

# Foundation of Deep Learning

## PytzMLS2018@IdabaX: CIVE UDOM Tanzania.

Anthony Faustine

PhD Fellow  
(IDLab research group-Ghent University)



4 April 2018



# Learning goal

- Understand the basic building block of deep learning model.
- Learn how to train deep learning models.
- Learn different techniques used in practise to train deep learning models.
- Understand different modern deep learning architectures and their application.
- Explore opportunities and research direction in deep learning.

# Outline

Introduction to Deep learning

Multilayer Perceptron

Training Deep neural networks

Deep learning in Practice

Modern Deep learning Architecture

Limitation and Research direction of deep neural networks

# What is Deep Learning

**Deep Learning** a subclass of machine learning algorithms that learn underlying features in data using multiple processing layers with multiple levels of abstraction.

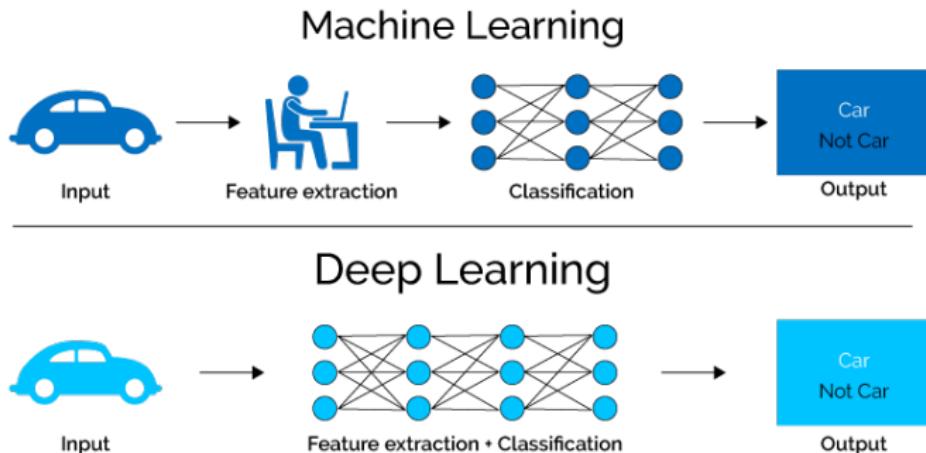
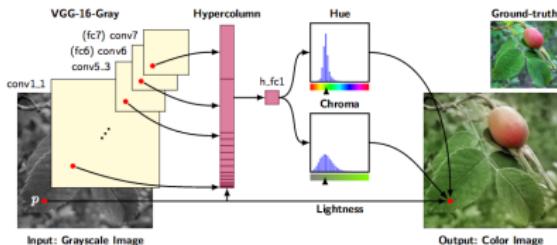


Figure 1: ML vs Deep learning: credit:

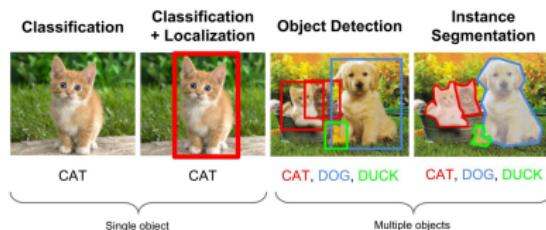
# Deep Learning Success

## Automatic Colorization



**Figure 2:** Automatic colorization

## Object Classification and Detection



**Figure 3:** Object recognition

## Image Captioning

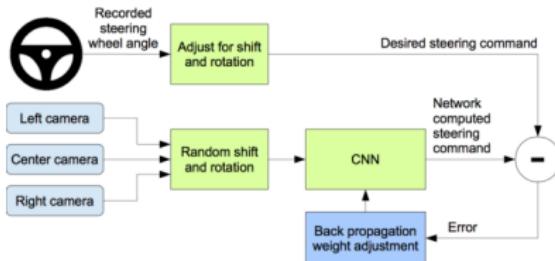


## Image Style Transfer



# Deep Learning Success

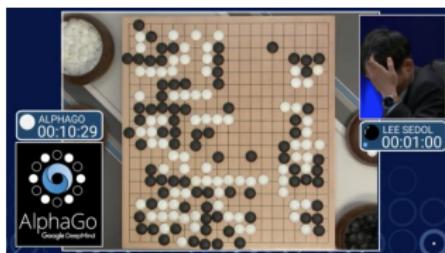
## Self driving car



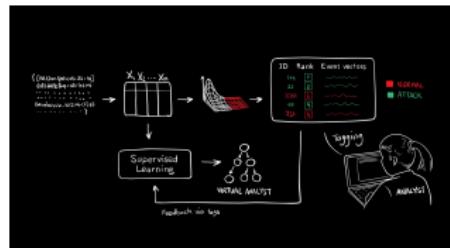
## Drones



## Game

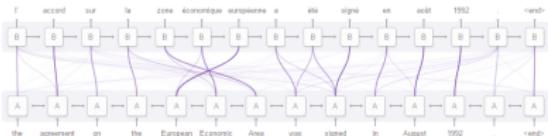


## Cyber attack prediction

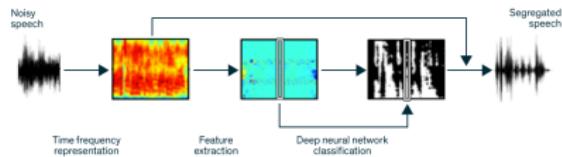


# Deep Learning Success

## Machine translation



## Speech Processing



## Automatic Text Generation

documents reveal iot-specific televisions can be used to secretly record conversations  
criminals who initiated the attack managed to commandeer a large number of internet-connected devices in current use.  
documents revealed that microwave ovens can spy on you - maybe if you personally don't suspect the sub-par security of the iot .

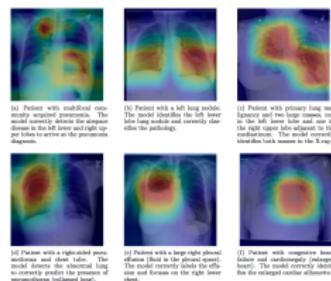
( IoT ) security breaches have been dominating the headlines lately . WikiLeaks 's leak of NSA documents revealed that internet-connected televisions can be used to secretly record conversations . Trump 's administration continues to increase the pressure on you - iot devices that were intentionally designed to be cameras so that they don't suffer the consequences of the sub-par security of the iot . connected gadgets may well be unwillingly cooperating with criminals . Last October , SYN came under an attack that disrupted access to popular websites . The cybercriminals who issued the attack managed to commandeer a large number of IoT devices . SYN is just one of many cybercriminals that have been using the iot to commit crimes . whether because of government regulation or greed it can expect increased investment in IoT security technologies . In its recently-released TechRadar report for security and risk professionals , Forrester Research identified relevant and important IoT security technologies , warning that " there is no single , magic security bullet that can easily fix all IoT security issues . " Basic IoT security measures such as strong authentication and secure communication protocols are important , but as the complexity of IoT security concerns is a bit more challenging than traditional network security because there is a wider range of communication protocols , standards , and device capabilities , all of increased complexity . Key capabilities include traditional endpoint security features such as antivirus and antimalware as well as other features such as fire and intrusion detection systems . Sample vendors : Bayshore Networks , Cisco , Darktrace , and Symantec . IoT authentication : Providing the ability for users to authenticate multiple users of a single device / such as a connected car . - comes from mobile static representations to move ahead in the identification mechanisms .

## Music composition

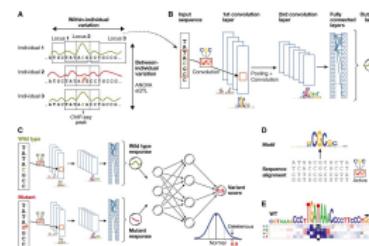


# Deep Learning Success

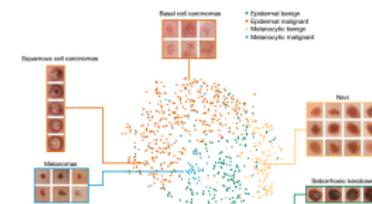
## Pneumonia Detection on Chest X-Rays



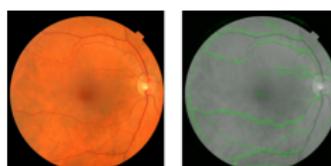
## Computational biology



## Diagnosis of Skin Cancer



Predict heart disease risk from eye scans



More stories



# Why Deep Learning and why now?

Why deep learning: Hand-Engineered Features vs. Learned features.

## Traditional ML

- Use engineered feature to extract useful patterns from data.
- Complex and difficult since different data sets require different feature engineering approach

## Deep learning

- Automatically discover and extract useful pattern from data.
- Allows learning complex features e.g speech and complex networks.

# Why Deep Learning and why now?

## Why Now?

### Big data availability

- Large datasets
- Easier collection and storage

### Increase in computational power

- Modern GPU architecture.

### Improved techniques

- Five decades of research in machine learning.

### Open source tools and models

- Tensorflow.
- Pytorch
- Keras

# Outline

Introduction to Deep learning

Multilayer Perceptron

Training Deep neural networks

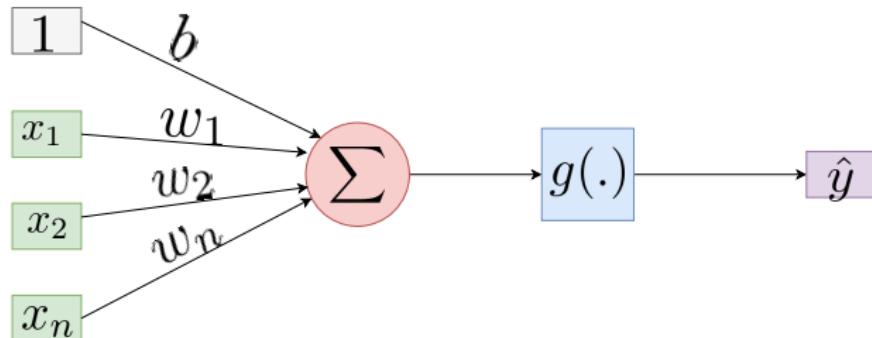
Deep learning in Practice

Modern Deep learning Architecture

Limitation and Research direction of deep neural networks

# The Perceptron

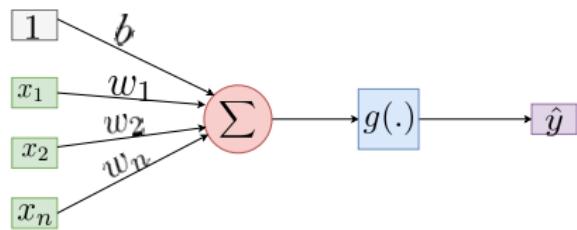
A perceptron is a simple model of a neuron.



The output:  $\hat{y} = f(x) = g(z(x))$  where

- $x, y$  input, output.
- $w, b$  weight and bias parameter  $\theta$
- activation function:  $g(\cdot)$
- pre-activation:  
$$z(x) = \sum_{i=1}^n w_i x_i + b$$

## Perceptron



$$\hat{y} = g(z(x))$$

$$\hat{y} = g(b + \sum_{i=1}^n w_i x_i))$$

$$\hat{y} = g(\mathbf{b} + \mathbf{w}\mathbf{x}))$$

# The Perceptron: Activation Function

## Why Activation Functions?

- Activation functions add non-linearity properties to neuro network function.
- Most real-world problems + data are non-linear.
- Activation function need to be differentiable.

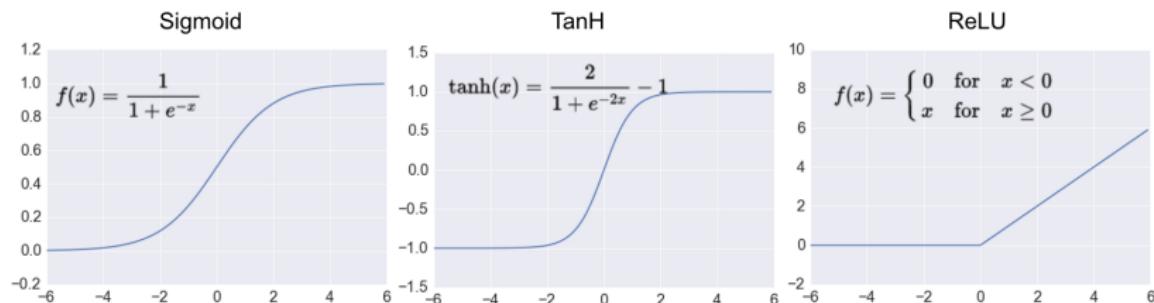
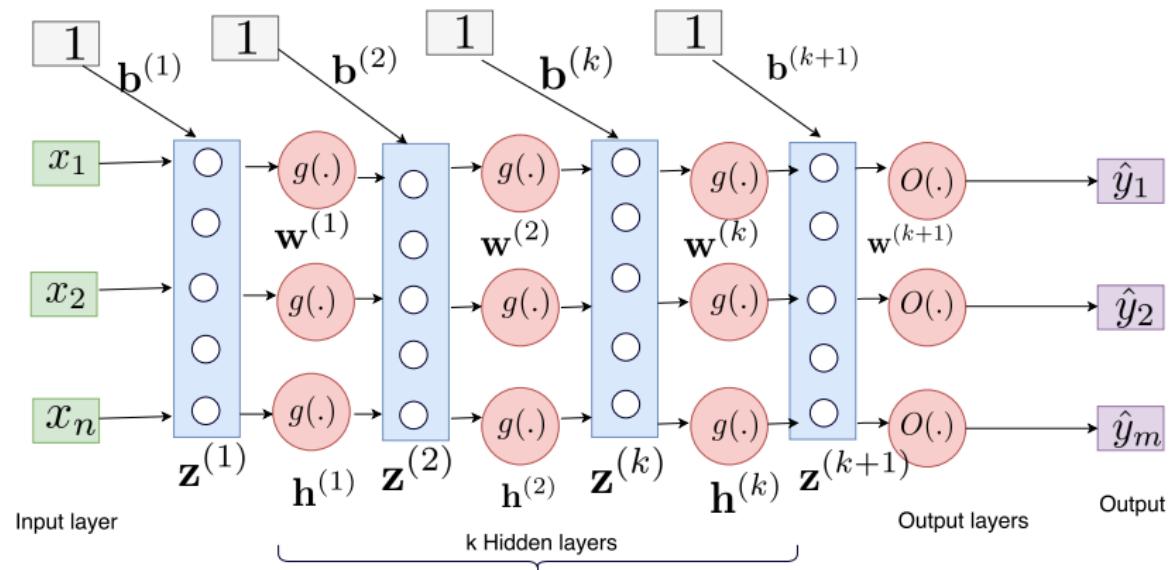


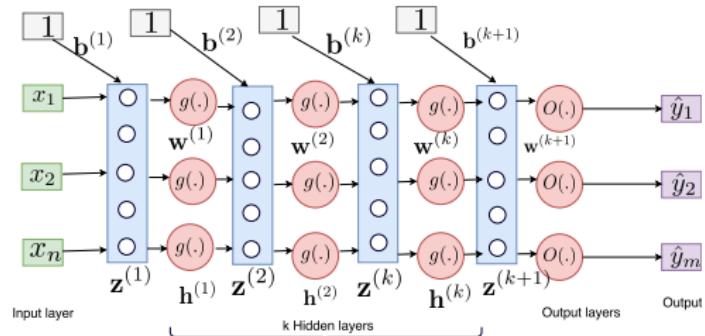
Figure 4: Activation function credit:kdnuggets.com

## Multilayer Perceptrons (MLP)

We can connect lots of perceptron units together into a directed acyclic graph.

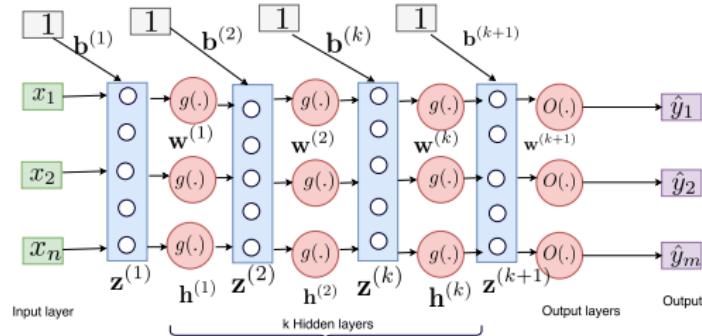


# Multilayer Perceptrons (MLP)



- Consists of  $L$  multiple layers ( $l_1, l_2 \dots l_L$ ) of perceptrons, interconnected in a feed-forward way.
- The first layer  $l_1$  is called the **input layer**  $\Rightarrow$  just pass the information to the next layer.
- The last layer is the **output layer**  $\Rightarrow$  maps to the desired output format.
- The intermediate  $k$  layers are **hidden layers**  $\Rightarrow$  perform computations and transfer the weights from the input layer.

# Multilayer Perceptrons (MLP)



- Input:

$$\mathbf{x} = \{x_1, x_2, \dots, x_d\} \in \mathbb{R}^{(d \times N)}$$

- Pre-activation:

$$\mathbf{z}^{(1)}(\mathbf{x}) = \mathbf{b}^{(1)} + \mathbf{w}^{(1)}(\mathbf{x})$$

where  $z(x)_i = \sum_j w_{i,j}^{(1)} x_j + b_i^{(1)}$

## Hidden layer 1

- Activation

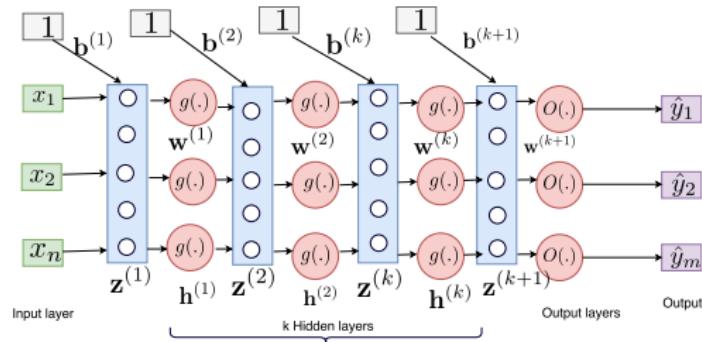
$$\mathbf{h}^{(1)}(\mathbf{x}) = g(\mathbf{z}^{(1)}(\mathbf{x}))$$

$$= g(\mathbf{b}^{(1)} + \mathbf{w}^{(1)}(\mathbf{x}))$$

- Pre-activation

$$\mathbf{z}^{(2)}(\mathbf{x}) = \mathbf{b}^{(2)} + \mathbf{w}^{(2)} \mathbf{h}^{(1)}(\mathbf{x})$$

# Multilayer Perceptrons (MLP)



## Hidden layer 2

- Activation

$$\begin{aligned}\mathbf{h}^{(2)}(\mathbf{x}) &= g(\mathbf{z}^{(2)}(\mathbf{x})) \\ &= g(\mathbf{b}^{(2)} + \mathbf{w}^{(2)}\mathbf{h}^{(1)}(\mathbf{x}))\end{aligned}$$

- Pre-activation

$$\mathbf{z}^{(3)}(\mathbf{x}) = \mathbf{b}^{(3)} + \mathbf{w}^{(3)}\mathbf{h}^{(2)}(\mathbf{x})$$

## Hidden layer $k$

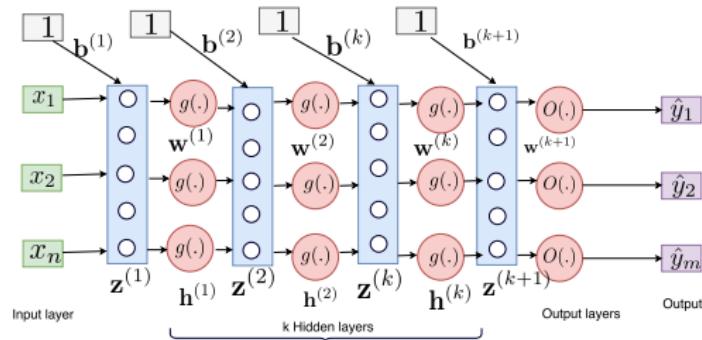
- Activation

$$\begin{aligned}\mathbf{h}^{(k)}(\mathbf{x}) &= g(\mathbf{z}^{(k)}(\mathbf{x})) \\ &= g(\mathbf{b}^{(k)} + \mathbf{w}^{(k)}\mathbf{h}^{(k-1)}(\mathbf{x}))\end{aligned}$$

- Pre-activation

$$\mathbf{z}^{(k+1)}(\mathbf{x}) = \mathbf{b}^{(k+1)} + \mathbf{w}^{(k+1)}\mathbf{h}^{(k)}(\mathbf{x})$$

# Multilayer Perceptrons (MLP)



## Output layer

- Activation

$$\begin{aligned} \mathbf{h}^{(k+1)}(\mathbf{x}) &= O(\mathbf{z}^{(k+1)}(\mathbf{x})) \\ &= O(\mathbf{b}^{(k+1)} + \mathbf{w}^{(k+1)} \mathbf{h}^{(k)}(\mathbf{x})) \\ &= \hat{\mathbf{y}} \end{aligned}$$

where  $O(\cdot)$  is output activation function

## Output activation function

- Binary classification:  $y \in \{0, 1\} \Rightarrow \text{sigmoid}$
- Multiclass classification:  $y \in \{0, K - 1\} \Rightarrow \text{softmax}$
- Regression:  $y \in \mathbb{R}^n \Rightarrow \text{identity}$  sometime RELU.

## Demo Playground



# MLP: Pytorch

```
import torch
model = torch.nn.Sequential(
    torch.nn.Linear(2, 16),
    torch.nn.ReLU(),
    torch.nn.Linear(16, 64),
    torch.nn.ReLU(),
    torch.nn.Linear(64, 1024),
    torch.nn.ReLU(),
    torch.nn.Linear(1024, 1),
    torch.nn.Sigmoid()
)
```

# MLP: Pytorch

```
import torch
from torch.nn import functional as F

class MLP(torch.nn.Module):
    def __init__(self):
        super(MLP, self).__init__()
        self.fc1 = torch.nn.Linear(2, 16)
        self.fc2 = torch.nn.Linear(16, 64)
        self.fc3 = torch.nn.Linear(64, 1024)
        self.out = torch.nn.Linear(1024, 1)

    def forward(self, x):
        x = F.relu(self.fc1(x))
        x = F.relu(self.fc2(x))
        x = F.relu(self.fc3(x))
        out = F.sigmoid(self.out(x))

    return x

model = MLP()
```



# Outline

Introduction to Deep learning

Multilayer Perceptron

Training Deep neural networks

Deep learning in Practice

Modern Deep learning Architecture

Limitation and Research direction of deep neural networks

# Training Deep neural networks

To train DNN we need:

- ① Define loss function:

$$\mathcal{L}(f(\mathbf{x}^{(i)} : \theta), \mathbf{y}^{(i)})$$

- ② A procedure to compute gradient  $\frac{\partial J_\theta}{\partial \theta}$
- ③ Solve optimisation problem.

# Training Deep neural networks: Define loss function

The type of **Loss** function is determined by the output layer of MLP.

## Binary classification

### Output

- Predict  $y \in \{0, 1\}$
- Use sigmoid  $\sigma(\cdot)$  activation function.

$$p(y = 1|x) = \frac{1}{1 + e^{-x}}$$

### Loss

- Binary cross entropy.

$$\mathcal{L}(\hat{y}, y) = y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

- `torch.nn.BCELoss()`

# Training Deep neural networks: Define loss function

## Mutli class classification

### Output

- Predict  $y \in \{1, k\}$
- Use softmax  $\sigma(\cdot)$  activation function.

$$p(y = i|x) = \frac{\exp(x_i)}{\sum_j^k}$$

### Loss

- Cross entropy.

$$\mathcal{L}(\hat{y}, y) = \sum_{i=1}^k y_i \log \hat{y}_i$$

- `torch.nn.CrossEntropyLoss()`

# Training Deep neural networks: Define loss function

## Regression

### Output

- Predict  $y \in \mathbb{R}^n$
- Use identity activation function and sometime ReLU activation.

### Loss

- Squared error loss.

$$\mathcal{L}(\hat{y}, y) = \frac{1}{2}(y_i - \hat{y}_i)^2$$

- `torch.nn.MSELoss()`

# Training Deep neural networks: Compute Gradients

**Backpropagation:** a procedure that is used to compute gradients of a loss function.

- It is based on the application of the chain rule and computationally proceeds 'backwards'.

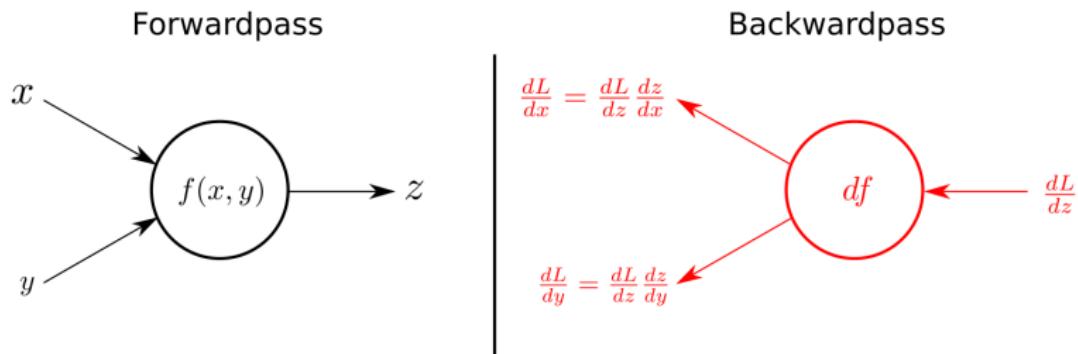
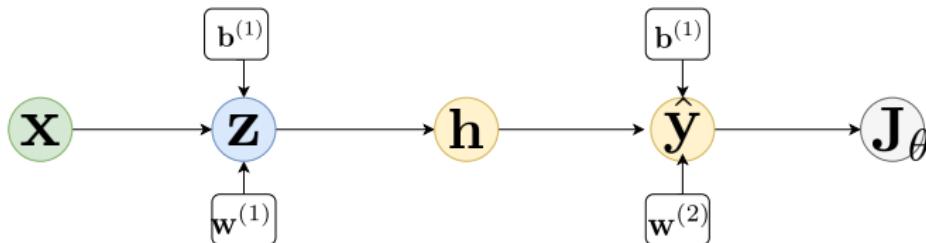


Figure 5: Back propagation: credit: Flair of Machine Learnin

# Training Deep neural networks: Backpropagation

Consider a following single hidden layer MLP.



Forward path

$$\mathbf{z} = \mathbf{w}^1 \mathbf{x} + \mathbf{b}^1$$

$$\mathbf{h} = g(\mathbf{z})$$

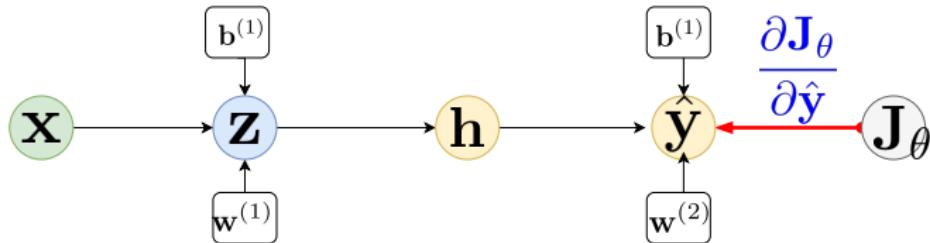
$$\hat{\mathbf{y}} = \mathbf{w}^2 \mathbf{h} + \mathbf{b}^2$$

$$J_\theta = \frac{1}{2} \|\mathbf{y} - \hat{\mathbf{y}}\|^2$$

We need to find:  $\frac{\partial J_\theta}{\partial \mathbf{w}^{(1)}}$ ,  $\frac{\partial J_\theta}{\partial \mathbf{b}^{(1)}}$ ,  $\frac{\partial J_\theta}{\partial \mathbf{w}^{(2)}}$   
and  $\frac{\partial J_\theta}{\partial \mathbf{b}^{(2)}}$

# Training Deep neural networks: Backpropagation

Back ward path

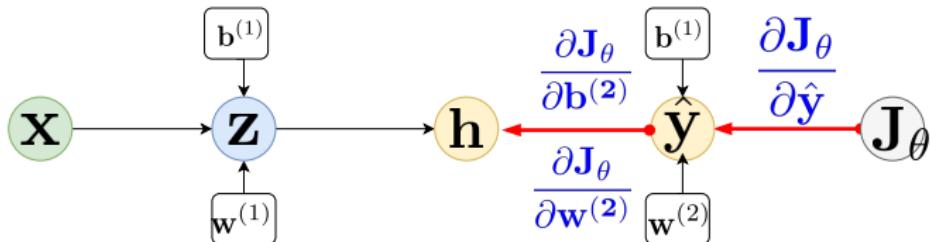


$$\mathbf{J}_\theta = \frac{1}{2} \|\mathbf{y} - \hat{\mathbf{y}}\|^2$$

$$\frac{\partial \mathbf{J}_\theta}{\partial \hat{\mathbf{y}}} = \|\mathbf{y} - \hat{\mathbf{y}}\|$$

# Training Deep neural networks: Backpropagation

## Backward path



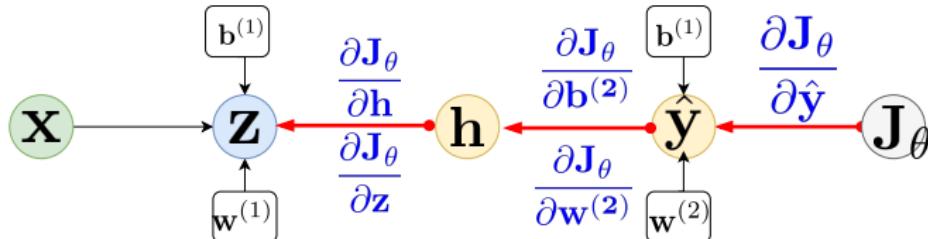
$$\hat{y} = \mathbf{w}^2 \mathbf{h} + \mathbf{b}^2$$

$$\frac{\partial \mathbf{J}_\theta}{\partial \mathbf{w}^{(2)}} = \frac{\partial \hat{y}}{\partial \mathbf{w}^{(2)}} \cdot \frac{\partial \mathbf{J}_\theta}{\partial \hat{y}} = \mathbf{h}^T \cdot \|\mathbf{y} - \hat{y}\|$$

$$\frac{\partial \mathbf{J}_\theta}{\partial \mathbf{b}^{(2)}} = \frac{\partial \hat{y}}{\partial \mathbf{b}^{(2)}} \cdot \frac{\partial \mathbf{J}_\theta}{\partial \hat{y}} = \|\mathbf{y} - \hat{y}\|$$

# Training Deep neural networks: Backpropagation

Backward path



$$\hat{y} = w^2 h + b^2$$

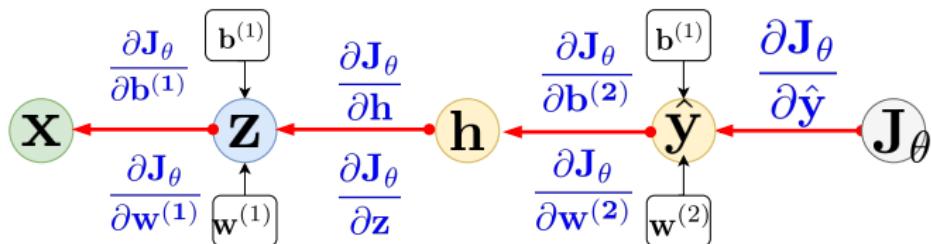
$$h = g(z)$$

$$\frac{\partial J_\theta}{\partial h} = \frac{\partial \hat{y}}{\partial h} \cdot \frac{\partial J_\theta}{\partial \hat{y}} = w^{(2)T} \cdot ||y - \hat{y}||$$

$$\frac{\partial J_\theta}{\partial z} = \frac{\partial h}{\partial z} \cdot \frac{\partial J_\theta}{\partial h} = g'((z)) \cdot \frac{\partial J_\theta}{\partial h}$$

# Training Deep neural networks: Backpropagation

Backward path



$$\mathbf{z} = \mathbf{w}^1 \mathbf{h} + \mathbf{b}^1$$

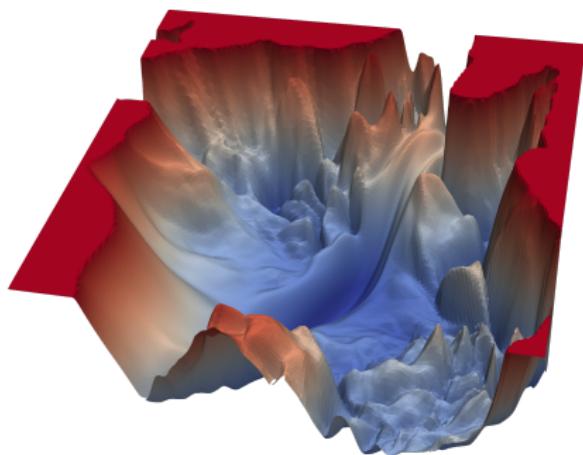
$$\frac{\partial J_\theta}{\partial w^{(1)}} = \frac{\partial \mathbf{z}}{\partial w^{(1)}} \cdot \frac{\partial J_\theta}{\partial \mathbf{z}} = \mathbf{x}^T \cdot \frac{\partial J_\theta}{\partial \mathbf{z}}$$

$$\frac{\partial J_\theta}{\partial b^{(1)}} = \frac{\partial \mathbf{z}}{\partial b^{(1)}} \cdot \frac{\partial J_\theta}{\partial \mathbf{z}} = \frac{\partial J_\theta}{\partial \mathbf{z}}$$

## Training Neural Networks: Solving optimisation problem

**Objective:** Find parameters  $\theta$  :  $\mathbf{w}$  and  $\mathbf{b}$  that minimize the cost function:

$$\arg \max_{\theta} \frac{1}{N} \sum_i \mathcal{L}(f(\mathbf{x}^{(i)} : \theta), \mathbf{y}^{(i)})$$



**Figure 6:** Visualizing the loss landscape of neural nets: *credit: Hao Li*

# Training Neural Networks: Gradient Descent

## Gradient Descent

① Initialize parameter  $\theta$ ,

② Loop until converge

    ① Compute gradient:

$$\frac{\partial J_{\theta}}{\partial \theta}$$

    ② Update parameters:

$$\theta^{t+1} = \theta^t - \alpha \frac{\partial J_{\theta}}{\partial \theta}$$

    ③ Retrn parameter  $\theta$

**Limitation:** Take time to compute

# Training Neural Networks: Stochastic Gradient Descent (SGD)

SGD consists of updating the model parameters  $\theta$  after every sample.

## SGD

Initialize  $\theta$  randomly.

For each training example:

- Compute gradients:  $\frac{\partial J_{i\theta}}{\partial \theta}$
- Update parameters  $\theta$  with update rule:

$$\theta^{(t+1)} := \theta^{(t)} - \alpha \frac{\partial J_{i\theta}}{\partial \theta}$$

Stop when reaching criterion

Easy to compute  $\frac{\partial J_{i\theta}}{\partial \theta}$  but **very noise**.



# Training Neural Networks: Mini-batch SGD training

Make update based on a min-batch  $B$  of example instead of single example  $i$

## Mini-batch SGD

- ① Initialize  $\theta$  randomly.
- ② For each mini-batch  $B$ :
  - Compute gradients:  $\frac{\partial J_\theta}{\partial \theta} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_{k(\theta)}}{\partial \theta}$
  - Update parameters  $\theta$  with update rule:  
$$\theta^{(t+1)} := \theta^{(t)} - \alpha \frac{\partial J_{i\theta}}{\partial \theta}$$
- ③ Stop when reaching criterion

Fast to compute  $\frac{\partial J_\theta}{\partial \theta} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_{k(\theta)}}{\partial \theta}$  and much **better** estimate of the true gradient.

Standard procedure for training deep learning.



# Training Neural Networks: Gradient Descent Issues

## Setting the learning rate $\alpha$

- Small learning rate: Converges slowly and gets stuck in false local minima.
- Large learning rate: Overshoot became unstable and diverge.
- Stable learning rate: Converges smoothly and avoid local minima.

How to deal with this ?

- ① Try lots of different learning rates and see what works for you.
  - Jeremy propose a technique to find stable learning rate
- ② Use an adaptive learning rate that adapts to the landscape of your loss function.

# Training Neural Networks: Adaptive Learning rates algorithm

- ① Momentum
- ② Adagrad
- ③ Adam
- ④ RMSProp

pytorch optimizer algorithms

# Outline

Introduction to Deep learning

Multilayer Perceptron

Training Deep neural networks

Deep learning in Practice

Modern Deep learning Architecture

Limitation and Research direction of deep neural networks

# Deep learning in Practice: Regularization

**Regularization:** Technique to help deep learning network perform better on unseen data.

- Constraints optimization problem to discourage complex model.

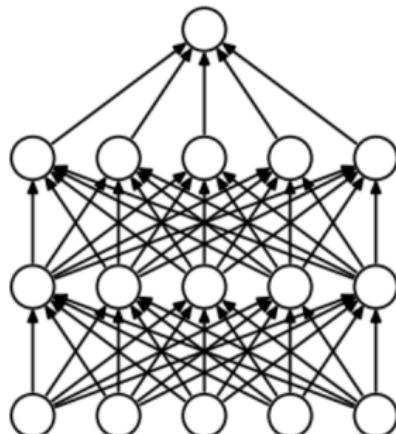
$$\arg \max_{\theta} \frac{1}{N} \sum_i \mathcal{L}(f(\mathbf{x}^{(i)} : \theta), \mathbf{y}^{(i)}) + \lambda \Omega(\theta)$$

- Improve generalization of deep learning model.

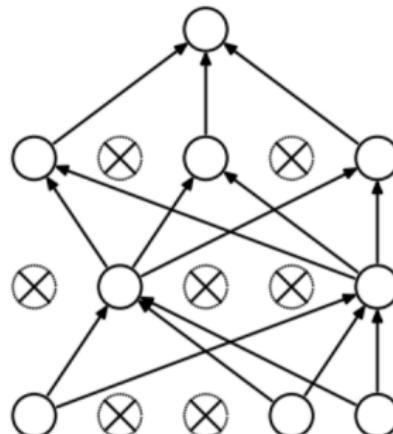
# Regularization 1: Dropout

**Dropout:** Randomly remove hidden unit from a layer during training step and put them back during test.

- Each hidden unit is set to 0 with probability  $p$ .
- Force network to not rely on any hidden node  $\Rightarrow$  prevent neural net from overfitting (improve performance).
- Any dropout probability can be used but 0.5 usually works well.



(a) Standard Neural Net



(b) After applying dropout.

# Regularization 1: Dropout

**Dropout**: in pytorch is implemented as `torch.nn.Dropout`

If we have a network:

```
model = torch.nn.Sequential(  
    torch.nn.Linear(1,100), torch.nn.ReLU(),  
    torch.nn.Linear(100,50), torch.nn.ReLU(),  
    torch.nn.Linear(50,2))
```

We can simply add dropout layers:

```
model = torch.nn.Sequential(  
    torch.nn.Linear(1,100), torch.nn.ReLU(),  
    torch.nn.Dropout()  
    torch.nn.Linear(100,50), torch.nn.ReLU(),  
    torch.nn.Dropout()  
    torch.nn.Linear(50,2))
```

**Note:** A model using dropout has to be set in `train` or `eval` mode.

# Regularization 1: Dropout

**Dropout**: in pytorch is implemented as `torch.nn.Dropout`

If we have a network:

```
model = torch.nn.Sequential(  
    torch.nn.Linear(1,100), torch.nn.ReLU(),  
    torch.nn.Linear(100,50), torch.nn.ReLU(),  
    torch.nn.Linear(50,2))
```

We can simply add dropout layers:

```
model = torch.nn.Sequential(  
    torch.nn.Linear(1,100), torch.nn.ReLU(),  
    torch.nn.Dropout()  
    torch.nn.Linear(100,50), torch.nn.ReLU(),  
    torch.nn.Dropout()  
    torch.nn.Linear(50,2))
```

**Note**: A model using dropout has to be set in `train` or `eval` mode.

## Regularization 2: Early Stopping

**Early Stopping:** Stop training before the model overfit.

- Monitor the deep learning training process from overfitting.
  - Stop training when validation error increases.

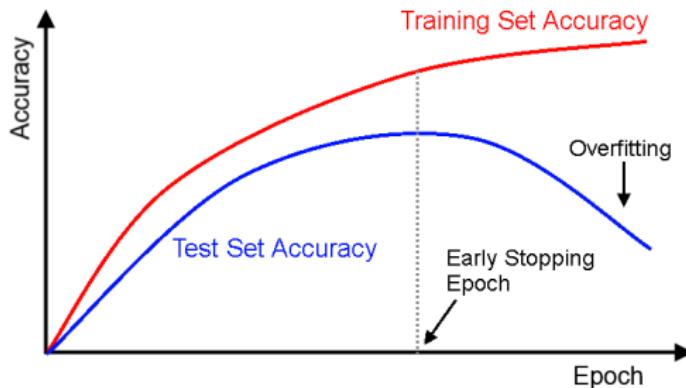


Figure 7: Early stopping: credit: Deeplearning4j.com

# Deep learning in Practice: Batch Normalization

**Batch normalisation:** A technique for improving the performance and stability of deep neural networks.

Training deep neural network is complicated

- The input of each layer changes as the parameter of the previous layer change.
- This slow down the training  $\Rightarrow$  require low learning rate and careful parameter initialization.
- Make hard to train models with saturation non-linearity.
- This phenomena is called **Covariate shift**

To address **covariate shift**  $\Rightarrow$  normalise the inputs of each layer for each mini-batch (Batch normalization)

- To have a mean output activation of zero and standard deviation of one.

# Deep learning in Practice: Batch Normalization

**Batch normalisation:** A technique for improving the performance and stability of deep neural networks.

Training deep neural network is complicated

- The input of each layer changes as the parameter of the previous layer change.
- This slow down the training  $\Rightarrow$  require low learning rate and careful parameter initialization.
- Make hard to train models with saturation non-linearity.
- This phenomena is called **Covariate shift**

To address **covariate shift**  $\Rightarrow$  normalise the inputs of each layer for each mini-batch (Batch normalization)

- To have a mean output activation of zero and standard deviation of one.

# Deep learning in Practice: Batch Normalization

If  $x_1, x_2, \dots, x_B$  are the sample in the batch with mean  $\hat{\mu}_b$  and variance  $\hat{\sigma}_b^2$ .

- During training batch normalization shift and rescale each component of the input according to batch statistics to produce output  $y_b$ :

$$y_b = \gamma \odot \frac{x_b - \hat{\mu}_b}{\sqrt{\hat{\sigma}_b^2 + \epsilon}} + \beta$$

where

- $\odot$  is the Hadamard component-wise product.
- The parameter  $\gamma$  and  $\beta$  are the desired moments which are either fixed or optimized during training.
- As for dropout the model behave differently during training and test.

# Deep learning in Practice: Batch Normalization

**Batch Normalization:** in pytorch is implemented as  
`torch.nn.BatchNorm1d`

If we have a network:

```
model = torch.nn.Sequential(  
    torch.nn.Linear(1,100), torch.nn.ReLU(),  
    torch.nn.Linear(100,50), torch.nn.ReLU(),  
    torch.nn.Linear(50,2))
```

We can simply add batch normalization layers:

```
model = torch.nn.Sequential(  
    torch.nn.Linear(1,100), torch.nn.ReLU(),  
    torch.nn.BatchNorm1d(100)  
    torch.nn.Linear(100,50), torch.nn.ReLU(),  
    torch.nn.BatchNorm1d(50)  
    torch.nn.Linear(50,2))
```

**Note:** A model using batch has to be set in `train` or `eval` mode.

# Deep learning in Practice: Batch Normalization

Batch Normalization: in pytorch is implemented as  
`torch.nn.BatchNorm1d`

If we have a network:

```
model = torch.nn.Sequential(  
    torch.nn.Linear(1,100), torch.nn.ReLU(),  
    torch.nn.Linear(100,50), torch.nn.ReLU(),  
    torch.nn.Linear(50,2))
```

We can simply add batch normalization layers:

```
model = torch.nn.Sequential(  
    torch.nn.Linear(1,100), torch.nn.ReLU(),  
    torch.nn.BatchNorm1d(100)  
    torch.nn.Linear(100,50), torch.nn.ReLU(),  
    torch.nn.BatchNorm1d(50)  
    torch.nn.Linear(50,2))
```

**Note:** A model using batch has to be set in `train` or `eval` mode.

# Deep learning in Practice: Batch Normalization

## When Applying Batch Normalization

- Carefully shuffle your sample.
- Learning rate can be greater.
- Dropout is not necessary.
- $L^2$  regularization influence should be reduced.

# Deep learning in Practice: Weight Initialization

Before training the neural network you have to initialize its parameters.

## Set all the initial weights to zero

- Every neuron in the network will compute the same output  $\Rightarrow$  same gradients.
- Not recommended

# Deep learning in Practice: Weight Initialization

## Random Initialization

- Initialize your network to behave like zero-mean standard gaussian function.

$$w_i \sim \mathcal{N} \left( \mu = 0, \sigma = \sqrt{\frac{1}{n}} \right)$$
$$b_i = 0$$

where  $n$  is the number of inputs.

# Deep learning in Practice: Weight Initialization

## Random Initialization: Xavier initialization

- Initialize your network to behave like zero-mean standard gaussian function such that

$$w_i \sim \mathcal{N} \left( \mu = 0, \sigma = \sqrt{\frac{1}{n_{in} + n_{out}}} \right)$$
$$b_i = 0$$

where  $n_{in}, n_{out}$  are the number of units in the previous layer and the next layer respectively. where  $n$  is the number of inputs.

# Deep learning in Practice: Weight Initialization

## Random Initialization: Kaiming

- Random initialization that take into account ReLU activation function.

$$w_i \sim \mathcal{N} \left( \mu = 0, \sigma = \sqrt{\frac{2}{n}} \right)$$

$$b_i = 0$$

- Recommended in practise.

# Deep learning in Practice: Pytorch Parameter Initialization

Consider the previous model:

```
model = torch.nn.Sequential(  
    torch.nn.Linear(1,100),  
    torch.nn.ReLU(),  
    torch.nn.BatchNorm1d(100)  
    torch.nn.Linear(100,50),  
    torch.nn.ReLU(),  
    torch.nn.BatchNorm1d(50)  
    torch.nn.Linear(50,2))
```

To apply weight initialization to nn.linear module.

```
def weights_init(m):  
    if isinstance(m, nn.Linear):  
        size = m.weight.size()  
        n_out = size[0]  
        n_in = size[1]  
        variance = np.sqrt(2.0/(n_in + n_out))  
        m.weight.data.normal_(0.0, variance)  
  
model.apply(weights_init)
```

# Outline

Introduction to Deep learning

Multilayer Perceptron

Training Deep neural networks

Deep learning in Practice

Modern Deep learning Architecture

Limitation and Research direction of deep neural networks

# Deep learning Architecture: Convolutional Neural Network

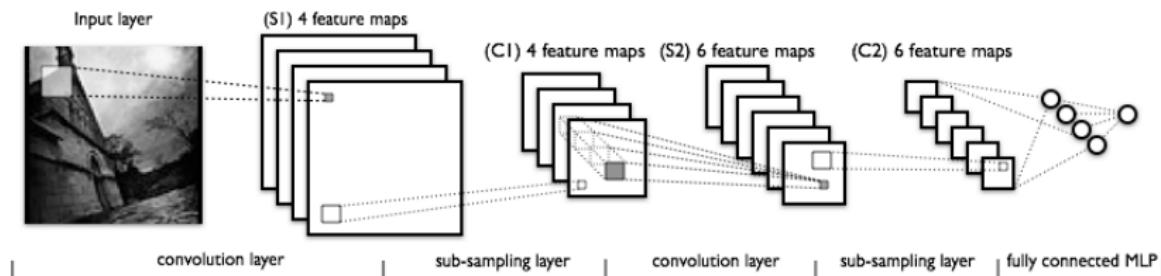
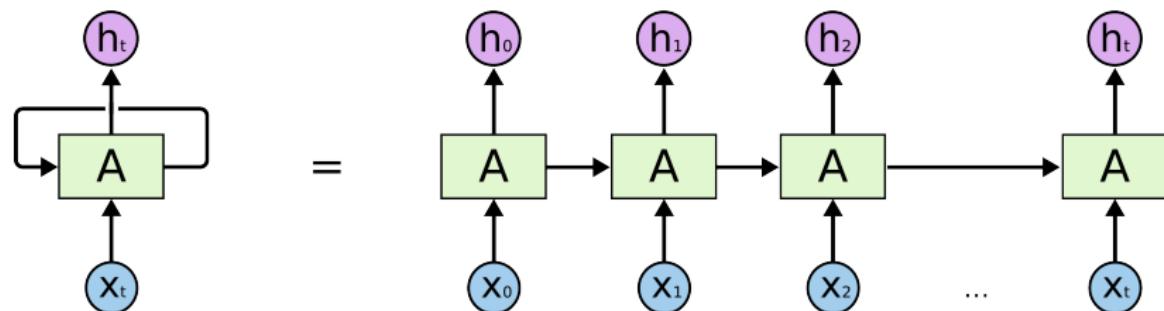


Figure 8: CNN [credit:deeplearning.net]

- Enhances the capabilities of MLP by inserting convolution layers.
- Composed of many “filters”, which convolve, or slide across the data, and produce an activation at every slide position
- Suitable for spatial data, object recognition and image analysis.

# Deep learning Architecture: Recurrent Neural Networks (RNN)

RNN are neural networks with loops in them, allowing information to persist.



- Can model a long time dimension and arbitrary sequence of events and inputs.
- Suitable for sequenced data analysis: time-series, sentiment analysis, NLP, language translation, speech recognition etc.
- Common type: LSTM and GRUs.

# Deep learning Architecture: Autoencoder

**Autoencoder:** A neural network where the input is the same as the output.

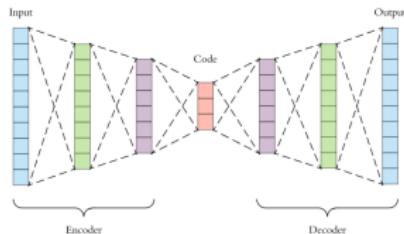


Figure 9: credit: Arden Dertat

- They compress the input into a lower-dimensional code and then reconstruct the output from this representation.
- It is an unsupervised ML algorithm similar to PCA.
- Several types exist: Denoising autoencoder, Sparse autoencoder.

# Deep learning Architecture: Auto-encoder

Autoencoder consists of components: **encoder**, **code** and **decoder**.

- The encoder compresses the input and produces the code,
- The decoder then reconstructs the input only using this code.

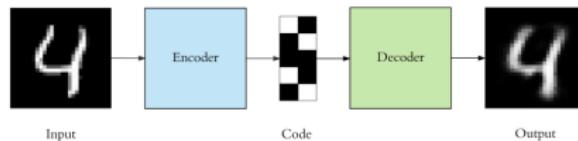


Figure 10: credit:Arden Dertat

# Deep learning Architecture: Deep Generative models

**Idea:** learn to understand data through generation → replicate the data distribution that you give it.

- Can be used to generate Musics, Speech, Language, Image, Handwriting, Language
- Suitable for unsupervised learning as they need lesser labelled data to train.

Two types:

- ① **Autoregressive models:** Deep NADE, PixelRNN, PixelCNN, WaveNet, ByteNet
- ② **Latent variable models:** VAE, GAN.

# Outline

Introduction to Deep learning

Multilayer Perceptron

Training Deep neural networks

Deep learning in Practice

Modern Deep learning Architecture

Limitation and Research direction of deep neural networks

# Limitation

- Very data hungry (eg. often millions of examples)
- Computationally intensive to train and deploy (tractably requires GPUs)
- Poor at representing uncertainty (how do you know what the model knows?)
- Uninterpretable black boxes, difficult to trust
- Difficult to optimize: non-convex, choice of architecture, learning parameters
- Often require expert knowledge to design, fine tune architectures

# Research Direction

- Transfer learning.
- Unsepevised machine learning.
- Computational efficiency.
- Add more reasoning (uncertatinity) abilities  $\Rightarrow$  Bayesian Deep learning
- Many applications which are under-explored especially in developing countries.

# Python Deep learning libraries

Tensorflow



Pytorch



Theano



Keras



Edward

→

Pyro



MXNET



# Lab 3: Introduction to Deep learning

**Part 1:** Feed-forward Neural Network (MLP):

**Objective:** Build MLP classifier to recognize handwritten digits using the MNIST dataset.

**Part 2:** Weight Initialization:

**Objective:** Experiments with different initialization techniques (zero, xavier, kaiming)

**Part 3:** Regularization:

**Objective:** Experiments with different regularization techniques (early stopping, dropout)

# References I

- Deep learning for Artificial Intelligence master course: TelecomBCN Barcelona(winter 2017)
- 6.S191 Introduction to Deep Learning: MIT 2018.
- Deep learning Specilization by Andrew Ng: Coursera
- Deep Learning by Russ Salakhutdinov: MLSS 2017
- Introducucion to Deep learning: CMU 2018
- Cs231n: Convolution Neural Network for Visual Recognition: Stanford 2018
- Deep learning in Pytorch, Francois Fleurent: EPFL 2018
- Advanced Machine Learning Specialization: Coursera