

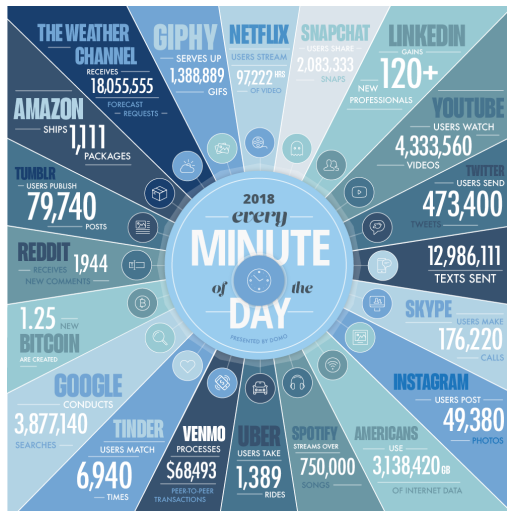
# Deep learning: an introduction

Gilles Gasso

INSA Rouen - ITI Department  
Laboratory LITIS

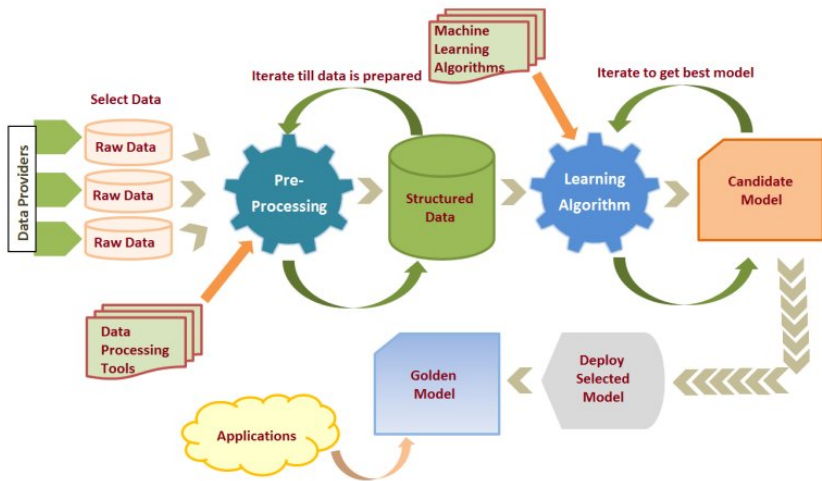
December 8, 2025

# From Data...

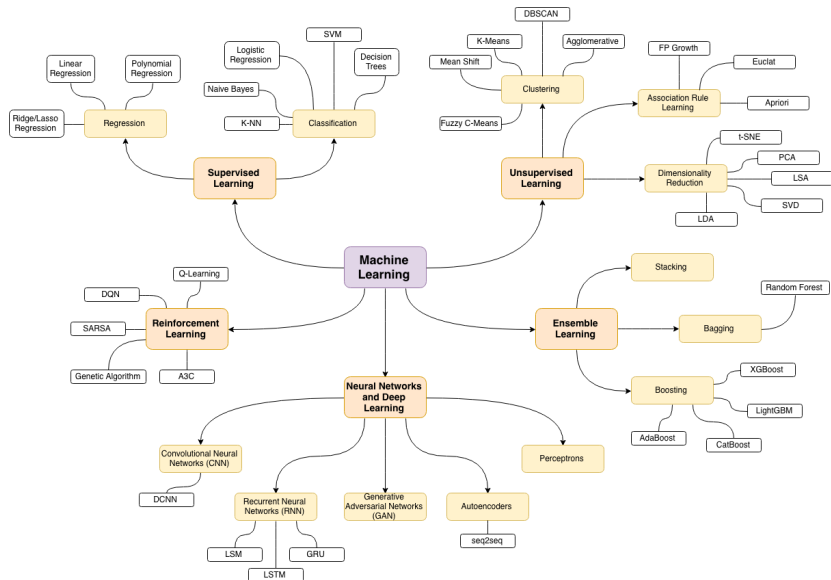


- Text
- Audio
- Images
- Videos
- Graphs ...

to processing...



# using algorithms

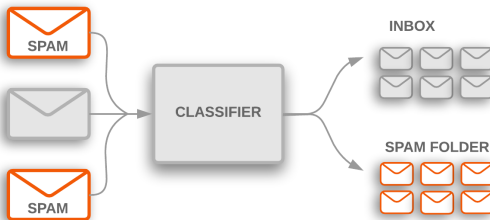


# Complex, structured data

Image classification : bus  $y = 1$  vs train  $y = 0$

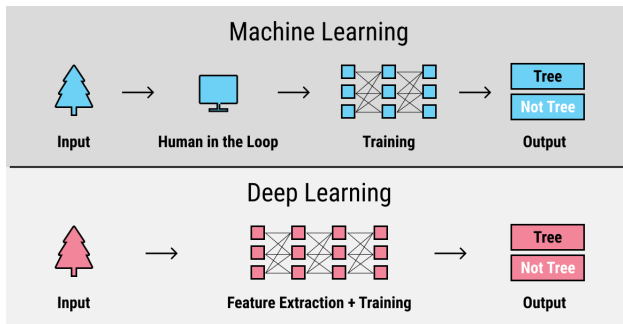


Text classification : spam vs non-spam



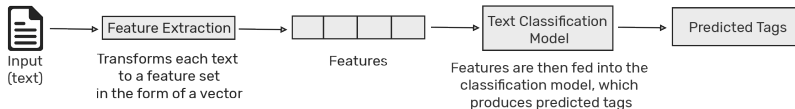
# Machine Learning vs Deep Learning

## Image

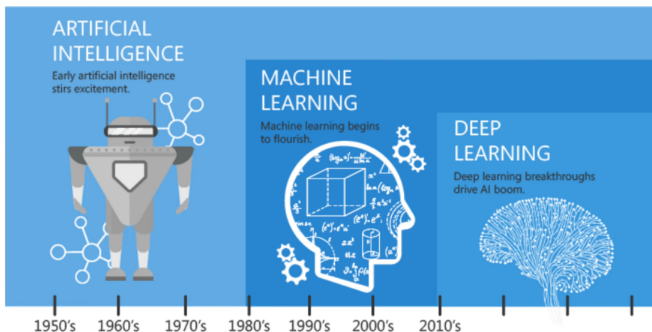


<https://labelyourdata.com/articles/machine-learning-and-training-data>

## Text



# Machine Learning vs Deep Learning



Today's IA is Deep Learning (a technique of Machine Learning)

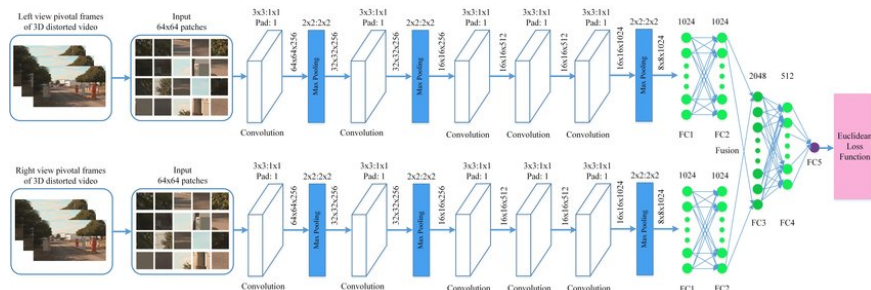
<https://blog.alore.io/machine-learning-and-artificial-intelligence/>

- Deep Learning: nowadays most popular paradigm in machine learning
- Its rise dates back to 2006 (ImageNet dataset challenge)
- Has introduced a paradigm shift in the way one will exploit data

# Deep Learning paradigm

## End-to-end learning

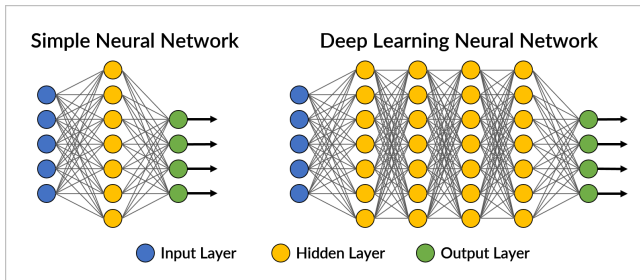
- Learn automatically (and simultaneously) the data representation  $\Phi(x)$  and the decision function  $f$



# Deep learning paradigm

## Principle

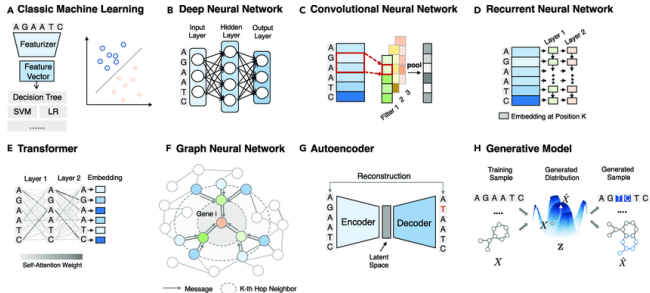
- Neural networks: models with many layers to learn hierarchical representations of data
- Composition of simple functions  $f = f_1 \circ f_2 \circ \dots \circ f_L$
- Learning based on data  $\mathcal{D} = \{(\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^N$
- Scales with data and compute



<https://www.go-rbcs.com/columns/deep-learning-to-the-rescue>

# Various deep learning models I

- Multilayer perceptron
- Convolutional networks
- Recurrent networks
- Transformers and Attention
- Diffusion models...

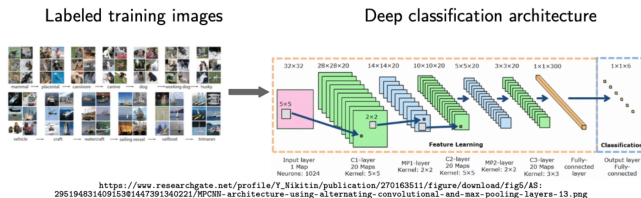


<https://www.researchgate.net/figure/illustrations-of-machine-learning-models-Details-about-each-model-can-be-found->

in\_fig3\_353783898

# Various deep learning models II

## Deep learning for computer vision tasks



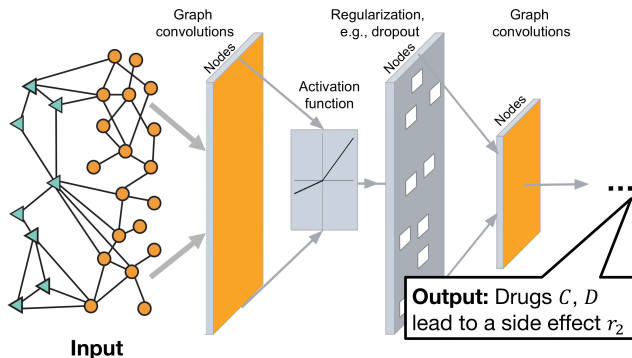
## Input images and predicted category



<https://adriancolyer.files.wordpress.com/2016/04/imagenet-fig41.png?w=656>

# Various deep learning models III

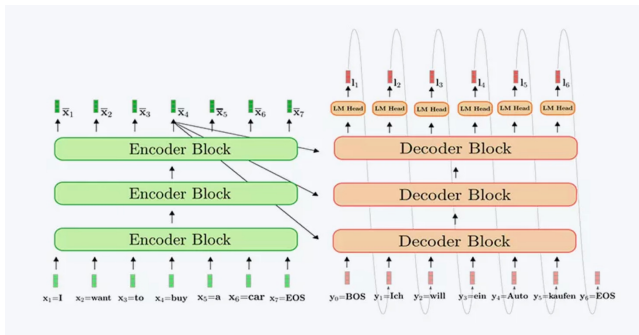
## Deep learning for graphs



<http://snap.stanford.edu/decagon/decagon-overview.png>

# Various deep learning models IV

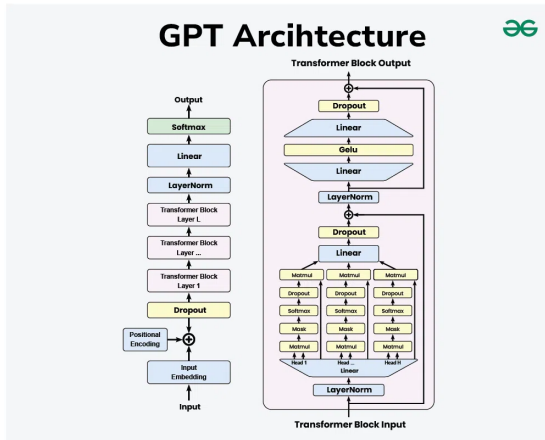
## Deep learning for texts



[https://raw.githubusercontent.com/patrickvonplaten/scientific\\_images/master/encoder\\_decoder/EncoderDecoder.png](https://raw.githubusercontent.com/patrickvonplaten/scientific_images/master/encoder_decoder/EncoderDecoder.png)

## Various deep learning models V

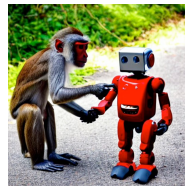
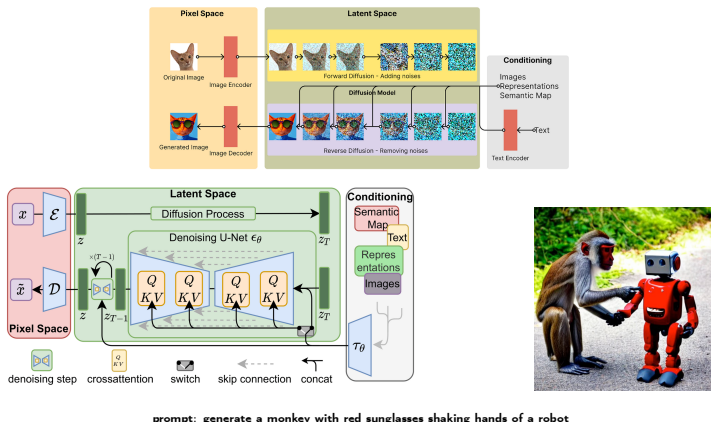
## Deep learning for texts



<https://www.geeksforgeeks.org/artificial-intelligence/introduction-to-generative-pre-trained-transformer-gpt/>

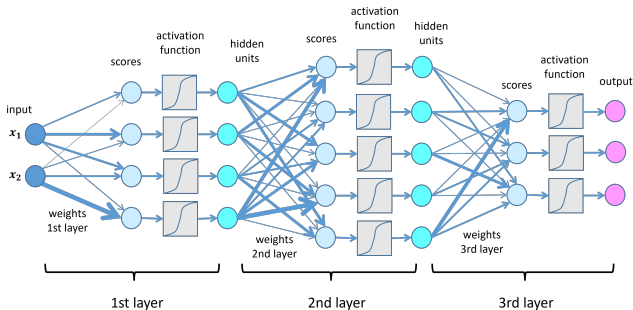
# Various deep learning models VI

## Deep generative models (image generation)



# Neural network

- Dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}_{i=1}^N$
- Goal: train a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  which prediction  $f(\mathbf{x})$  approximates  $y$  as best as possible
- $f$  is designed as a composition of functions, inspiring from human brains:  
$$f(\mathbf{x}) = f_1(\mathbf{x}) \circ f_2(\mathbf{x}) \circ \dots \circ f_L(\mathbf{x})$$
- Each function  $f_\ell$  represents a layer



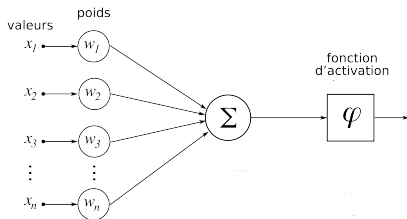
# Artificial neuron [McCulloch et Pitts, 1943]

## Formal neuron

- Input:  $\mathbf{x} \in \mathbb{R}^d$ , Output:  $y$
- Input-output relationship

$$f(\mathbf{x}) = \varphi(\mathbf{w}^\top \mathbf{x} + b)$$

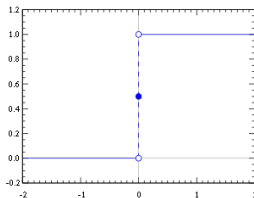
- $\varphi$ : (non-linear) activation function



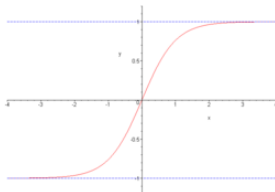
# Activation functions

## Examples

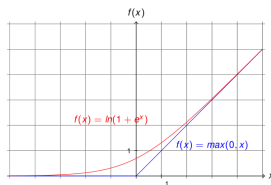
- Identity function:  $\varphi(z) = z$
- Heaviside:  $\varphi(z) = 0$  si  $z < 0$ , 1 sinon
- sigmoid:  $\varphi(z) = \frac{1}{1+e^{-z}}$
- tanh:  $\varphi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \frac{e^{2z} - 1}{e^{2z} + 1}$
- ReLU:  $\varphi(z) = \max(0, z)$



heaviside



tanh



ReLU

# The formal neuron as a learning machine (Perceptron 1958)

- Dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}\}_{i=1}^N$
- Decision function:  $f(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$

---

## Algorithm 1 Perceptron Algorithm

---

Initialize  $\mathbf{w}^0$ ,  $b^0$ , set the learning rate  $\eta$ ,  $t = 0$

**repeat**

    Draw randomly a sample  $(\mathbf{x}_j, y_j)$

**if**  $y_j ((\mathbf{w}^t)^\top \mathbf{x}_j + b^t) \leq 0$  **then**

        Update:  $\begin{pmatrix} \mathbf{w}^{t+1} \\ b^{t+1} \end{pmatrix} = \begin{pmatrix} \mathbf{w}^t \\ b^t \end{pmatrix} + \eta \times y_j \times \begin{pmatrix} \mathbf{x}_j \\ 1 \end{pmatrix}, \quad t = t + 1$

**end if**

**until** convergence

---

- The learning rule is a **stochastic gradient algorithm** for minimizing the number of wrongly predicted labels
- Under some mild conditions, the algorithm is guaranteed to converge after a finite number of iterations (Novikof, 1962)

# Perceptron algorithm and stochastic gradient descent (SGD)

- Objective function:  $J(\mathbf{w}) = \sum_{i=1}^N -\mathbf{1}_{y_i(\mathbf{w}^\top \mathbf{x}_i + b) \leq 0} y_i(\mathbf{w}^\top \mathbf{x}_i + b) = \sum_{i=1}^N \text{cost}(y_i, \mathbf{x}_i)$
- Goal: minimize the overall classification error
- Gradient  $\nabla_{\mathbf{w}} J = - \sum_{i=1}^N \mathbf{1}_{y_i(\mathbf{w}^\top \mathbf{x}_i + b) \leq 0} y_i \mathbf{x}_i$
- Learning scheme
  - Perform a gradient descent on a sample at a time
    - Pick a sample  $(\mathbf{x}_j, y_j)$  at random
    - Update  $\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla_{\mathbf{w}} \text{cost}(y_j, \mathbf{x}_j)$ .  $b$  is updated similarly
  - Perceptron algorithm performs a **stochastic gradient descent (SGD)**

## Adaline (Widrow - Hoff 1959)

Similar algorithm has been derived for least squares problem

$$J(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \varphi(\mathbf{w}^\top \mathbf{x}_i + b))^2$$

# Batch vs Stochastic Gradient Descent

$$\text{Global loss } J(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \text{cost}(y_i, \mathbf{x}_i)$$

Batch gradient

$$\mathbf{g} = \frac{1}{N} \sum_{i=1}^N \nabla_{\mathbf{w}} \text{cost}(y_i, \mathbf{x}_i)$$

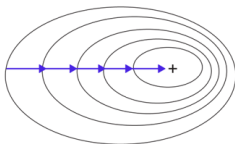
Stochastic gradient

$$\mathbf{g}_i = \frac{1}{N} \nabla_{\mathbf{w}} \text{cost}(y_i, \mathbf{x}_i)$$

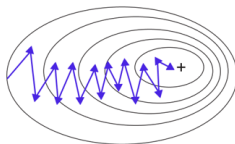
Mini-batch gradient

$$\mathbf{g}_{\mathcal{B}} = \frac{1}{N} \sum_{i \in \mathcal{B}} \nabla_{\mathbf{w}} \text{cost}(y_i, \mathbf{x}_i)$$

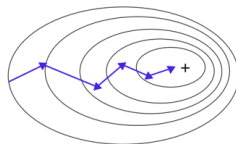
Batch Gradient Descent



Stochastic Gradient Descent

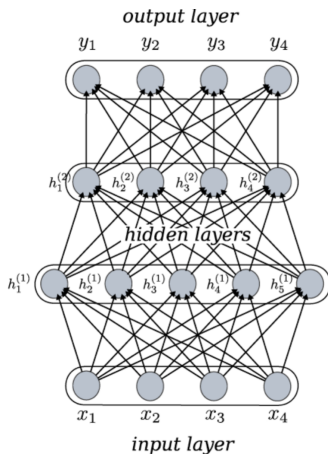


Mini-Batch Gradient Descent



<https://alwaysai.co/blog/what-is-gradient-descent>

# Multi-layer perceptron (MLP)



$$y = \varphi(W_3 h^{(2)})$$

$$h^{(2)} = \varphi(W_2 h^{(1)})$$

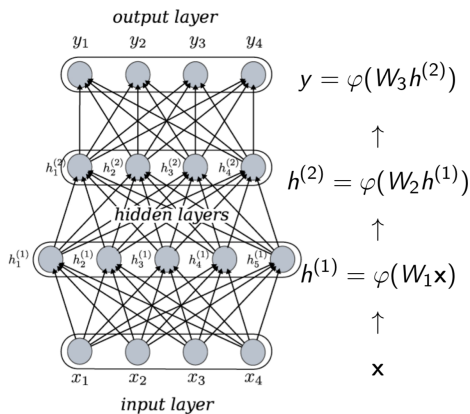
$$h^{(1)} = \varphi(W_1 \mathbf{x})$$

$\mathbf{x}$

## Learning

Use backpropagation algorithm to compute the parameters (weights)  $\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3$

# Forward and Backward Propagation



$$\nabla_{W_3} J = (y - y_a) \varphi'(W_3 h^{(2)}) h^{(2)}$$

↓

$$\nabla_{W_2} J = \nabla_{h^{(2)}} J \varphi'(W_2 h^{(1)}) h^{(1)}$$

↓

$$\nabla_{W_1} J =$$

## Update of the parameters

$$\mathbf{W}_\ell^{t+1} \leftarrow \mathbf{W}_\ell^t - \eta \nabla_{\mathbf{W}_\ell^t} J(\mathbf{W}_\ell^t) \quad \text{pour } \ell = 1, 2, \dots, L$$

# Training algorithm

---

**Algorithm 2** Backpropagation algorithm

---

Initialize the weights  $\mathbf{W}_\ell$ , set  $L$ ,  $\eta$  the learning rate, the mini-batch size  $B$

**while** non convergence **do**

**for**  $t = 1 \rightarrow \text{round}(N/B)$  **do**

        Draw randomly a set of  $B$  samples  $\mathcal{B}_t = \{(\mathbf{x}_i, y_i)\}_{i=1}^B$

        #Forward pass

        Compute the hidden vectors  $\mathbf{h}^{(\ell)} = \varphi(\mathbf{W}^t \mathbf{h}^{(\ell-1)})$ ,  $\forall \ell = 1, \dots, L$  and the loss  $J(\mathbf{W}^t, \mathcal{B}_t)$  based on mini-batch  $\mathcal{B}_t$

        #Backward pass

        Compute the gradients  $\nabla_{\mathbf{W}_\ell^t} J(\mathbf{W}^t, \mathcal{B}_t) \quad \forall \ell$  based on mini-batch  $\mathcal{B}_t$

        Update weights  $\mathbf{W}_\ell^{t+1} \leftarrow \mathbf{W}_\ell^t - \eta \nabla_{\mathbf{W}_\ell} J(\mathbf{W}^t, \mathcal{B}_t) \quad \forall \ell = 1, \dots, L$

**end for**

**end while**

---

## Learning rate

Several methods exist to set up the learning rate  $\eta$ : fixed, adaptive, momentum...

# Optimization in deep learning: optimizer I

## Optimizer Adagrad: SGD algorithms with Adaptive learning rate

AdaGrad adapts learning rates for each parameter  $w_k$  at each step  $t$  based on historical squared gradients. Let  $g_k^t$  the gradient of  $J$  w.r.t  $w_k$ , one computes:

$$G_k^t = G_k^{t-1} + (g_k^t)^2$$

with update rule:

$$w_k^{t+1} = w_k^t - \frac{\eta}{\sqrt{G_k^t + \epsilon}} g_k^t$$

Pros:

- Good for sparse data
- Per-parameter adaptive learning rate

Cons:

- Learning rate decays too aggressively.

# Optimization in deep learning: optimizer II

## Optimizer RMS Prop

RMSProp uses an exponential decaying average of past gradients:

$$E[g_k^2]_t = \rho E[g_k^2]_{t-1} + (1 - \rho)(g_k^t)^2$$

Update rule:

$$w_k^{t+1} = w_k^t - \frac{\eta}{\sqrt{E[g_k^2]_t + \epsilon}} g_k^t$$

Characteristics:

- Controls the unbounded growth of AdaGrad's accumulator.
- Works well for deep learning tasks where gradients vary over time.

# Optimization in deep learning: optimizer III

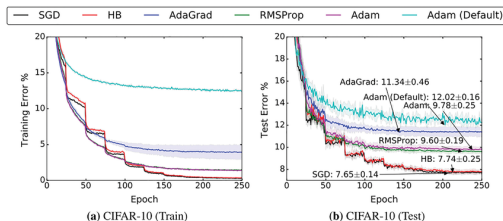
## Optimizer Adam

Adam combines Momentum and RMSProp:

$$m_k^t = \beta_1 m_k^{t-1} + (1 - \beta_1) g_k^t$$

$$v_k^t = \beta_2 v_k^{t-1} + (1 - \beta_2) (g_k^t)^2$$

Bias-corrected estimates:  $\hat{m}_k^t = \frac{m_k^t}{1 - \beta_1}$ ,  $\hat{v}_k^t = \frac{v_k^t}{1 - \beta_2}$ , Update:  $w_k^{t+1} = w_k^t - \eta \frac{\hat{m}_k^t}{\sqrt{\hat{v}_k^t + \epsilon}}$



<https://www.researchgate.net/publication/368951831/figure/fig71/AS:>

11431281233495765@1712057231036/Standard-SGD-and-SGD-with-momentum-vs-AdaGrad-RMSProp-Adam-on-CIFAR-10-dataset.tif

# Optimization in deep learning: regularization

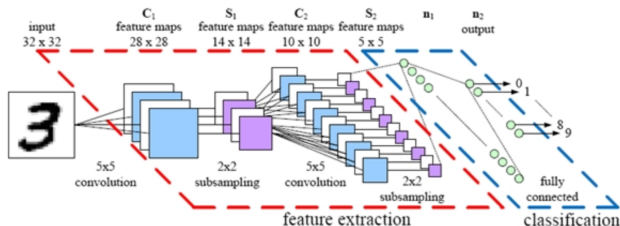
## Regularization schemes

- Explicit regularization: rather minimize  $J(\mathbf{W}) + \lambda \|\mathbf{W}\|^2$  with  $\lambda > 0$
- Dropout: Randomly drop some weights at training time
  - Parameter: dropout percentage  $p$  (percentage of parameters to deactivate)
  - Each weight is dropped with probability  $p$  (Bernoulli) i.e. is inactive in the forward and backward pass
- Batch Normalization

$$\hat{x} = (x - \mu_B) \oslash (\sqrt{\sigma_B^2 + \epsilon})$$

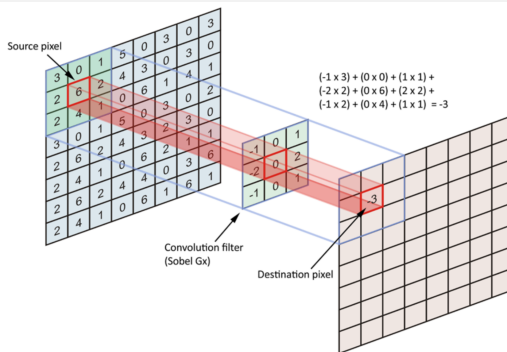
$\oslash$  elementwise division

# What if the inputs are images?



Use Convolutional Neural Network (CNN)

# More on CNN I



Convolution: Given input  $I$  and filter  $K$ ,

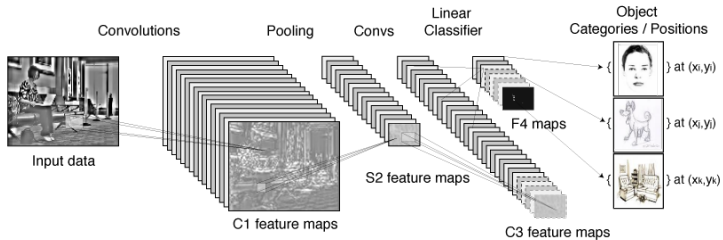
$$(I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n)$$

Followed by pooling (max pooling):  $\text{pool}(i, j) = \max_{(m, n) \in \text{window}} I(i + m, j + n)$

# More on CNN II

## Stacking several layers

- Each layer includes Convolution and Pooling to summarize spatial information
- Last layers are fully connected layers (MLP-like) to yield the output



animation : <http://cs231n.github.io/convolutional-networks/>

# Recurrent networks

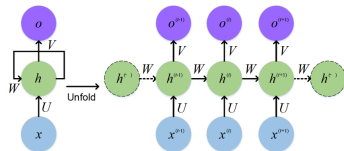
Many tasks involve sequential or temporal structure:

- Natural language: sentences, documents
- Time-series prediction: stock prices, weather
- Audio processing: speech, music
- Video: sequence of frames

RNNs introduce **recurrence**: each output depends on the current input *and* previous hidden state.

$$h_t = f(h_{t-1}, x_t)$$

This allows information to persist across time.



[https://www.researchgate.net/publication/](https://www.researchgate.net/publication/318332317/figure/fig1/AS:61430956243766401523474221928/)

318332317/figure/fig1/AS:

61430956243766401523474221928/

The-standard-RNN-and-unfolded-RNN.png

# Packages

- TensorFlow (Google, python) <https://www.tensorflow.org>
- Keras (Google, python) <https://keras.io/> (+ Theano ou TF)
- Pytorch