and, for each sub-region, it outputs a single value (ex: max pooling, mean pooling)



²³<http://ufldl.stanford.edu/tutorial/supervised/Pooling>

CNN - Fully Connected Layers

Fully Connected Layer

CNN chains together convolutional(filtering) , pooling (downsampling) and then fully connected(MLP) layers.

- After processing with convolutions and pooling, use fully connected layers for classification
- Architecture allows capturing local structure in convolutions, and long range structure in later stage convolutions and fully connected layers

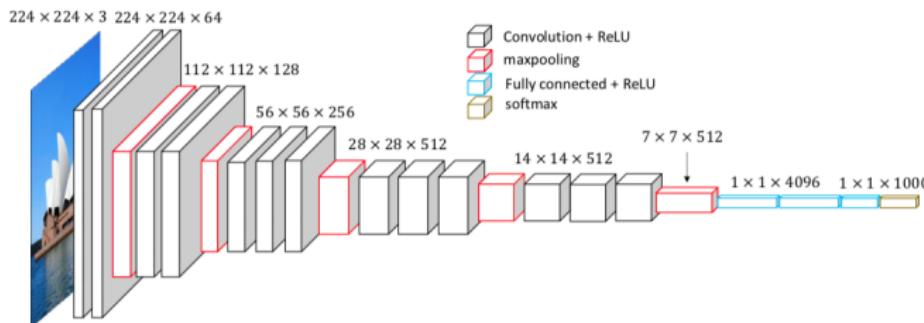
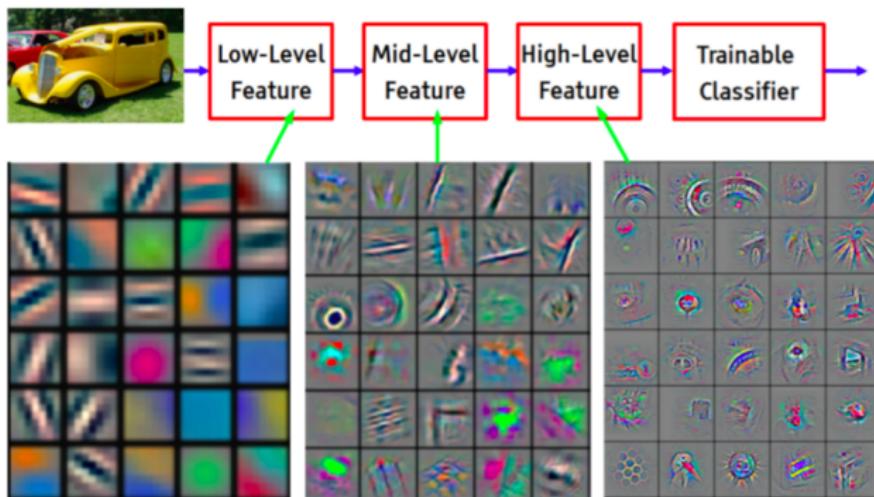


Figure 2: The architecture of VGG16 model .

CNN Feature Visualization

Each CNN layer is responsible for capturing a different level of features as can be seen from ImagiNet²⁴

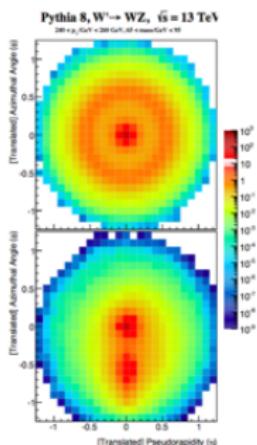
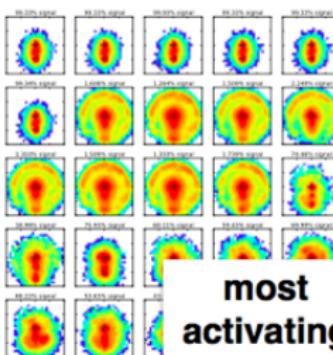
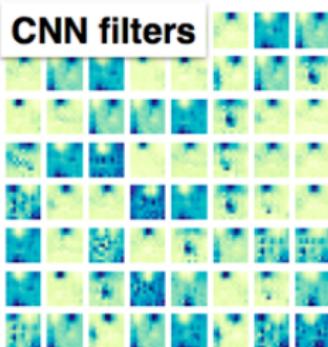


Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

²⁴<https://arxiv.org/pdf/1311.2901.pdf>

CNN Application in HEP: Jet ID

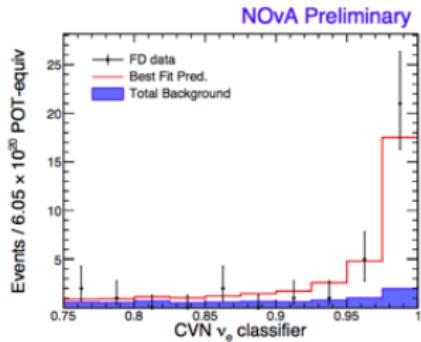
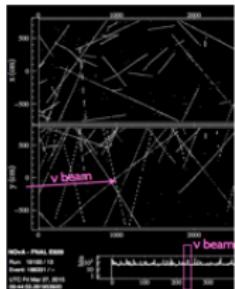
Jet images with convolutional nets



L. de Oliveira et al., 2015

CNN Application HEP: Neutrino ID

Neutrinos with convolutional nets

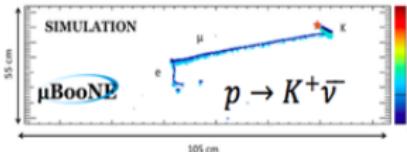
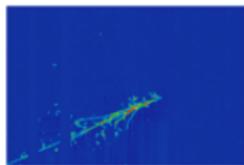


76% Purit
73% Effici

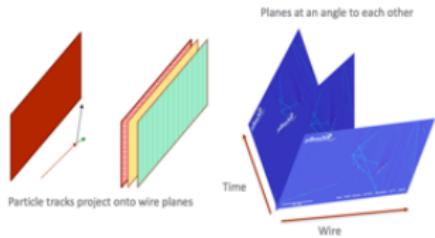


An equivalent increased exposure of 30%

Aurisiano et al. 2016



μBooNE



Sequential Data

Sequential Data

Sequential data is an interdependent data stream (ex: text , audio and video)

The food was good, not bad at all.

vs.

The food was bad, not good at all.

Sequences and Intelligence

- Sequences are very important for intelligence
- Brain is all the time predicting what comes ahead ...
- If prediction is wrong the brain learns and corrects
- We memorize sequences for alphabet, words, phone numbers and not just symbols !

Feed forward networks can't learn correlation between previous and current input !

Recurrent Neural Network(RNN)

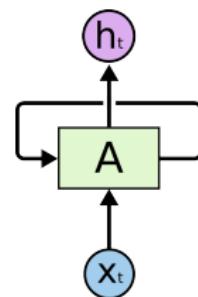
Recurrent Neural Networks (RNN)

RNNs are networks that use feedback loops to process sequential data²⁵. These feedbacks allows information to persist, which is an effect often described as memory.

RNN Cell (Neuron²⁶)

The hidden state depends not only on the current input, but also on the entire history of past inputs

- ① **Hidden State:** $h^{[t]} = F(W_{xh}x^{[t]} + W_{hh}h^{[t-1]})$
- ② **Output:** $y^{[t]} = W_{hy}h^{[t]}$



²⁵<https://eli.thegreenplace.net/2018/understanding-how-to-implement-a-character-based-rnn/>

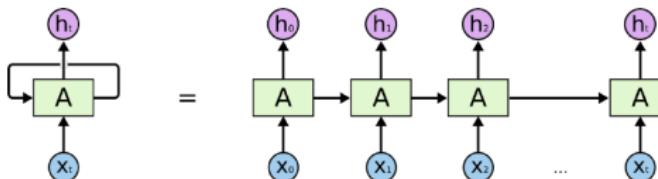
²⁶<https://www.analyticsvidhya.com/blog/2017/12/introduction-to-recurrent-neural-networks/>

Recurrent Neural Network(RNN)

RNNs²⁷ process an input sequence element at a time, maintaining in their hidden units a ‘state vector’ that contains information about all the past elements history.

RNN Unrolling

A RNN can be thought of as multiple copies of the same network, each passing a message to a successor. Unrolling is a visualization tool which views a RNN as a sequence of unit cells.



Backpropagation Through Time (BPTT)

Backpropagation through time is just a fancy buzz word for backpropagation on an unrolled RNN

Unrolled RNN can lead to very deep networks \Rightarrow bias to capture only short term dependencies !

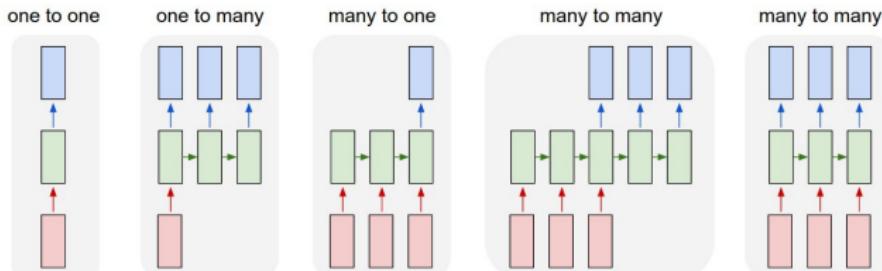
²⁷ <http://www.wildml.com/2015/09/recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/>

Recurrent Neural Network(RNN)

RNN Properties

- It can take as input variable size data sequences
- RNN allows one to operate over sequences of vectors in the input, the output or both

Bellow, input vectors are in red, output vectors are in blue and green vectors hold the RNN's state²⁸

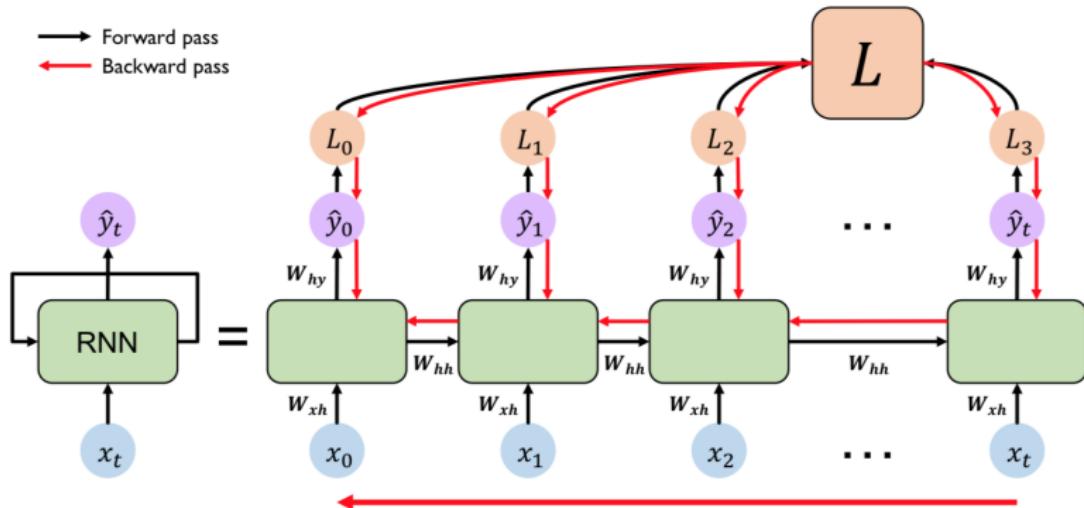


There are no constraints on the sequences lengths because the recurrent transformation can be applied as many times as necessary

²⁸<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

Recurrent Neural Network(RNN)

The RNN computational graph and information flow



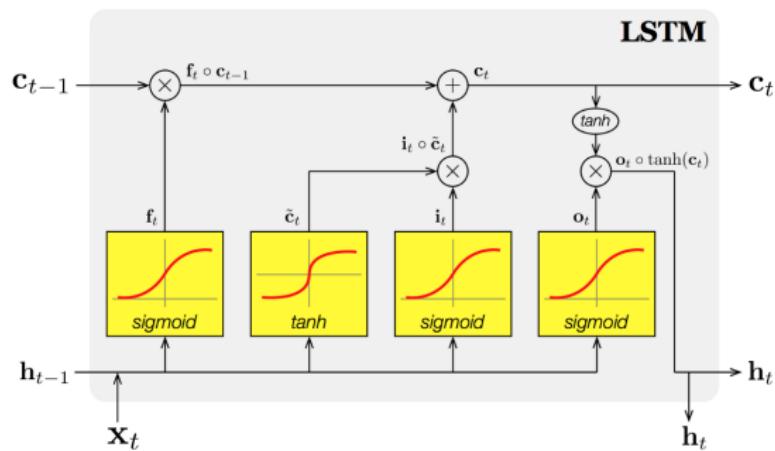
The same set of weights are used to compute the hidden state and output for all time steps

Long Short Term Memory(LSTM)

Long Short-Term Memory (LSTM) Network

LSTM²⁹ is a gated RNNs capable of learning long-term dependencies. It categorize data into **short** and **long** term, deciding its importance and what to remember or forget.

The LSTM unit cell has an **input** and a **forget** gate. The **input** defines how much of the newly computed state for the current input is accepted, while the **forget** defines how much of the previous state is accepted.



Gating variables

$$f_t = \sigma(\mathbf{W}_f[h_{t-1}, x_t] + \mathbf{b}_f)$$

$$i_t = \sigma(\mathbf{W}_i[h_{t-1}, x_t] + \mathbf{b}_i)$$

$$o_t = \sigma(\mathbf{W}_o[h_{t-1}, x_t] + \mathbf{b}_o)$$

Candidate (memory) cell state

$$\tilde{c}_t = \tanh(\mathbf{W}_c[h_{t-1}, x_t] + \mathbf{b}_c)$$

Cell & Hidden state

$$c_t = f_t \circ c_{t-1} + i_t \circ \tilde{c}_t$$

$$h_t = o_t \circ \tanh(c_t)$$

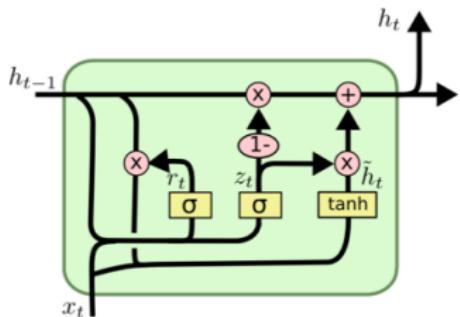
²⁹<http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Gated Recurrent Network (GRU)

Gated Recurrent Network (GRU)

GRU³⁰ networks have been proposed as a simplified version of LSTM, which also avoids the vanishing gradient problem and is even easier to train.

In a GRU the **reset** gate determines how to combine the new input with the previous memory, and the **update** gate defines how much of the previous memory is kept



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

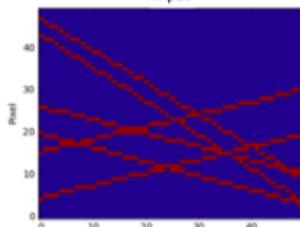
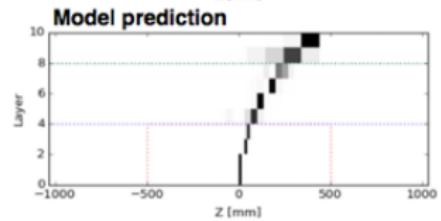
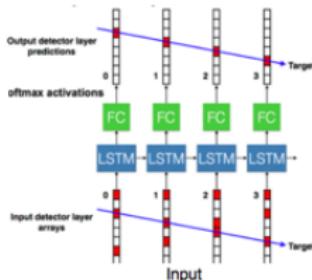
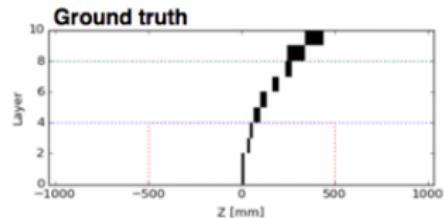
$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

³⁰<https://isaacchanghau.github.io/post/lstm-gru-formula>

LSTM Application in HEP: Tracking

Tracking with recurrent neural networks (LSTM)³¹

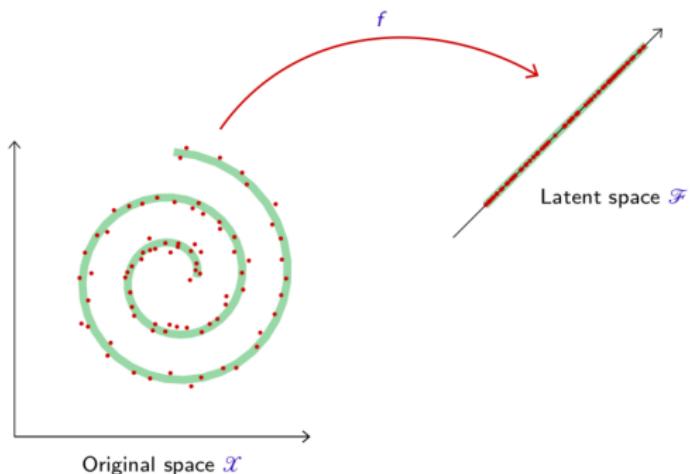


**Time dimension
(state memory)**

³¹ <https://heptrkx.github.io>

Autoencoder

Many applications such as data compression, denoising and data generation require to go beyond classification and regression problems. This modeling usually consists of finding “meaningful degrees of freedom”, that can describe high dimensional data in terms of a smaller dimensional representation



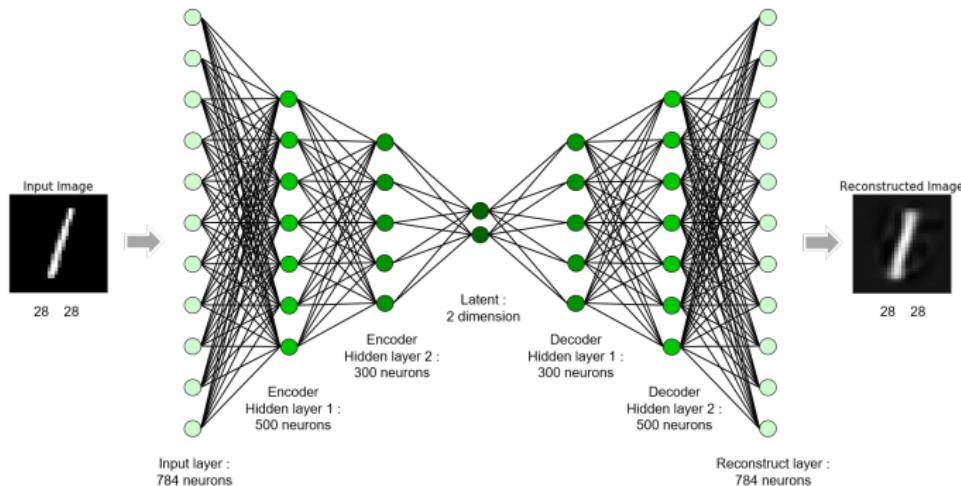
Traditionally, autoencoders were used for dimensionality reduction and denoising. Recently autoencoders are being used also in generative modeling

Autoencoder(AE)

Autoencoder(AE)

An AE is a neural network that is trained to attempt to copy its input to its output in an **self-supervised way**. In doing so, it learns a representation(encoding) of the data set features / in a low dimensional latent space.

It may be viewed as consisting of two parts: an encoder $I = f(x)$ and a decoder $y = g(I)$.

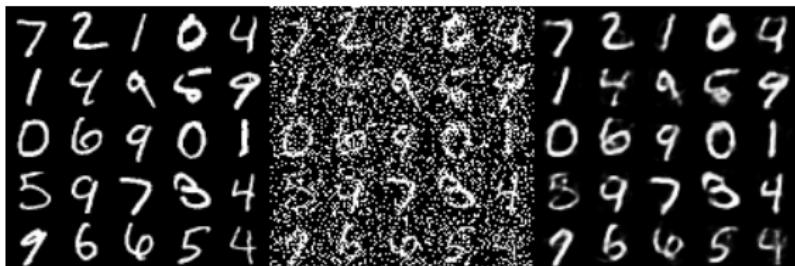


⇒ Understanding is data reduction

Denoising Autoencoder (DAE)

Denoising Autoencoder (DAE)

The DAE is an extension of a classical autoencoder where one corrupts the original image on purpose by adding random noise to it's input. The autoencoder is trained to reconstruct the input from a corrupted version of it and then used as a tool for noise extraction

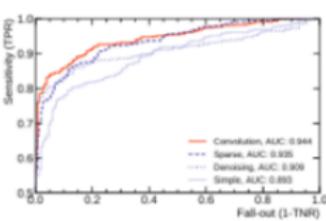
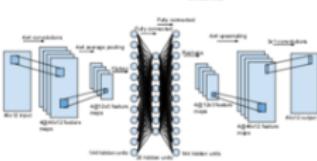
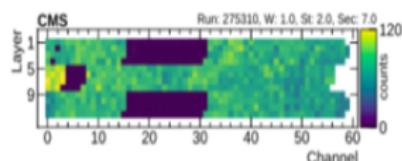
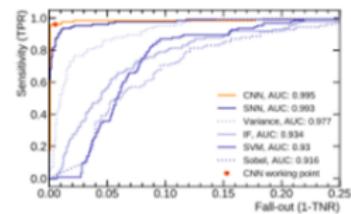
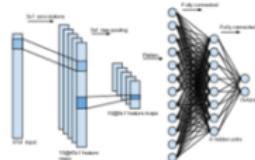
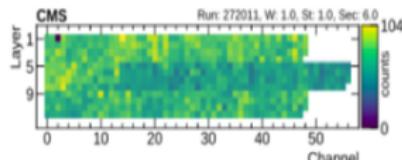


Autoencoder Application in HEP: DQM

AE can be used for anomaly detection by training on a single class , so that every anomaly gives a large reconstruction error

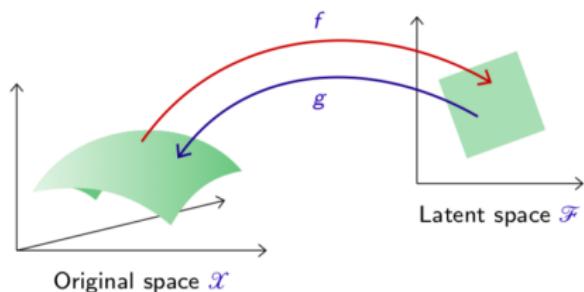
Detector Quality Monitoring (DQM)

Monitoring the CMS data taking to spot failures (anomalies) in the detector systems



Generative Autoencoder

An autoencoder combines an encoder f from the original space \mathcal{X} to a latent space \mathcal{F} , and a decoder g to map back to \mathcal{X} , such that the composite map $g \circ f$ is close to the identity when evaluated on data.



Autoencoder Loss Function

$$L = \| X - f \circ g(X) \|^2$$

Autoencoder as a Generator

One can train an AE on images and save the encoded vector to reconstruct (generate) it later by passing it through the decoder. The problem is that two images of the same number (ex: 2 written by different people) could end up far away in latent space !

Variational Autoencoder(VAE)

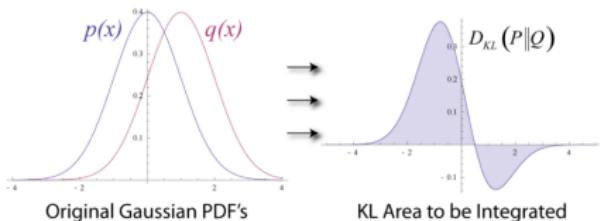
Variational Autoencoder(VAE)

A VAE³² is an autoencoder with a loss function penalty Kullback–Leibler (KL)divergence that forces it to generate latent vectors that follows a unit gaussian distribution. To generate images with a VAE one just samples a latent vector from a unit gaussian and pass it through the decoder.

The KL divergence, or relative entropy, is a measure of how one probability distribution differs from a second, reference distribution

Kullback–Leibler Divergence (Relative Entropy)

$$D_{KL}(p||q) = \int_{-\infty}^{+\infty} dx p(x) \log \left(\frac{p(x)}{q(x)} \right)$$



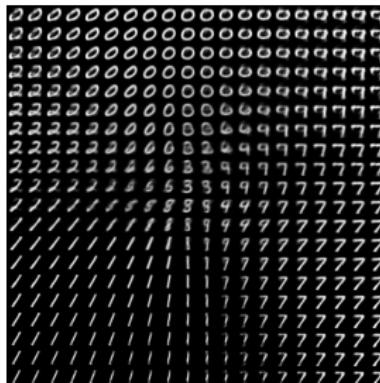
Variational Autoencoder(VAE)

VAE Loss

The VAE loss ³³ function is composed of a mean squared error (generative) loss that measures the reconstruction accuracy, and a KL divergence (latent) loss that measures how close the latent variables match a gaussian.

$$L = \| x - f \circ g(x) \|^2 + D_{KL}(p(z|x) | q(z|x))$$

Bellow we have an example of a set of a VAE generated numbers obtained by gaussian sampling a 2D latent space

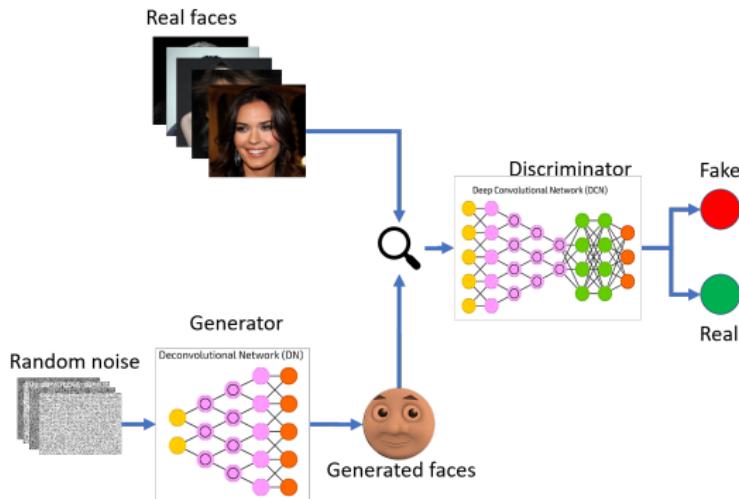


³³<https://tiao.io/post/tutorial-on-variational-autoencoders-with-a-concise-keras-implement/>

Generative Adversarial Network(GAN)

Generative Adversarial Networks (GAN)

GANs are composed by two NN, where one generates candidates and the other classifies them. The generator learns a map from a latent space to data, while the classifier discriminates generated data from real data.

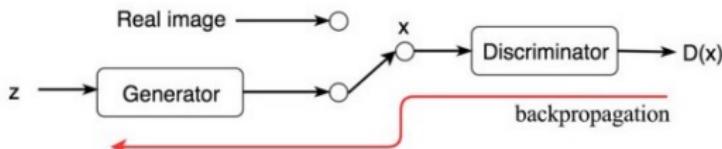


Generative Adversarial Network(GAN)

GAN adversarial training ³⁴ works by the two neural networks competing and training each other. The generator tries to "fool" the discriminator, while the discriminator tries to uncover the fake data.

GAN Training

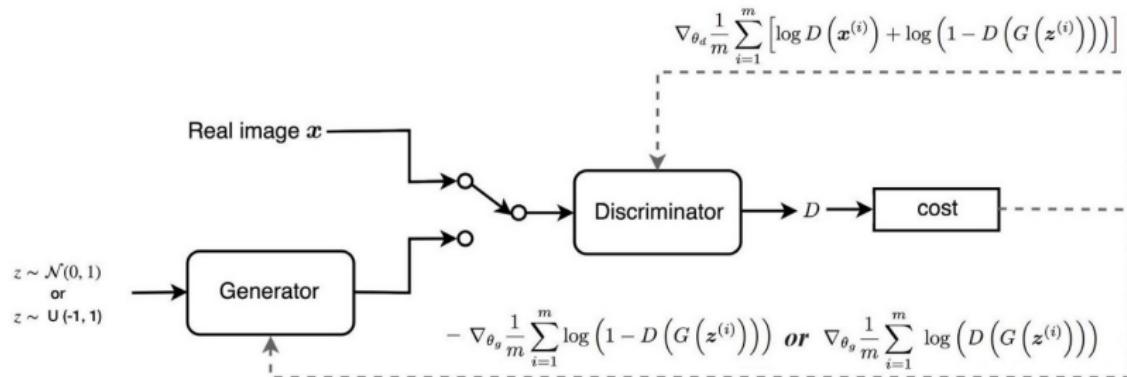
- ① The discriminator receives as input samples synthesized by the generator and real data . It is trained just like a classifier, so if the input is real, we want output=1 and if it's generated, output=0
- ② The generator is seeded with a randomized input that is sampled from a predefined latent space (ex: multivariate normal distribution)
- ③ We train the generator by backpropagating this target value all the way back to the generator
- ④ Both networks are trained in alternating steps and in competition



³⁴https://medium.com/@jonathan_hui/gan-whats-generative-adversarial-networks-and-its-applications-1f7d8c3e3d0d

Generative Adversarial Network(GAN)

GAN training algorithm is illustrated in more detail by the diagram below³⁵



³⁵https://medium.com/@jonathan_hui

Generative Adversarial Network(GAN)

GANs are quite good on faking celebrities images³⁶ or Monet style paintings³⁷ !



Training Data



Sample Generator



³⁶https://research.nvidia.com/publication/2017-10_Progressive-Growing-of

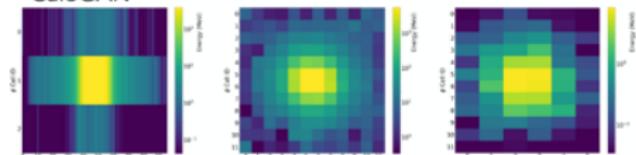
³⁷<https://github.com/junyanz/pytorch-CycleGAN-and-pix2pix>

GAN Application in HEP: MC Simulation

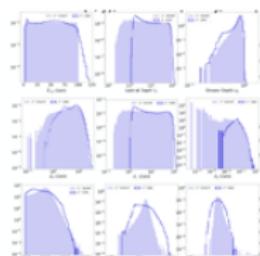
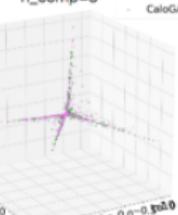
CaloGAN

Simulating 3D high energy particle showers in multi-layer electromagnetic calorimeters with a GAN³⁸

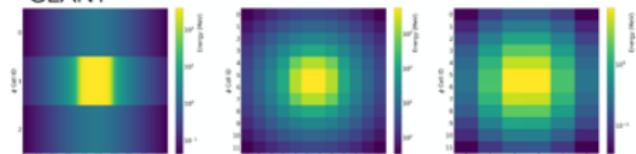
- CaloGAN



kernel=poly
n_comp=3



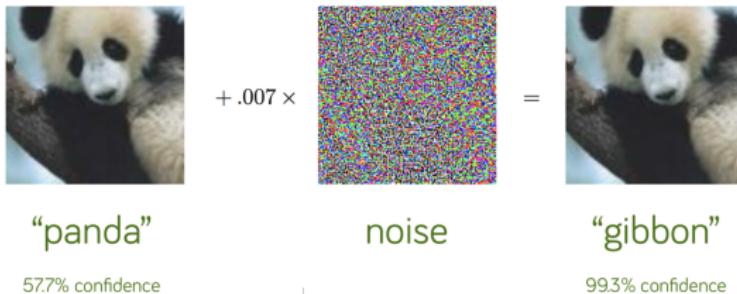
- GEANT



³⁸<https://github.com/hep-lbdl/CaloGAN>

Adversarial Attacks

In 2014 researchers at Google and NYU found that it was far too easy to fool ConvNets



This is like an optical illusion for the neural network. Our brain can clearly tell that both the images look like pandas !

Small perturbation in individual pixels can cause a dramatic change in the NN dot products, leading to a 'point' in high dimensional input space that our networks have never seen before. This space is very sparse and data is concentrated in small regions. ReLu has a non-zero gradient everywhere to the right of 0, making NN more stable and faster to train. That also makes it possible to push the ReLu activation function to arbitrarily high values, leading to a trade-off between trainability and robustness to adversarial attacks.



Residual Neural Network (ResNet)

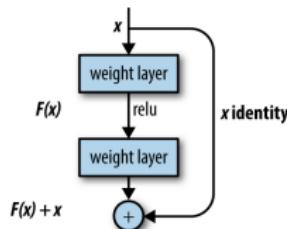
Training very deep networks is difficult because of the vanishing gradients problem. ResNets³⁹ try avoid this by reusing activations from a previous layer until the layer next to the current one have learned its weights.

Residual Neural Network (ResNet)

A Residual Neural Network is a network inspired on the pyramidal cells of cerebral cortex. The network skip connections or short-cuts to jump over some , using as inputs not only the previous neuron output, but also it's input.

ResNet Neurons

$$a^{(l)} = F(z_i^{(l)}) + a^{(l-2)}, \text{ where } z^{(l)} = W^{(l-1,l)} a^{(l-1)}$$



In the initial training stage the weights will adapt to mute layers, collapsing the network into fewer layers and making it easier to learn. Then as it learns it gradually expands the layers.

³⁹ <https://towardsdatascience.com/an-overview-of-resnet-and-its-variants-528fe2f56035>

Additional Topics:

- Hyperparameter Optimization: Grid search , Random search and Bayesian Optimization
- New architectures: Capsule Networks , Graph Neural Networks
- New computing frameworks: GPU , TPU , FPGA (Amazon EC2)
- Brain reverse engineering
- What's the NN learning ? Debuging a CNN ?
 - Which pixels of a picture most contributed to a prediction?
 - Which pixels were in contradiction to prediction?
 - "Sensitivity analysis and heatmaps (based on "input gradient" df/dx_i)".

THE END

Machine Learning Software

- General ML library (Python):

① <https://scikit-learn.org/stable>

- Deep learning libraries:

① <https://www.tensorflow.org> (TensorFlow)

② <https://pytorch.org> (PyTorch)

③ <https://www.microsoft.com/en-us/cognitive-toolkit> (CNTK)

- High level deep learning API:

① <https://keras.io> (Keras)

② <https://docs.fast.ai> (FastAI)

Bibliography

- Neural Networks and Deep Learning, M.Nielsen :
<http://neuralnetworksanddeeplearning.com>
- EPFL EE559 class, F.Fleuret : <https://fleuret.org/ee559/>
- Stanford University CS231 class, A.Karpathy : <http://cs231n.stanford.edu>
- Stanford University CS229 class, A.Ng : <http://cs229.stanford.edu>
- Introduction to machine learning, Murray:
http://videolectures.net/bootcamp2010_murray_iml
- Machine Learning HEP School, A.Rogozhnikov:
<https://indico.cern.ch/event/497368>
- Machine Learning , M.Kagan: <https://indico.cern.ch/event/726959>
- INSIGHTS Workshop at CERN, S.Gleizer:
<https://indico.cern.ch/event/747653/contributions>
- Pisa School on Future Colliders, S.Valecrosa: <https://indico.cern.ch/event/669093>
- VBSCOST Ljubljana ML Training, M.Pierini:
<https://indico.cern.ch/event/775229/contributions>
- ML4HEP at Univ.Zurich, G.Kasieczka:
<https://indico.cern.ch/event/757837/contributions>
- MIT Introduction do Deep Learning Course: <http://introtodeeplearning.com>