

Deep Learning

Outline:

- Neural networks
- Convolutional Neural Networks
- Deep Generative Learning



Co-funded by
the European Union

What is Deep Learning?

The hierarchical relationship shows how deep learning builds upon previous AI advances while adding its own unique capabilities

ARTIFICIAL INTELLIGENCE

Any technique that enables computers to mimic human behavior



MACHINE LEARNING

Ability to learn without explicitly being programmed



DEEP LEARNING

Extract patterns from data using neural networks



Teaching computers how to **learn a task** directly from **raw data**

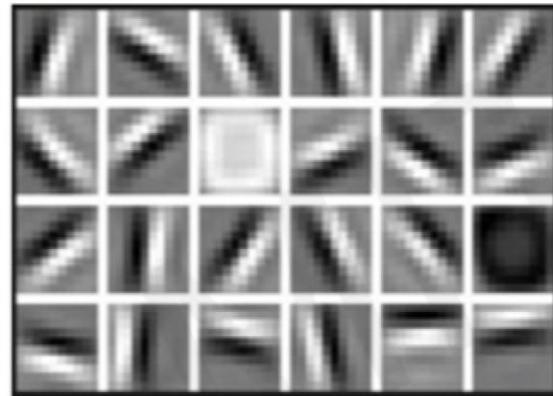
Why Deep Learning?

- Deep learning automatically learns relevant features at multiple levels of abstraction
- Lower-level features combine to form increasingly complex and meaningful representations

Hand engineered features are time consuming, brittle, and not scalable in practice

Can we learn the **underlying features** directly from data?

Low Level Features



Lines & Edges

Mid Level Features



Eyes & Nose & Ears

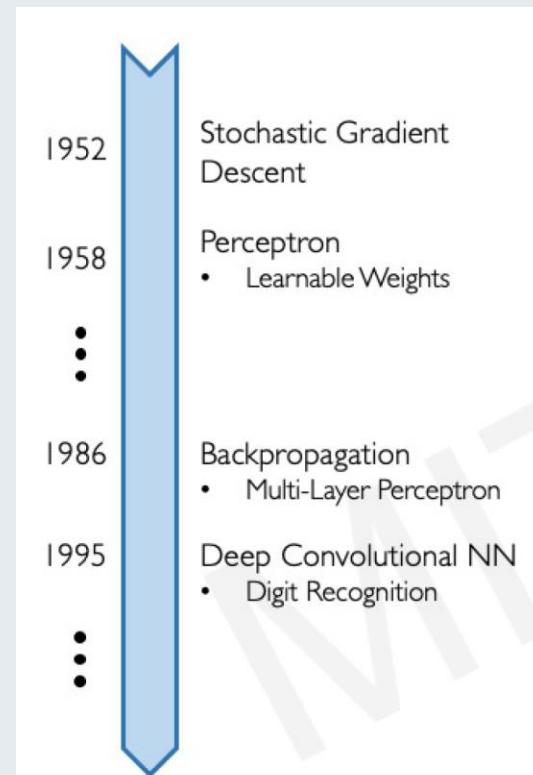
High Level Features



Facial Structure

Why Now?

- The field of neural networks has existed since the 1950s but only recently became practically viable
- Three key developments enabled modern deep learning: availability of big data, powerful GPU hardware, and improved software tools



Neural Networks date back decades, so why the resurgence?

I. Big Data

- Larger Datasets
- Easier Collection & Storage



2. Hardware

- Graphics Processing Units (GPUs)
- Massively Parallelizable



3. Software

- Improved Techniques
- New Models
- Toolboxes



The Perceptron

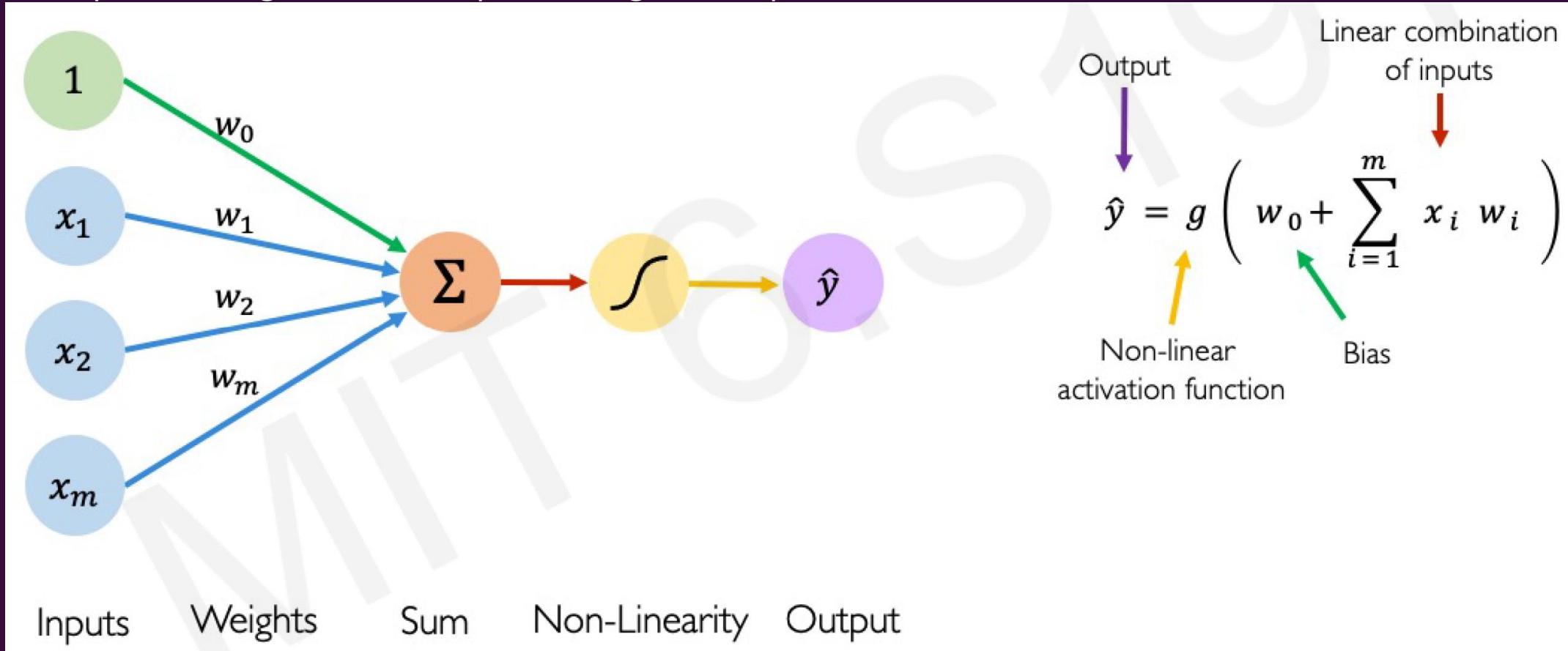
The structural building block of deep learning



Co-funded by
the European Union

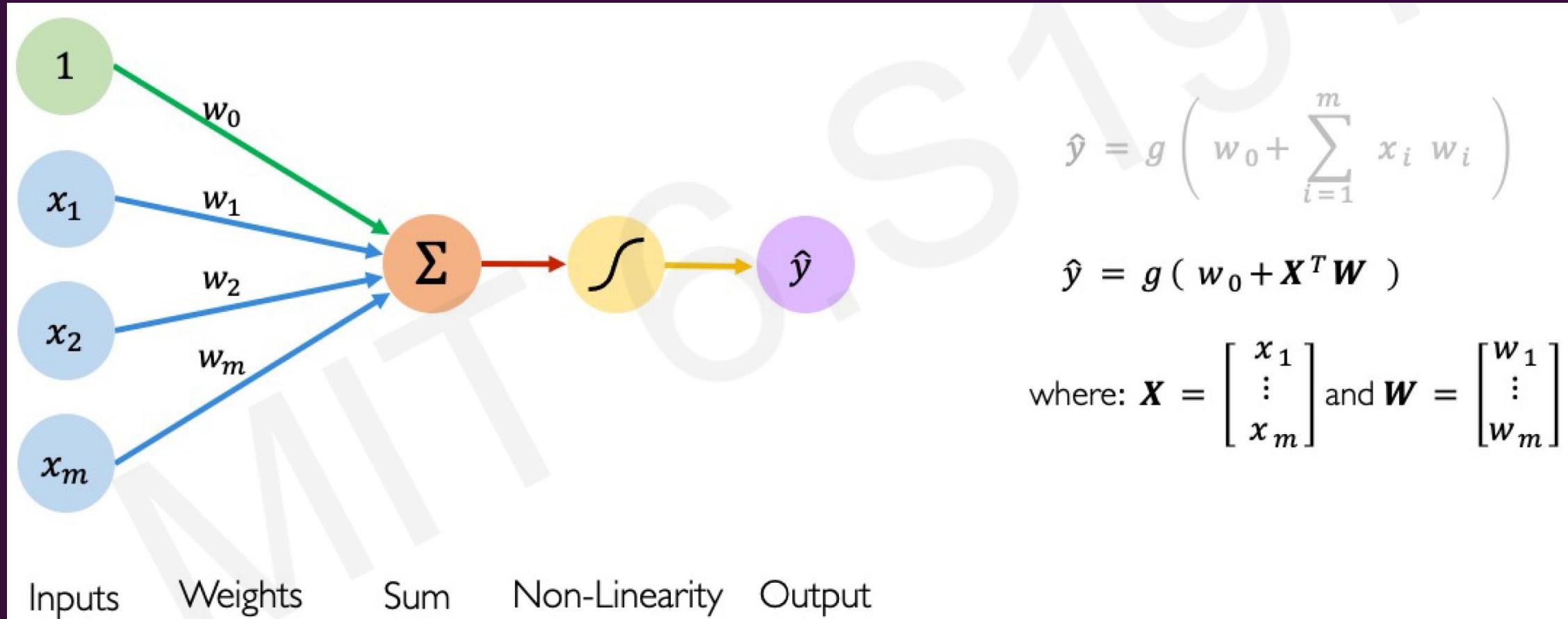
The Perceptron: Forward Propagation

Think of it as a mathematical model of a biological neuron - taking inputs, processing them, and producing an output



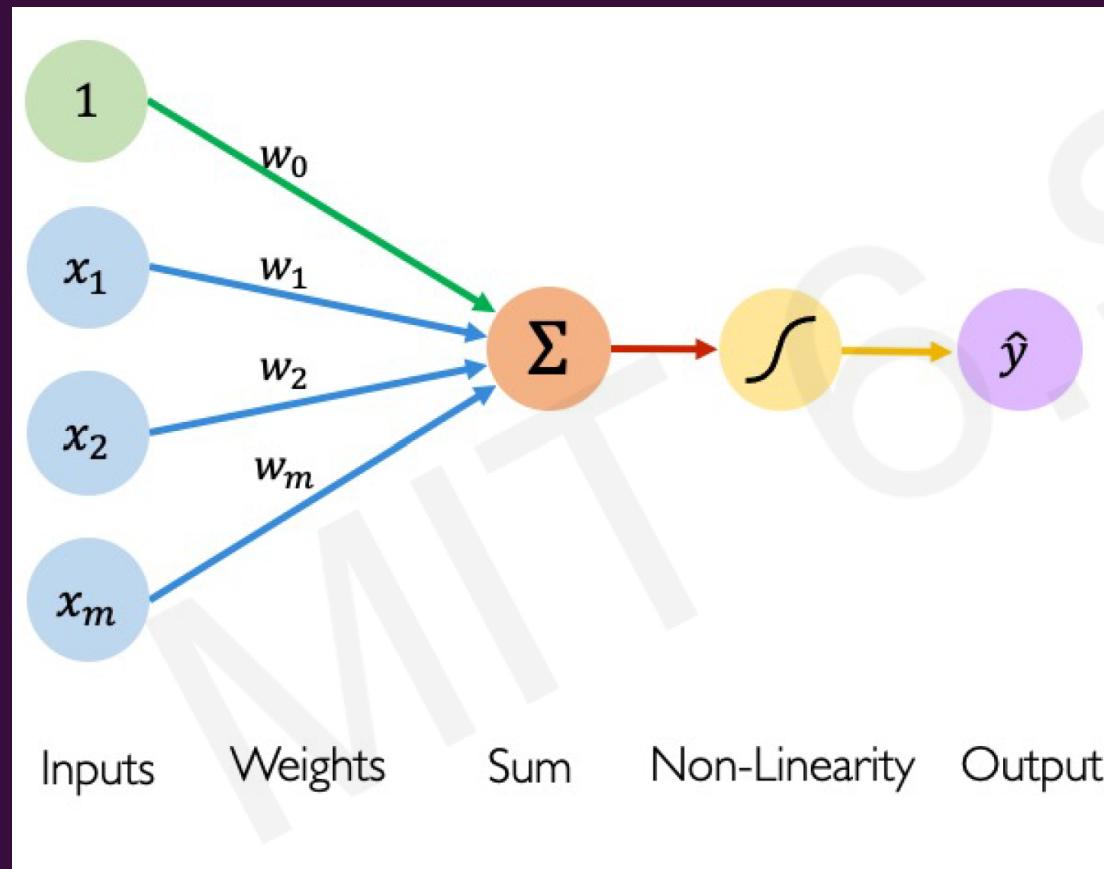
The Perceptron: Forward Propagation

- Each perceptron computes a weighted sum of its inputs plus a bias term



The Perceptron: Forward Propagation

- The sum passes through an activation function that introduces non-linearity into the network

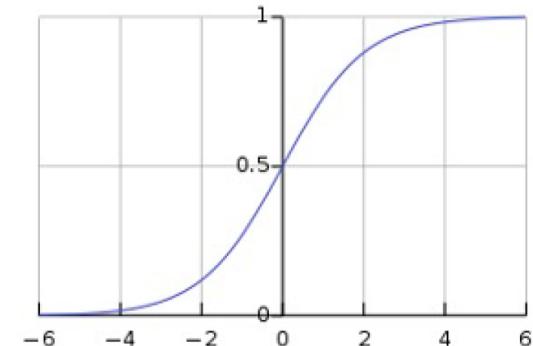


Activation Functions

$$\hat{y} = g(w_0 + \mathbf{X}^T \mathbf{W})$$

- Example: sigmoid function

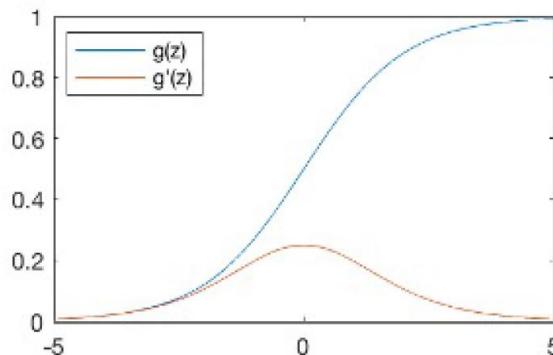
$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$



Common Activation Functions

Activation functions determine whether and how strongly a perceptron "fires" based on its inputs

Sigmoid Function

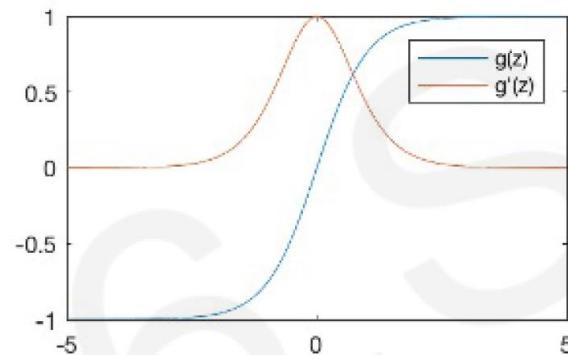


$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

 `tf.math.sigmoid(z)`

Hyperbolic Tangent

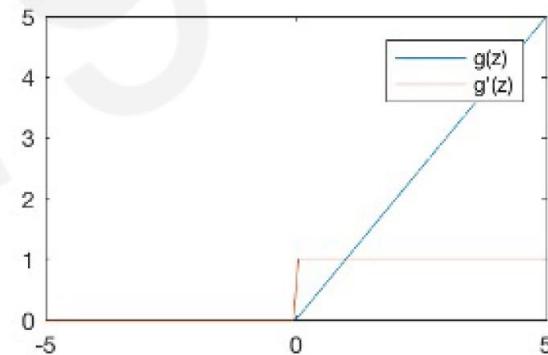


$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

 `tf.math.tanh(z)`

Rectified Linear Unit (ReLU)



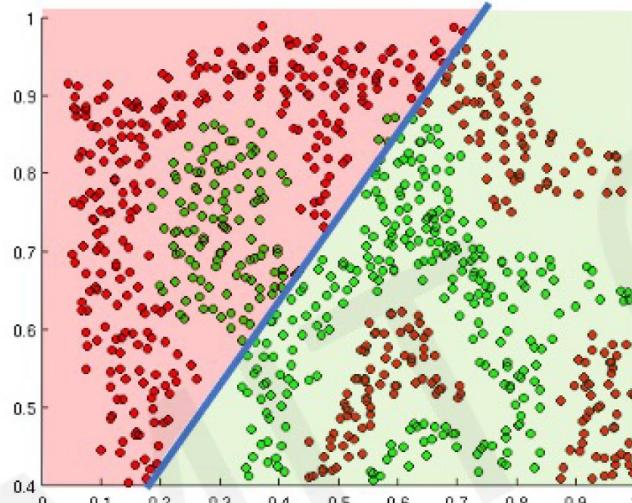
$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

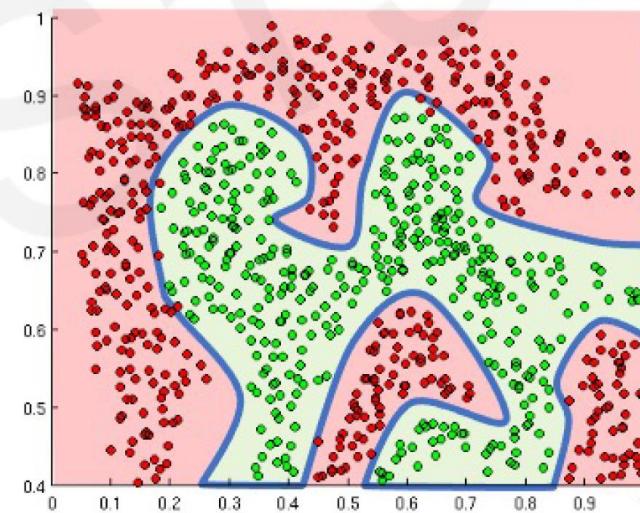
 `tf.nn.relu(z)`

Importance of Activation Functions

*The purpose of activation functions is to **introduce non-linearities** into the network*



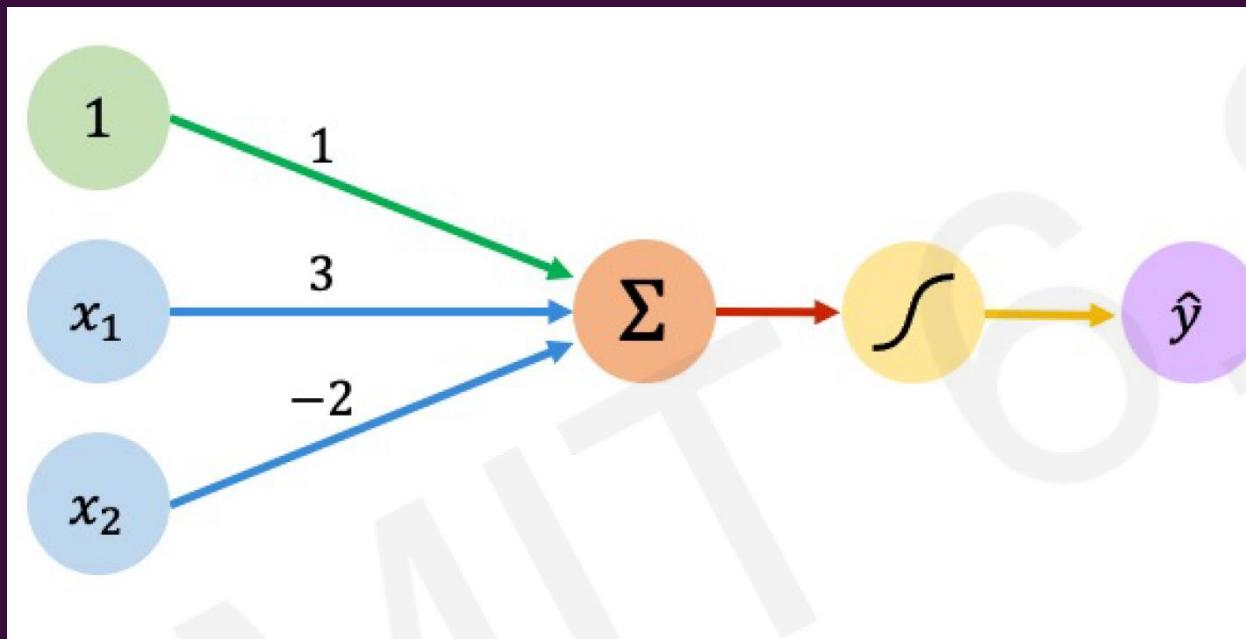
Linear activation functions produce linear decisions no matter the network size



Non-linearities allow us to approximate arbitrarily complex functions

The Perceptron: Example

- In this practical example, we have two input features (x_1, x_2) and specific weights (3, -2)
- The bias term ($w_0 = 1$) helps shift the decision boundary without changing its orientation



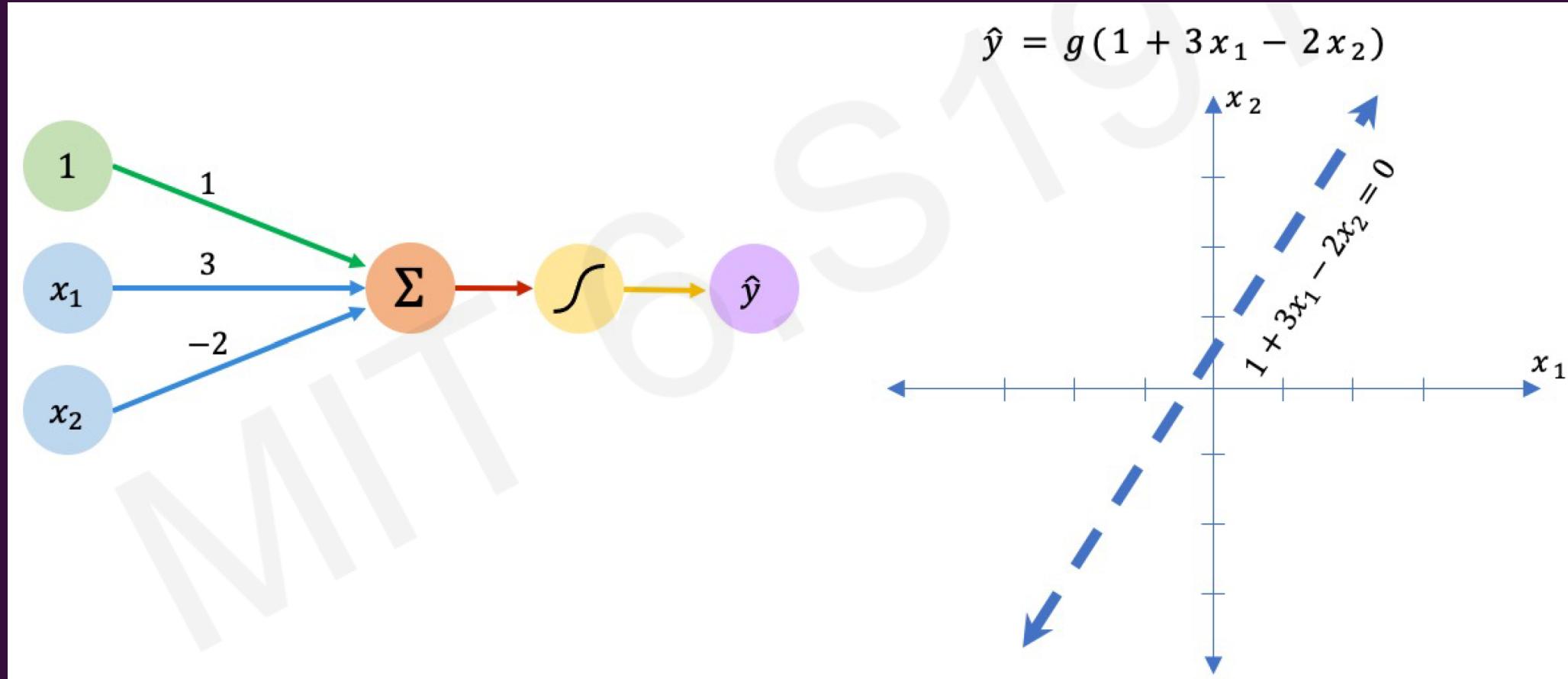
We have: $w_0 = 1$ and $\mathbf{w} = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$

$$\begin{aligned}\hat{y} &= g(w_0 + \mathbf{X}^T \mathbf{w}) \\ &= g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right) \\ \hat{y} &= g\left(1 + 3x_1 - 2x_2\right)\end{aligned}$$

This is just a line in 2D!

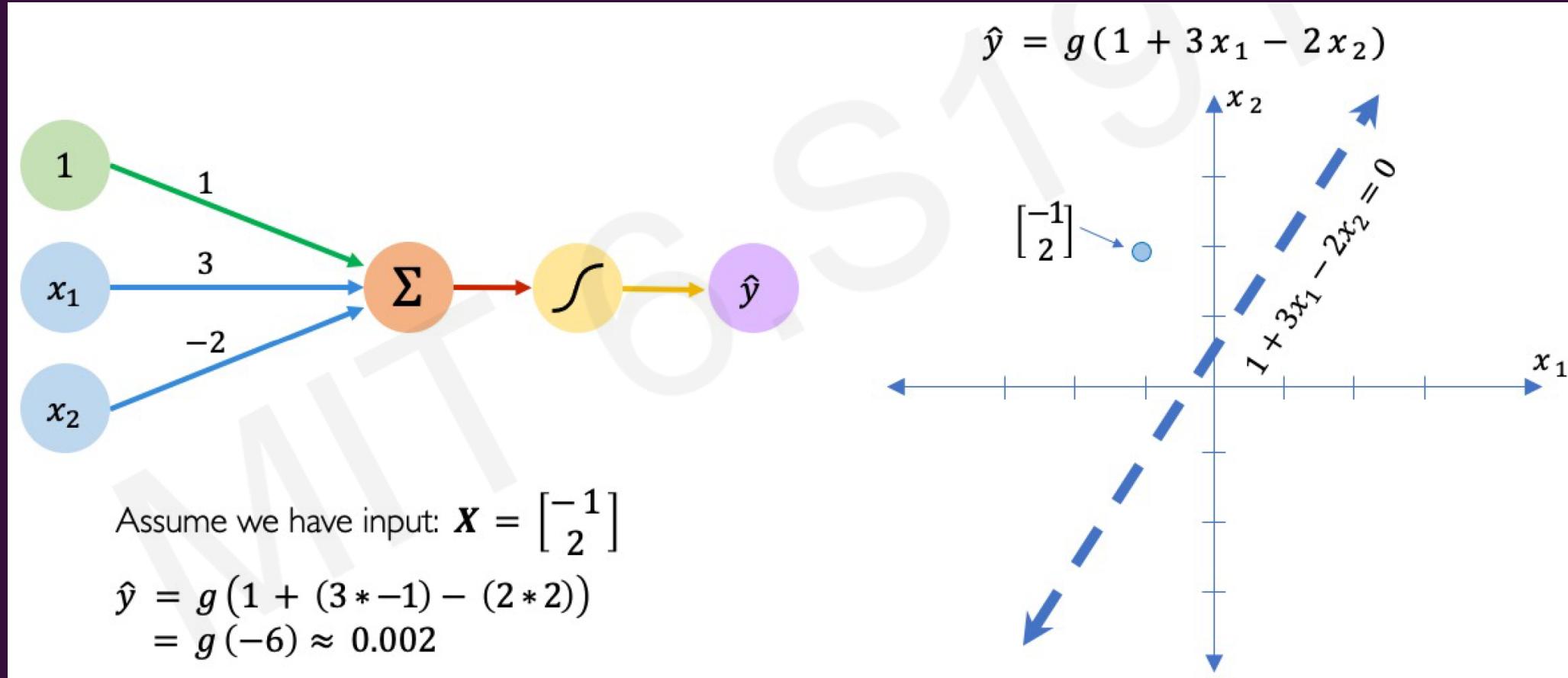
The Perceptron: Example

When mapped in 2D space, the perceptron creates a linear decision boundary



The Perceptron: Example

Understanding how numbers flow through this simple example helps grasp more complex networks

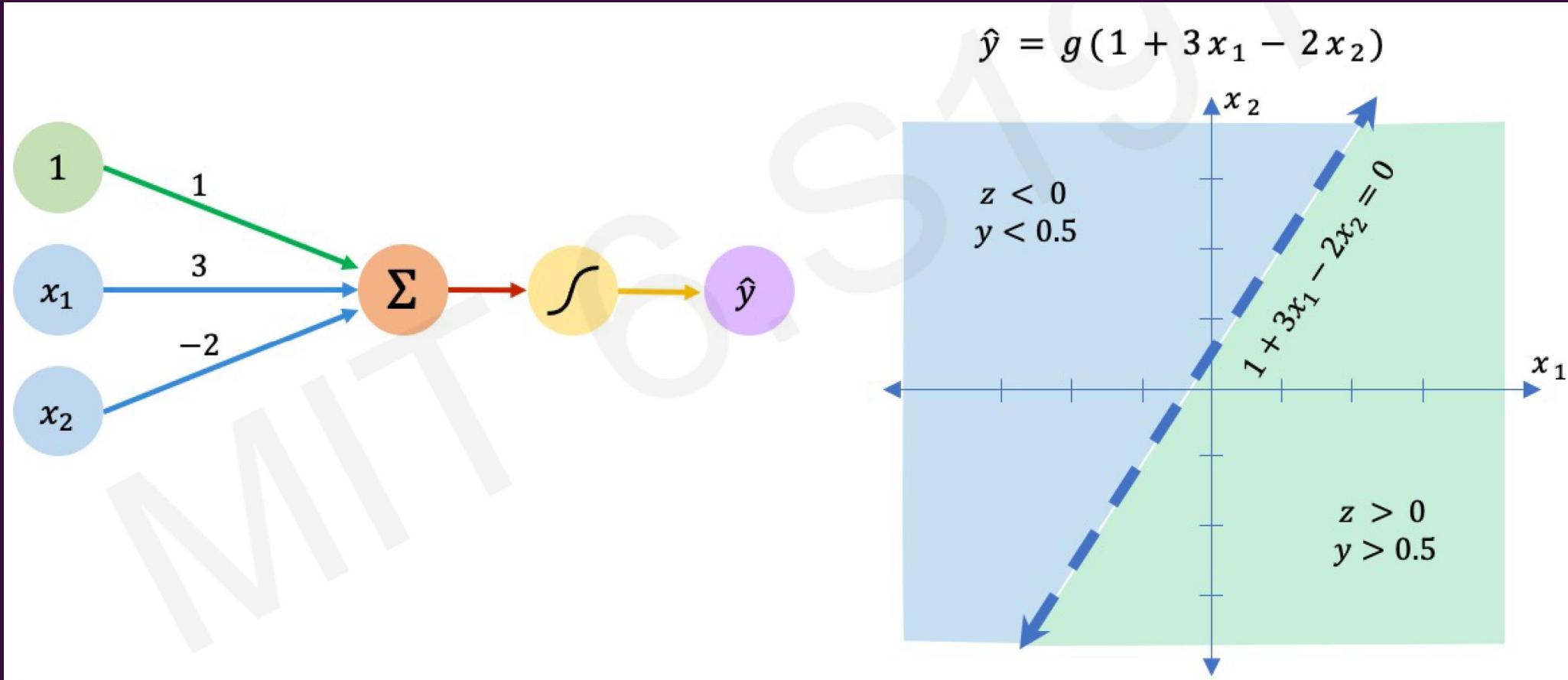


Assume we have input: $X = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$

$$\begin{aligned}\hat{y} &= g(1 + (3 * -1) - (2 * 2)) \\ &= g(-6) \approx 0.002\end{aligned}$$

The Perceptron: Example

Understanding how numbers flow through this simple example helps grasp more complex networks

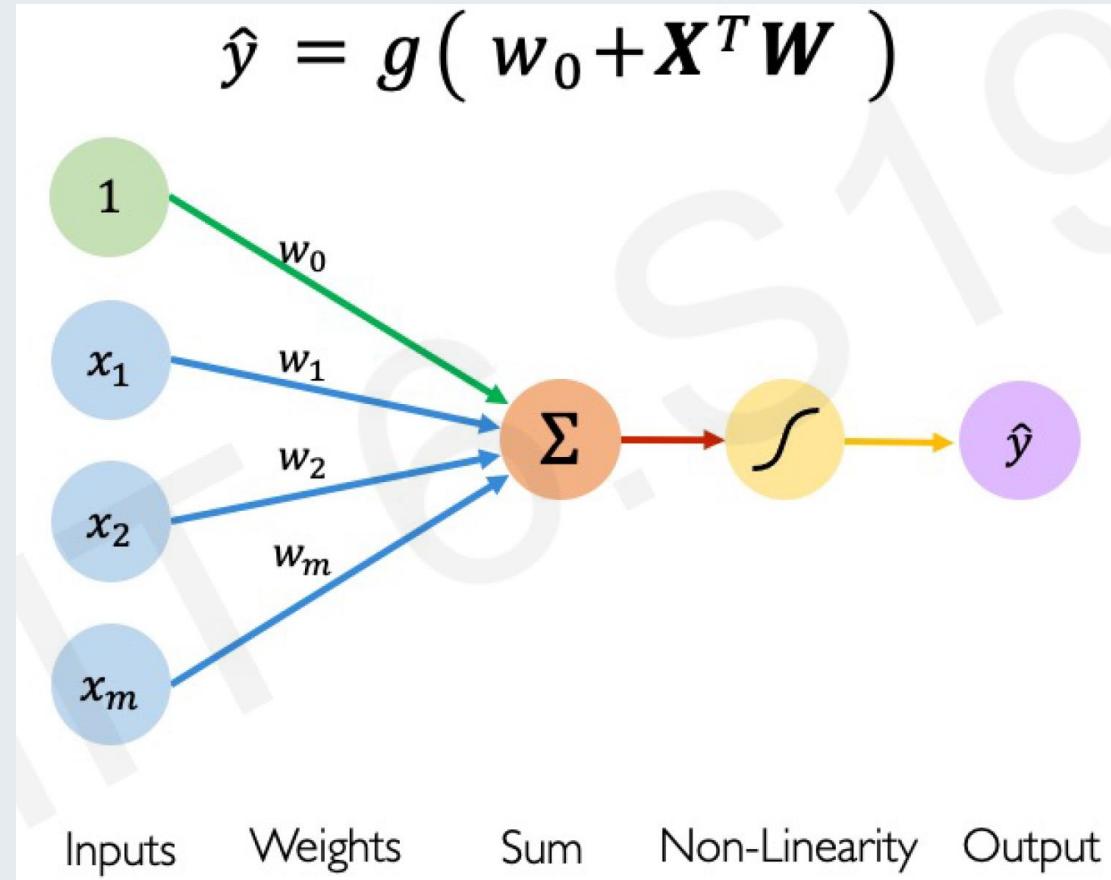


Building Neural Networks with Perceptrons

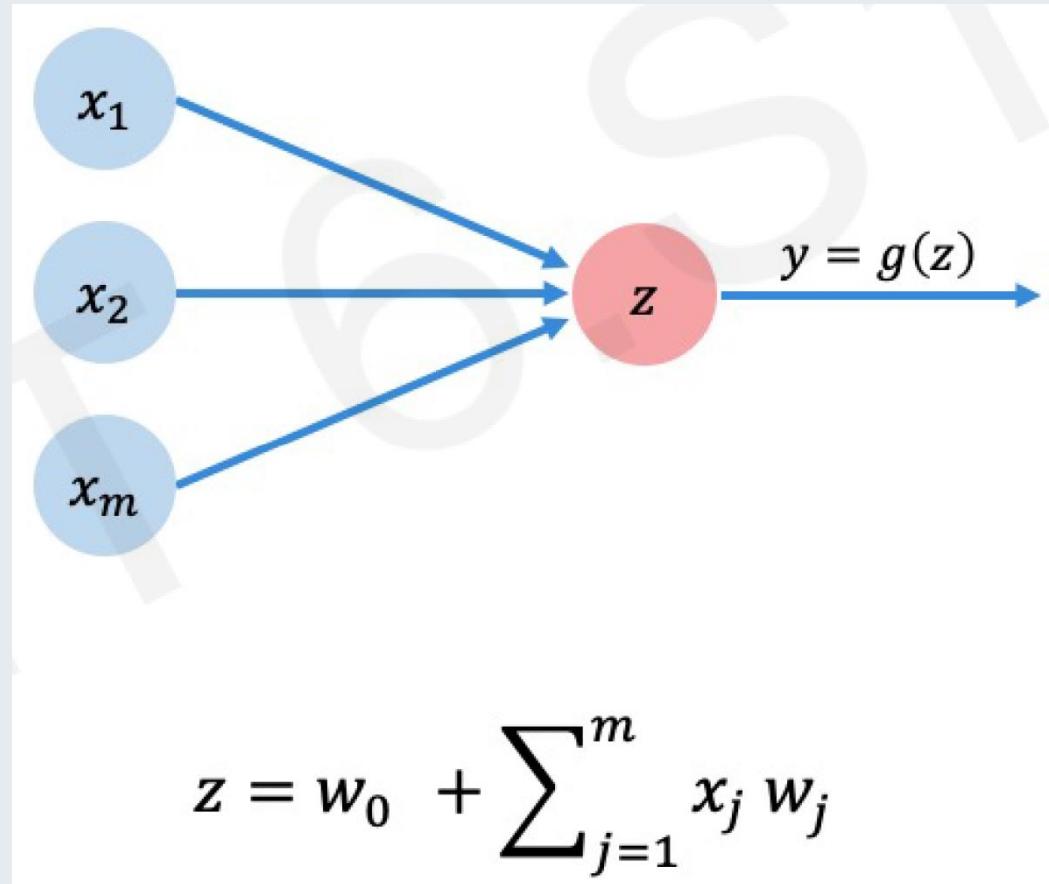


Co-funded by
the European Union

The Perceptron Simplified



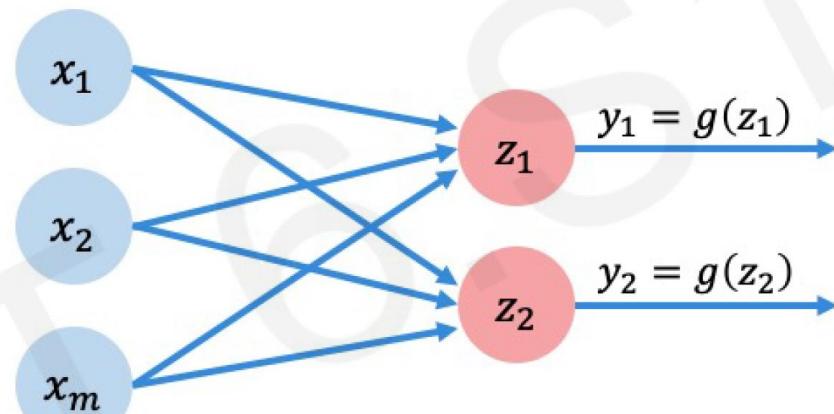
The Perceptron Simplified



Multi Output Perceptron

Multiple perceptrons combine to form layers of a neural network

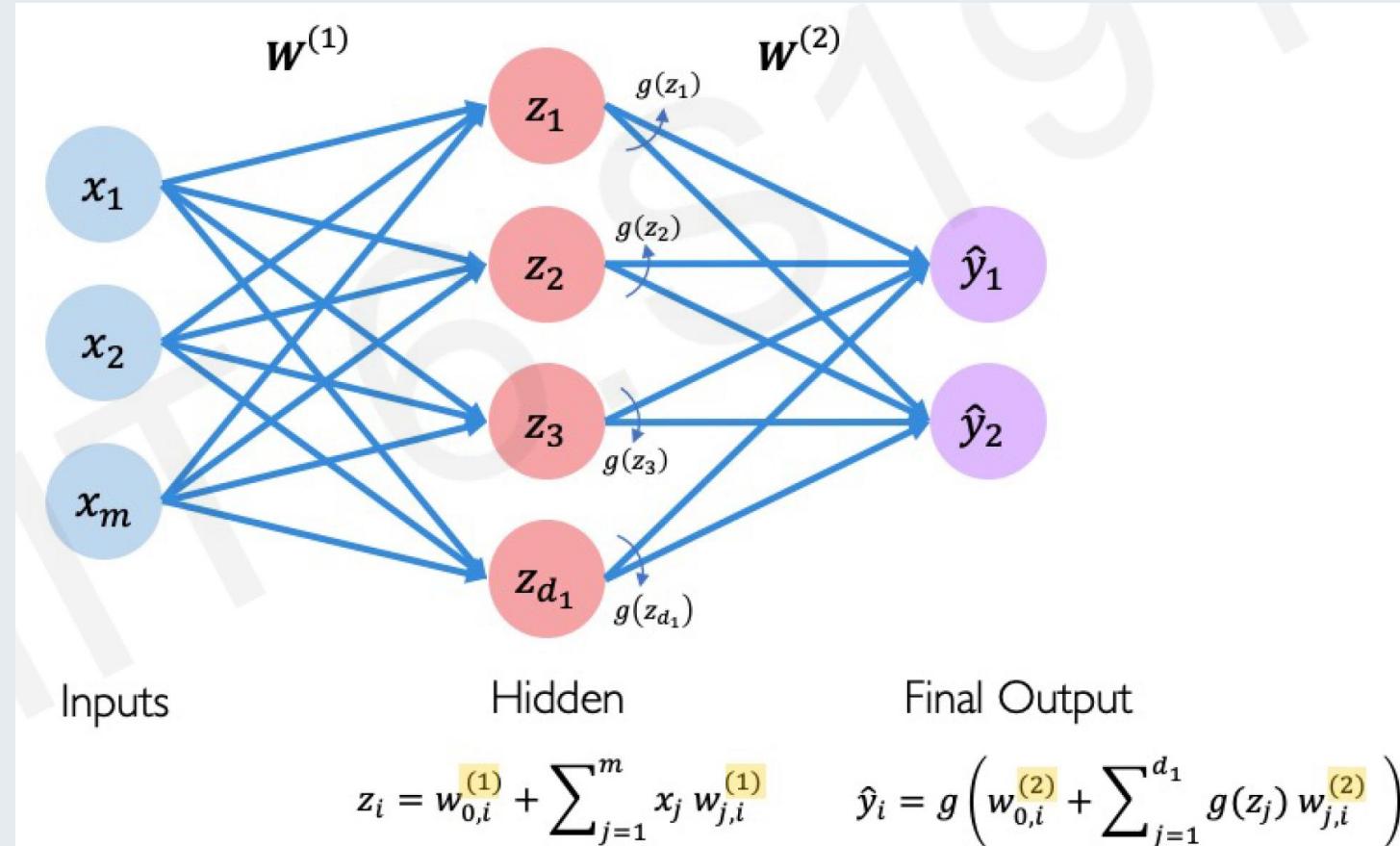
Because all inputs are densely connected to all outputs, these layers are called **Dense** layers



$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

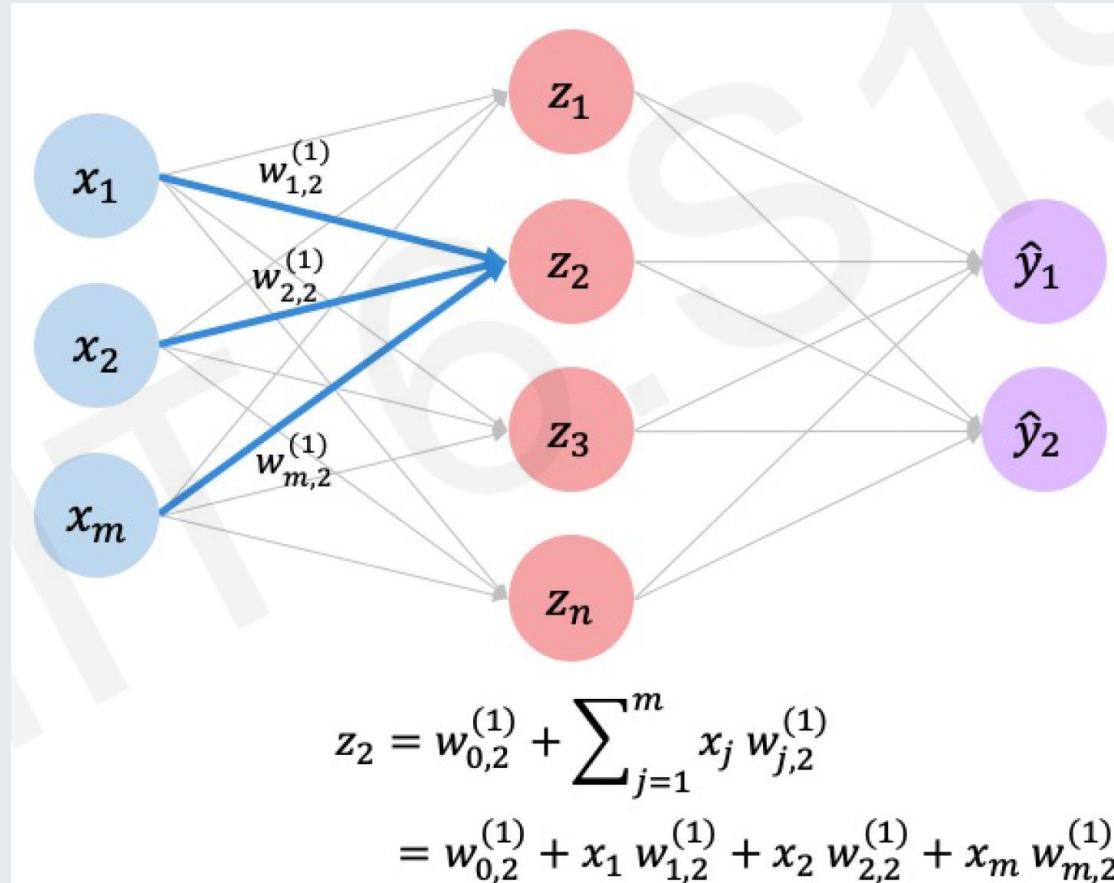
Single Layer Neural Network

Each connection between neurons has its own weight that gets learned during training



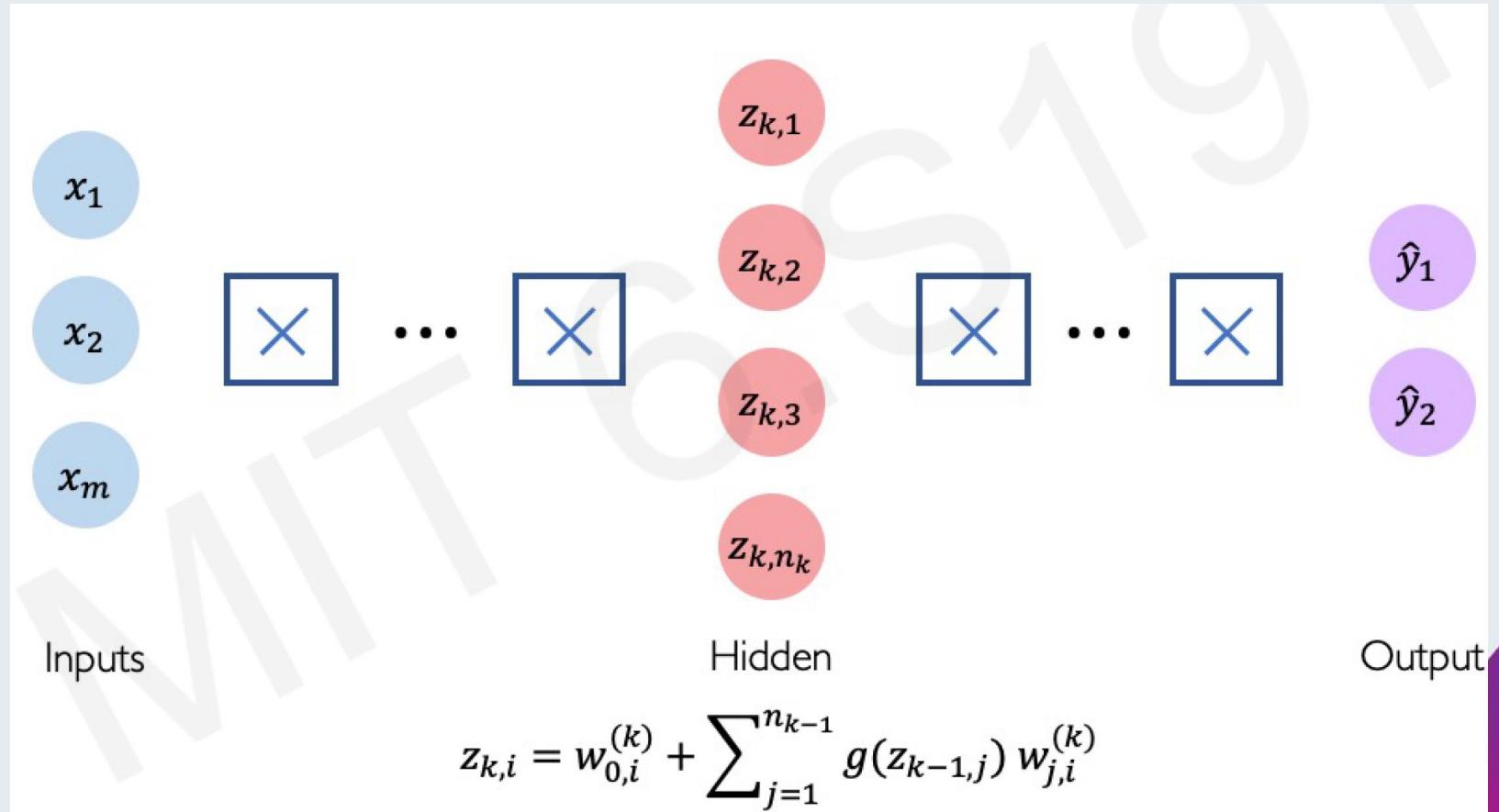
Single Layer Neural Network

Each connection between neurons has its own weight that gets learned during training



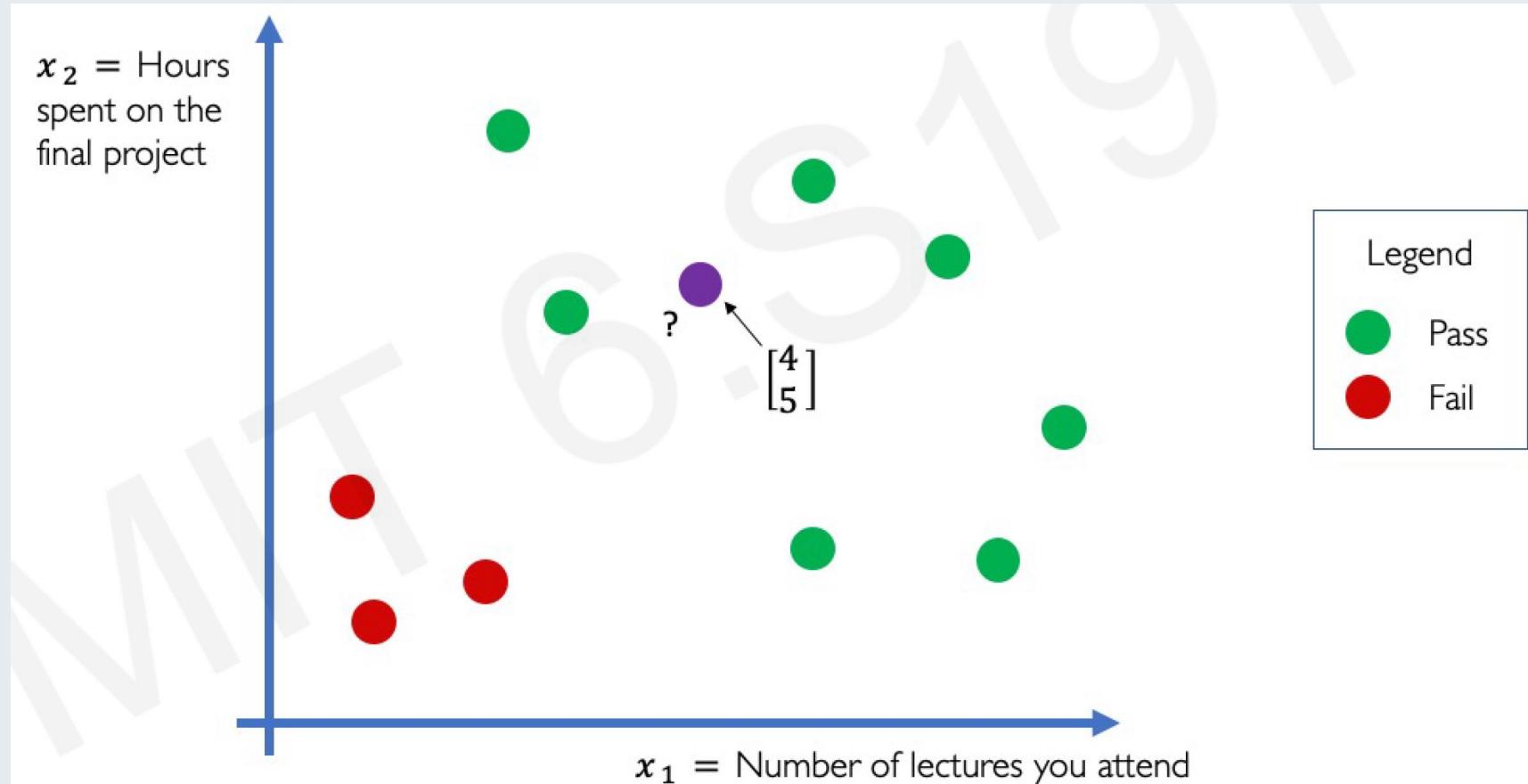
Deep Neural Network

- A single-layer neural network can model only simple linear patterns, whereas a deep neural network can capture highly complex patterns.
- Each connection between neurons has a weight, which is learned through training using optimization techniques like gradient descent.



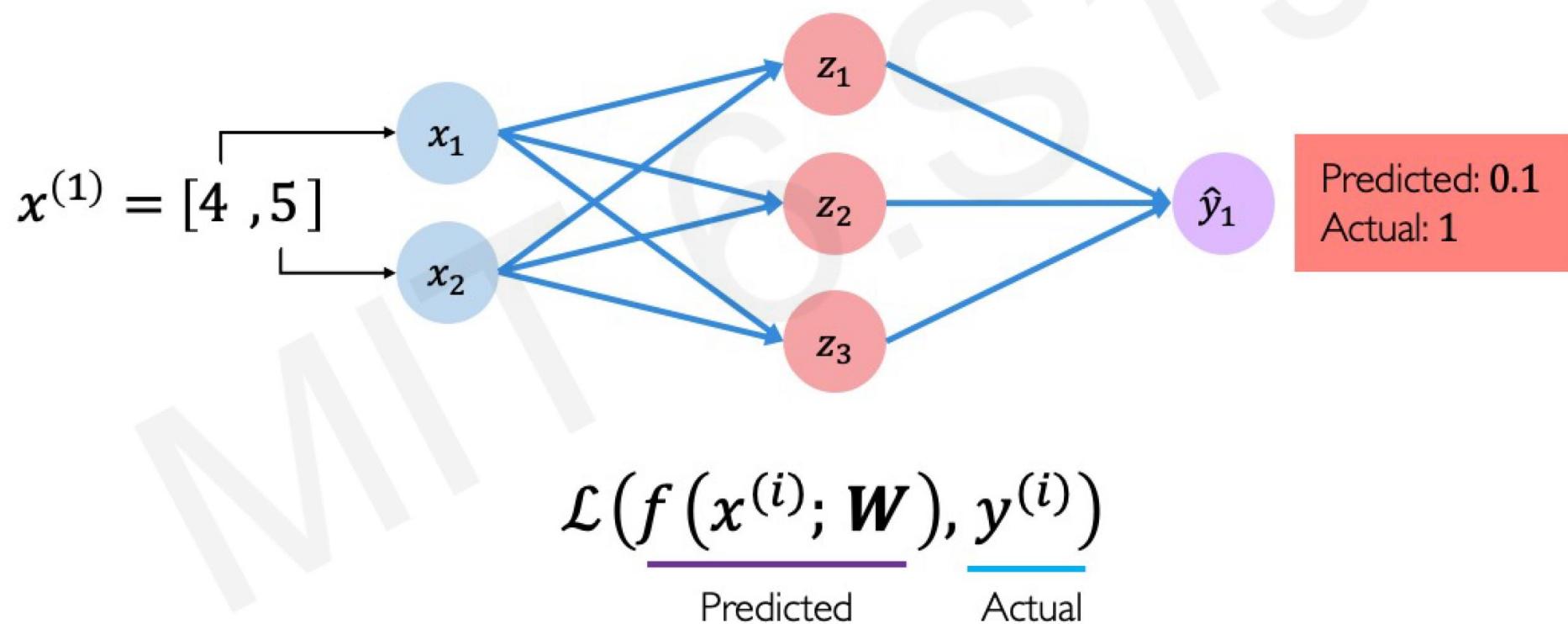
Example Problem: Will I Pass This Class?

The network's depth (number of layers) determines how complex patterns it can learn



Example Problem: Will I Pass This Class?

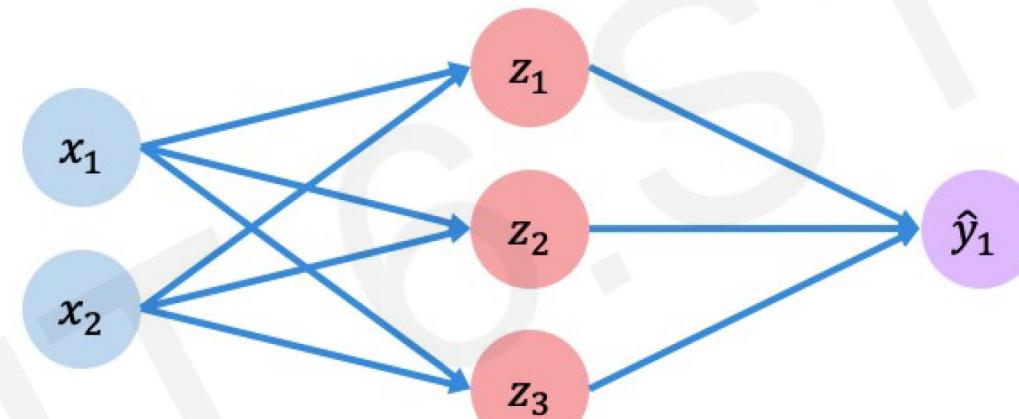
The **loss** of our network measures the cost incurred from incorrect predictions



Empirical Loss

The **empirical loss** measures the total loss over our entire dataset

$$\mathbf{X} = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$



$f(x)$	y
0.1	✗
0.8	✗
0.6	✓
\vdots	\vdots

Also known as:

- Objective function
- Cost function
- Empirical Risk

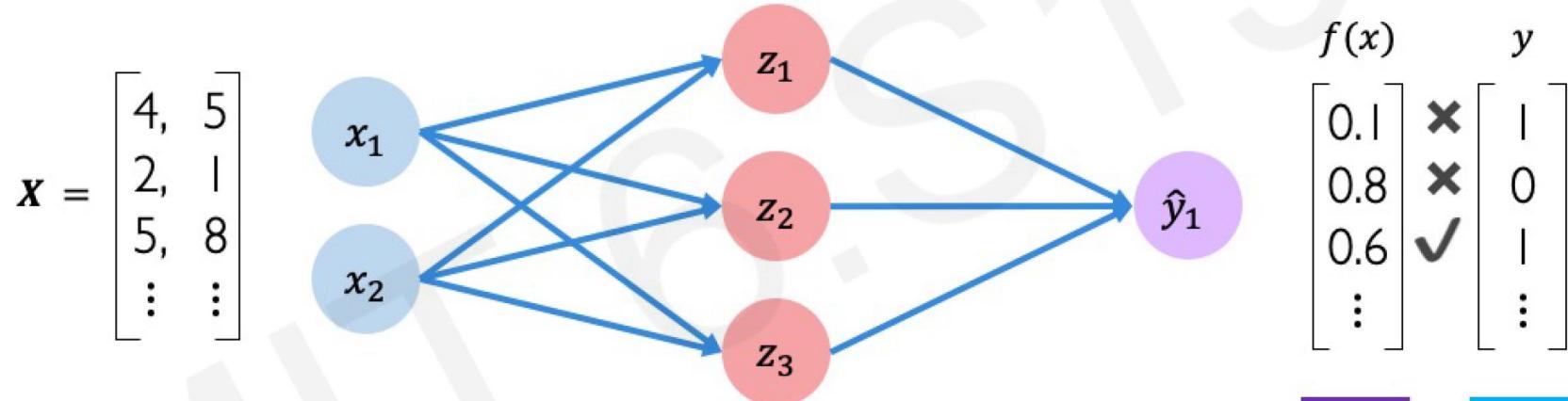
$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(\mathbf{x}^{(i)}; \mathbf{W}), \mathbf{y}^{(i)})$$

Predicted Actual



Binary Cross Entropy Loss

Cross entropy loss can be used with models that output a probability between 0 and 1



$$J(\mathbf{W}) = -\frac{1}{n} \sum_{i=1}^n \underbrace{y^{(i)} \log(f(x^{(i)}; \mathbf{W}))}_{\text{Actual}} + \underbrace{(1 - y^{(i)}) \log(1 - f(x^{(i)}; \mathbf{W}))}_{\text{Predicted}}$$

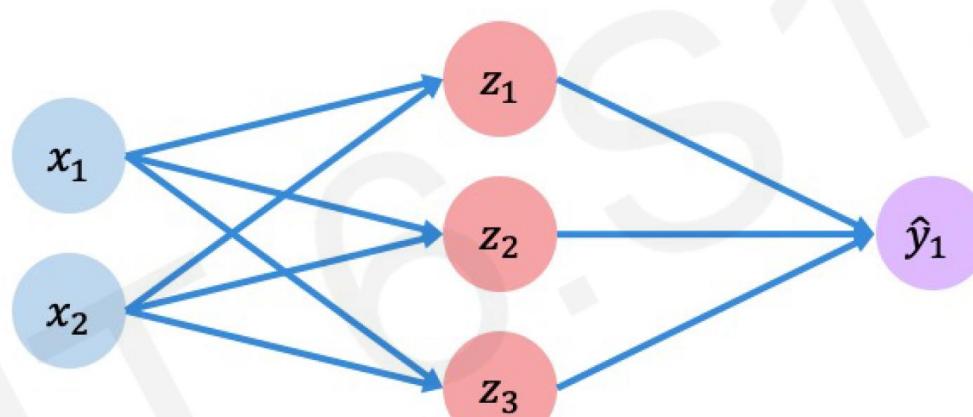


```
loss = tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(y, predicted))
```

Mean Squared Error Loss

Mean squared error loss can be used with regression models that output continuous real numbers

$$\mathbf{x} = \begin{bmatrix} 4, & 5 \\ 2, & 1 \\ 5, & 8 \\ \vdots & \vdots \end{bmatrix}$$



$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \left(\underline{y^{(i)}} - \underline{f(x^{(i)}; \mathbf{W})} \right)^2$$

Actual Predicted

$f(x)$	y
30	✗ 90
80	✗ 20
85	✓ 95
⋮	⋮

Final Grades
(percentage)



```
loss = tf.reduce_mean(tf.square(tf.subtract(y, predicted)) )
loss = tf.keras.losses.MSE( y, predicted )
```

Training Neural Networks



Co-funded by
the European Union

Loss Optimization

We want to find the network weights that **achieve the lowest loss**

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$

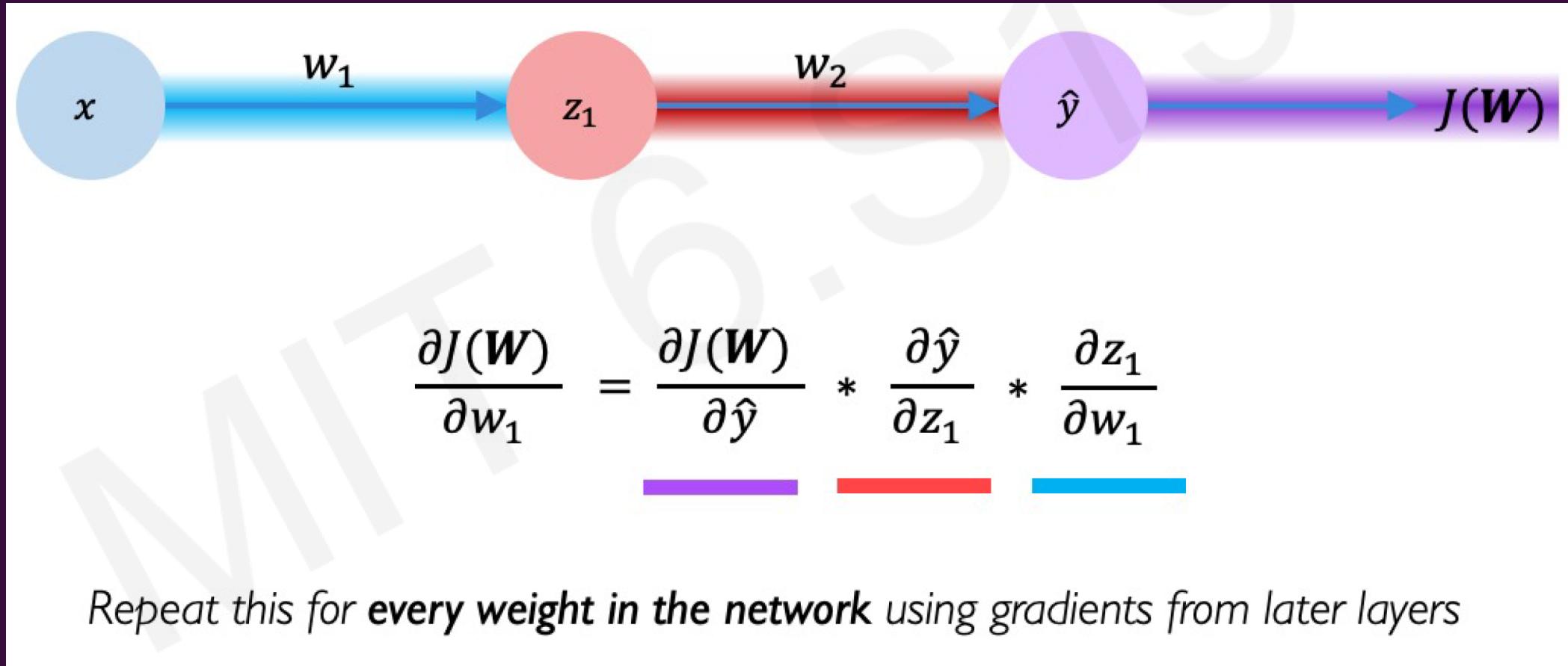


Remember:

$$\mathbf{W} = \{\mathbf{W}^{(0)}, \mathbf{W}^{(1)}, \dots\}$$

Computing Gradients: Backpropagation

Backpropagation calculates how much each weight contributed to the error

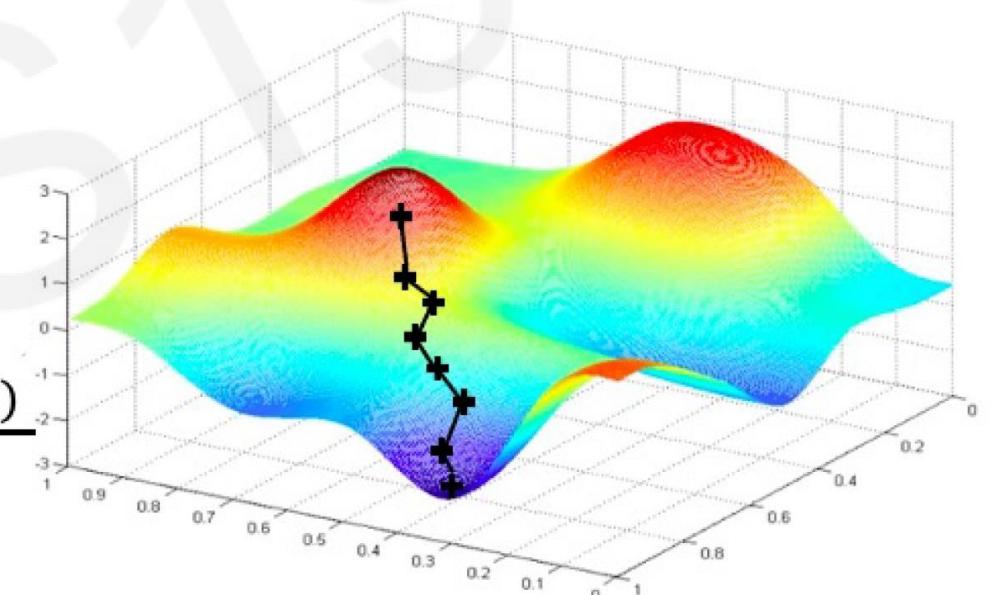


Stochastic Gradient Descent

Gradient descent updates weights in the direction that reduces error

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights

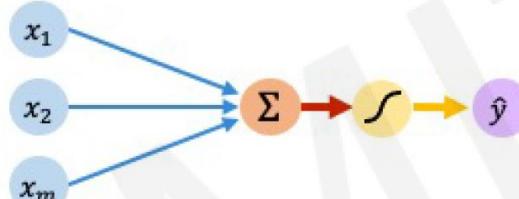


Neural Networks Summary

Training requires finding the right balance between learning rate, batch size, and number of epochs

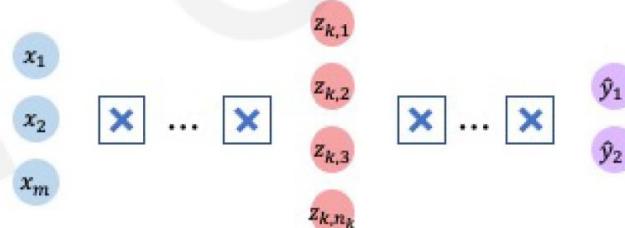
The Perceptron

- Structural building blocks
- Nonlinear activation functions



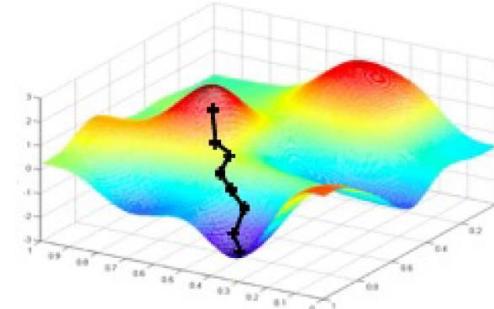
Neural Networks

- Stacking Perceptrons to form neural networks
- Optimization through backpropagation



Training in Practice

- Adaptive learning
- Batching
- Regularization



Deep Learning for Computer Vision

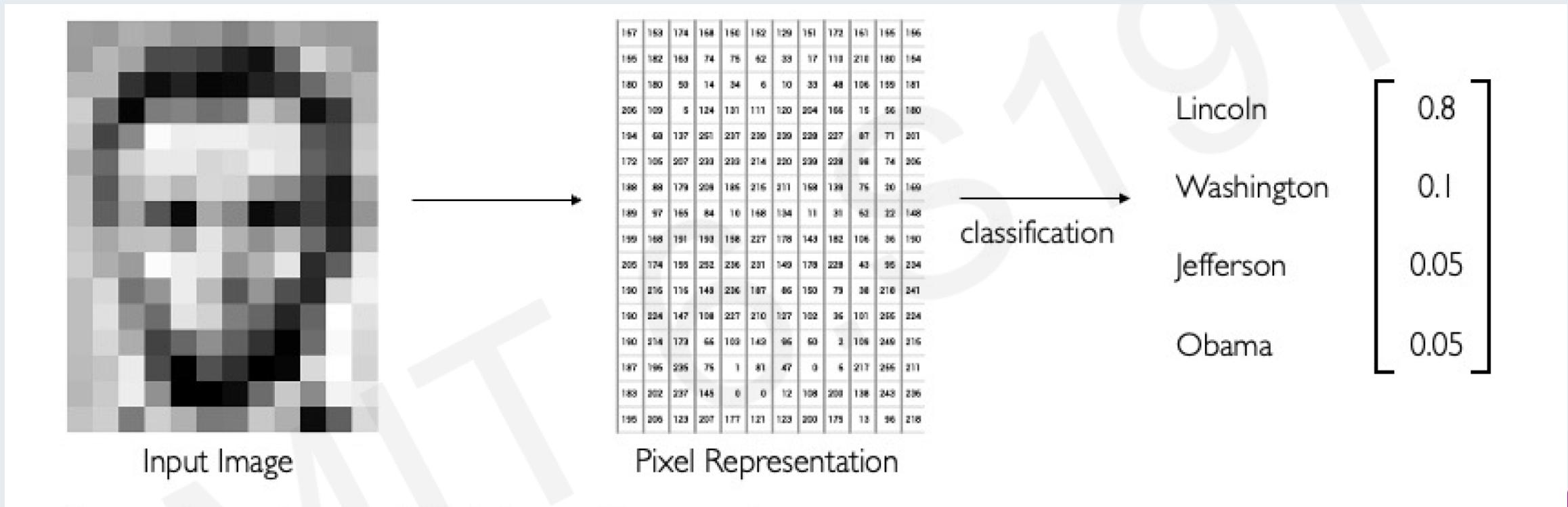


Co-funded by
the European Union



Tasks in Computer Vision

Computer vision tasks require understanding spatial relationships in image data



- **Regression:** output variable takes continuous value
- **Classification:** output variable takes class label. Can produce probability of belonging to a particular class



Manual Feature Extraction

Domain knowledge

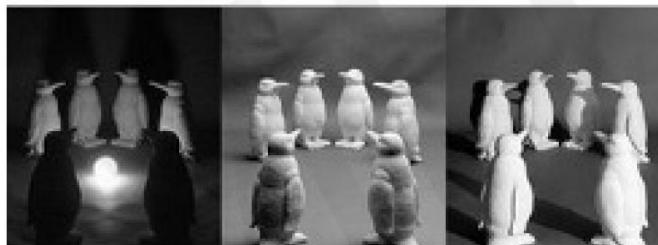
Define features

Detect features
to classify

Viewpoint variation



Illumination conditions



Scale variation



Deformation



Background clutter



Occlusion



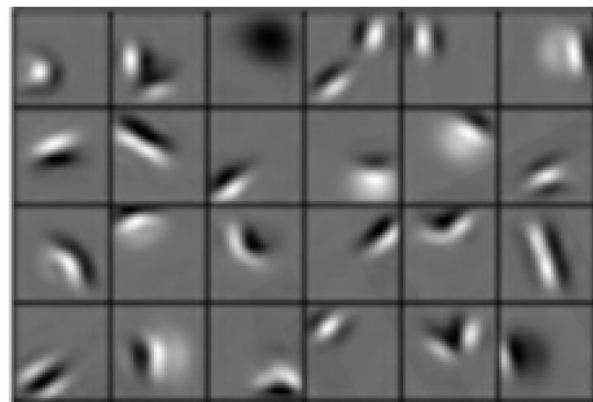
Intra-class variation



Learning Feature Representations

Can we learn a **hierarchy of features** directly from the data instead of hand engineering?

Low level features



Edges, dark spots

Mid level features



Eyes, ears, nose

High level features



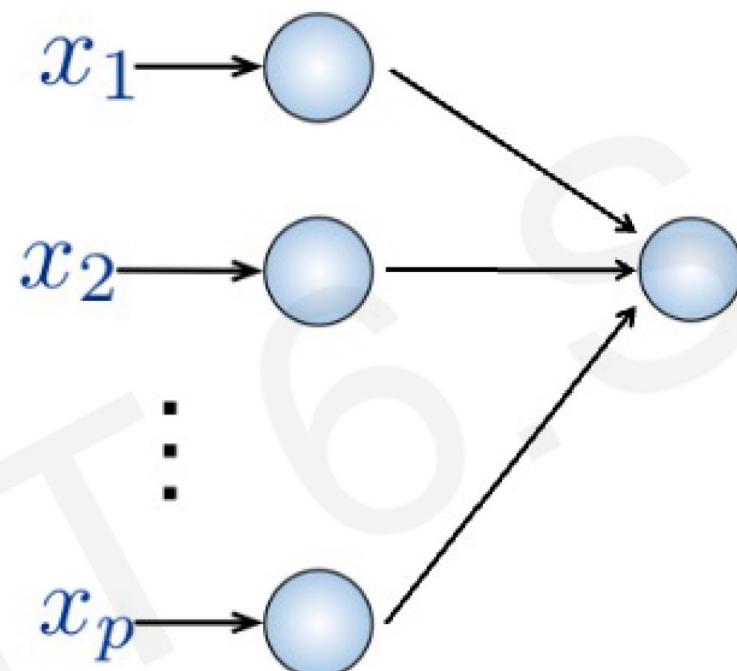
Facial structure

Fully Connected Neural Network

Traditional neural networks don't handle spatial data efficiently

Input:

- 2D image
- Vector of pixel values



Fully Connected:

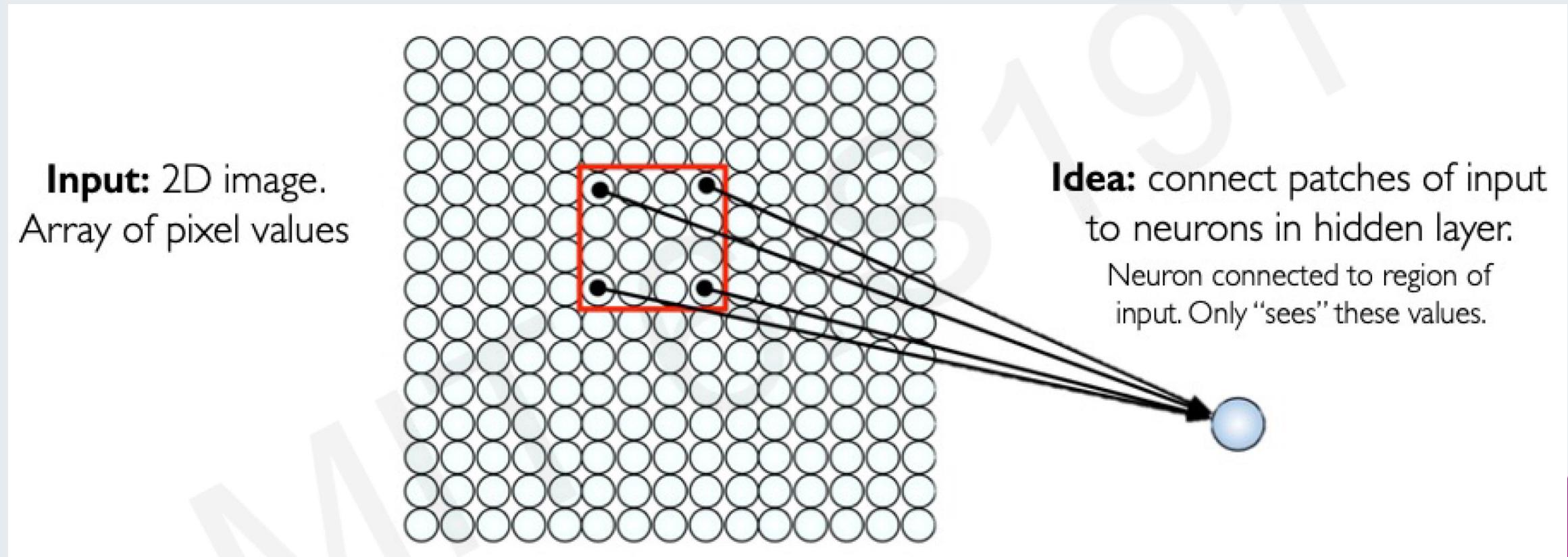
- Connect neuron in hidden layer to all neurons in input layer
- No spatial information!
- And many, many parameters!

How can we use **spatial structure** in the input to inform the architecture of the network?



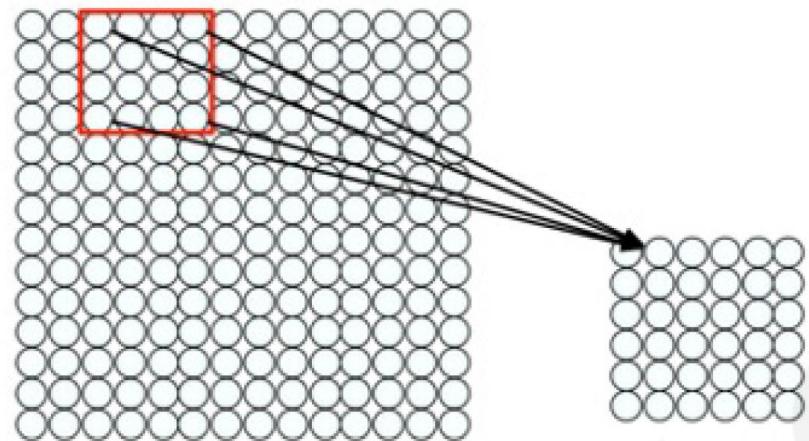
Using Spatial Structure

Convolutional layers learn filters that detect features like edges, textures, and patterns



Feature Extraction with Convolution

The network progressively learns more complex visual features in deeper layers



- Filter of size 4×4 : 16 different weights
- Apply this same filter to 4×4 patches in input
- Shift by 2 pixels for next patch

This “patchy” operation is **convolution**

- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) **Spatially share** parameters of each filter

Convolutional Neural Networks

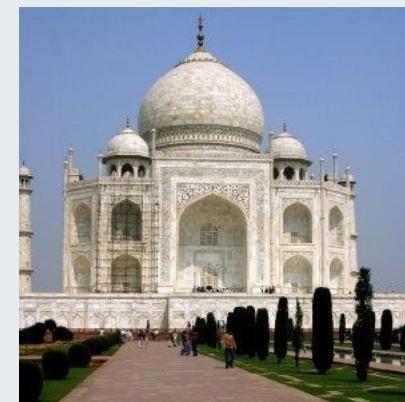
Key components include filters (kernels), pooling layers, and fully connected layers for classification.

1 <small>x1</small>	1 <small>x0</small>	1 <small>x1</small>	0	0
0 <small>x0</small>	1 <small>x1</small>	1 <small>x0</small>	1	0
0 <small>x1</small>	0 <small>x0</small>	1 <small>x1</small>	1	1
0	0	1	1	0
0	1	1	0	0

Image

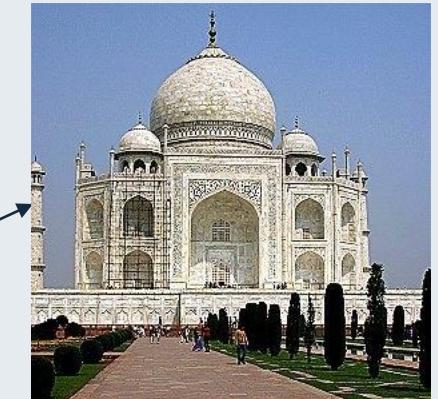
4		

Convolved Feature



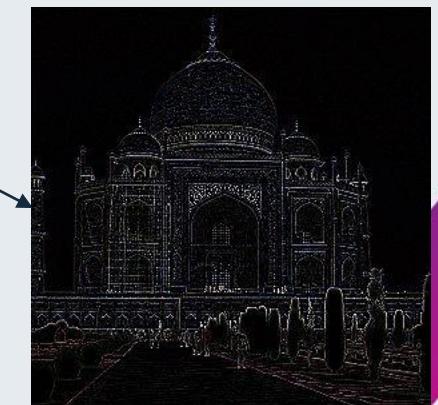
sharpen

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



0	1	0
1	-4	1
0	1	0

detect edges



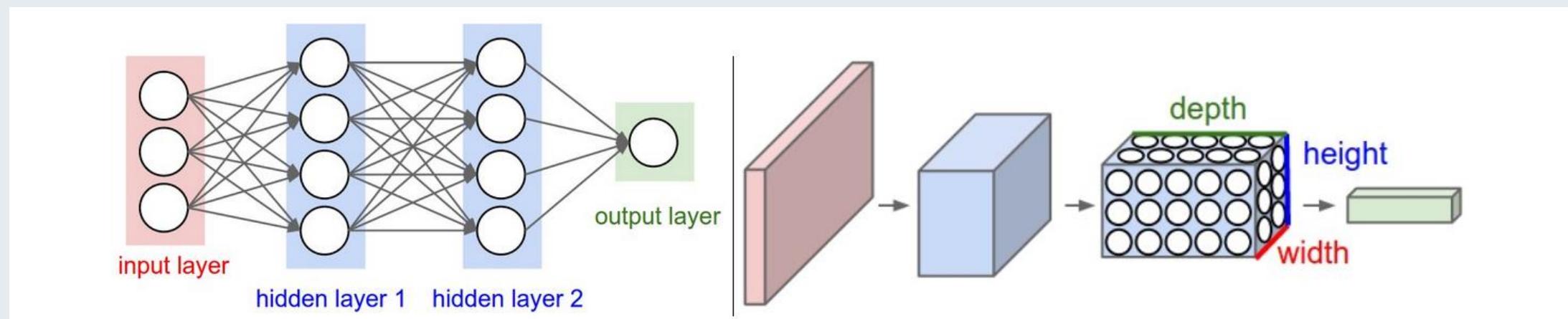
<http://www.wildml.com/2015/11/understanding-convolutional-neural-networks-for-nlp/>



Co-funded by
the European Union

Convolutional Neural Networks

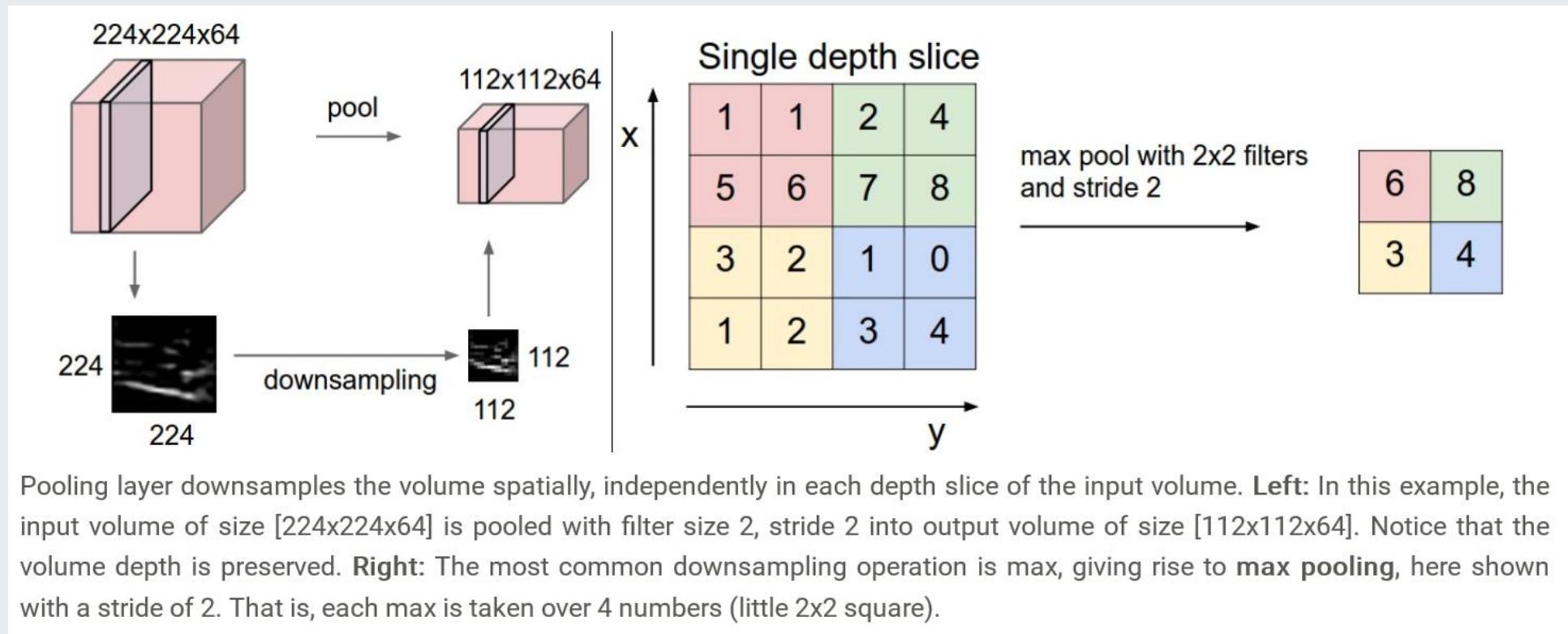
- <http://cs231n.github.io/convolutional-networks/>



Left: A regular 3-layer Neural Network. Right: A ConvNet arranges its neurons in three dimensions (width, height, depth), as visualized in one of the layers. Every layer of a ConvNet transforms the 3D input volume to a 3D output volume of neuron activations. In this example, the red input layer holds the image, so its width and height would be the dimensions of the image, and the depth would be 3 (Red, Green, Blue channels).

Convolutional Neural Networks

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



Supervised vs Unsupervised

- Supervised learning uses labeled data for training, while unsupervised learning identifies patterns in unlabeled data.
- Deep generative models, such as Autoencoders and Generative Adversarial Networks (GANs), are commonly used for unsupervised learning.

Supervised Learning

Data: (x, y)

x is data, y is label

Goal: Learn function to map
 $x \rightarrow y$

Examples: Classification,
 regression, object detection,
 semantic segmentation, etc.

Unsupervised Learning

Data: x

x is data, no labels!

Goal: Learn the *hidden or underlying structure* of the data

Examples: Clustering, feature or dimensionality reduction, etc.

Generative models learn the underlying distribution of data to create new samples that resemble the original dataset.

Goal: Take as input training samples from some distribution and learn a model that represents that distribution

Density Estimation



Sample Generation



Input samples

Training data $\sim P_{data}(x)$



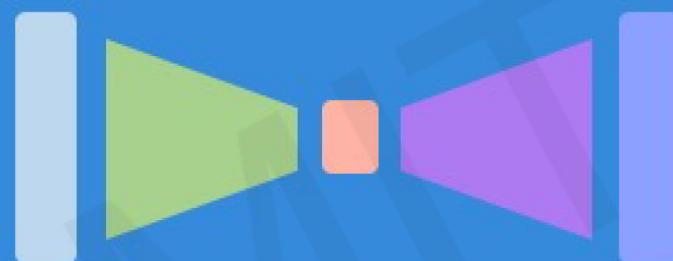
Generated samples

Generated $\sim P_{model}(x)$

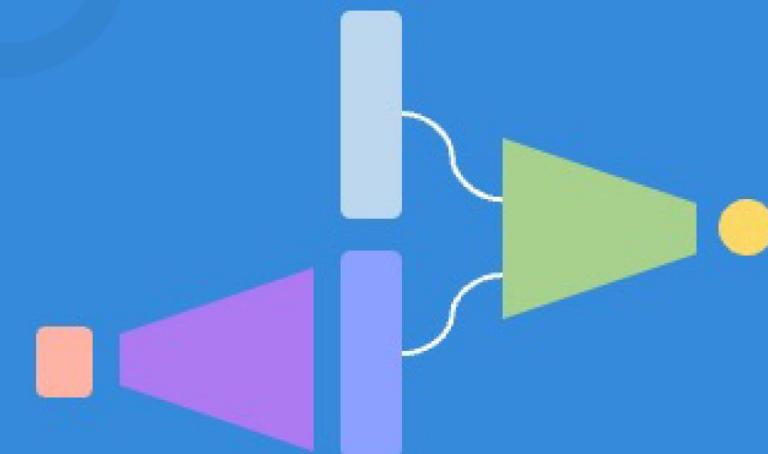
How can we learn $P_{model}(x)$ similar to $P_{data}(x)$?

Latent Variable Models

Autoencoders and Variational Autoencoders (VAEs)



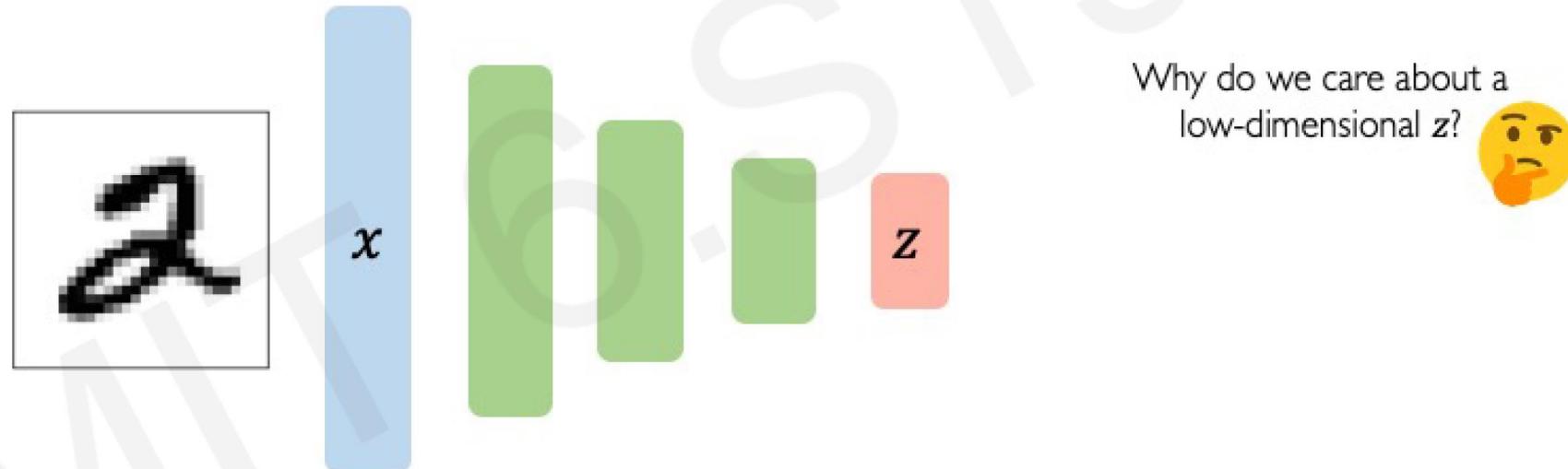
Generative Adversarial Networks (GANs)



Autoencoders: Background

Autoencoders compress input data into a smaller representation (latent space) and reconstruct it, often used for anomaly detection and feature learning.

Unsupervised approach for learning a **lower-dimensional** feature representation from unlabeled training data

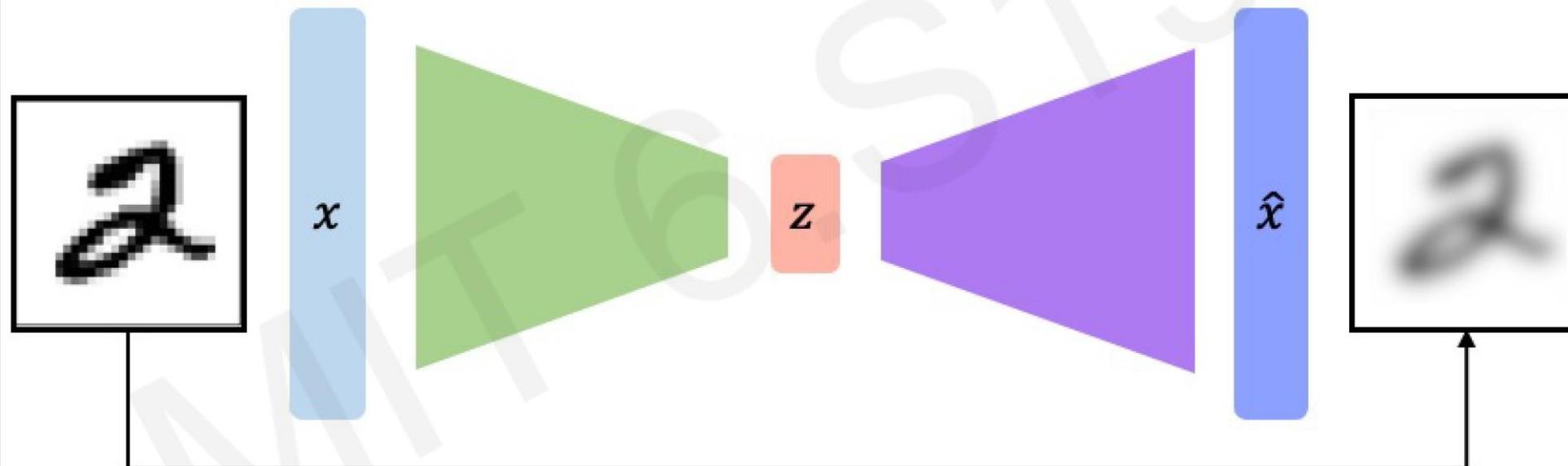


"Encoder" learns mapping from the data, x , to a low-dimensional latent space, z

Autoencoders: Background

How can we learn this latent space?

Train the model to use these features to **reconstruct the original data**



$$\mathcal{L}(x, \hat{x}) = \|x - \hat{x}\|^2$$

Loss function doesn't
use any labels!!

Dimensionality of Latent Space -> Reconstruction Quality

Autoencoding is a form of compression!
Smaller latent space will force a larger training bottleneck

2D latent space

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

5D latent space

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

Ground Truth

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

Variational Autoencoder (VAE)

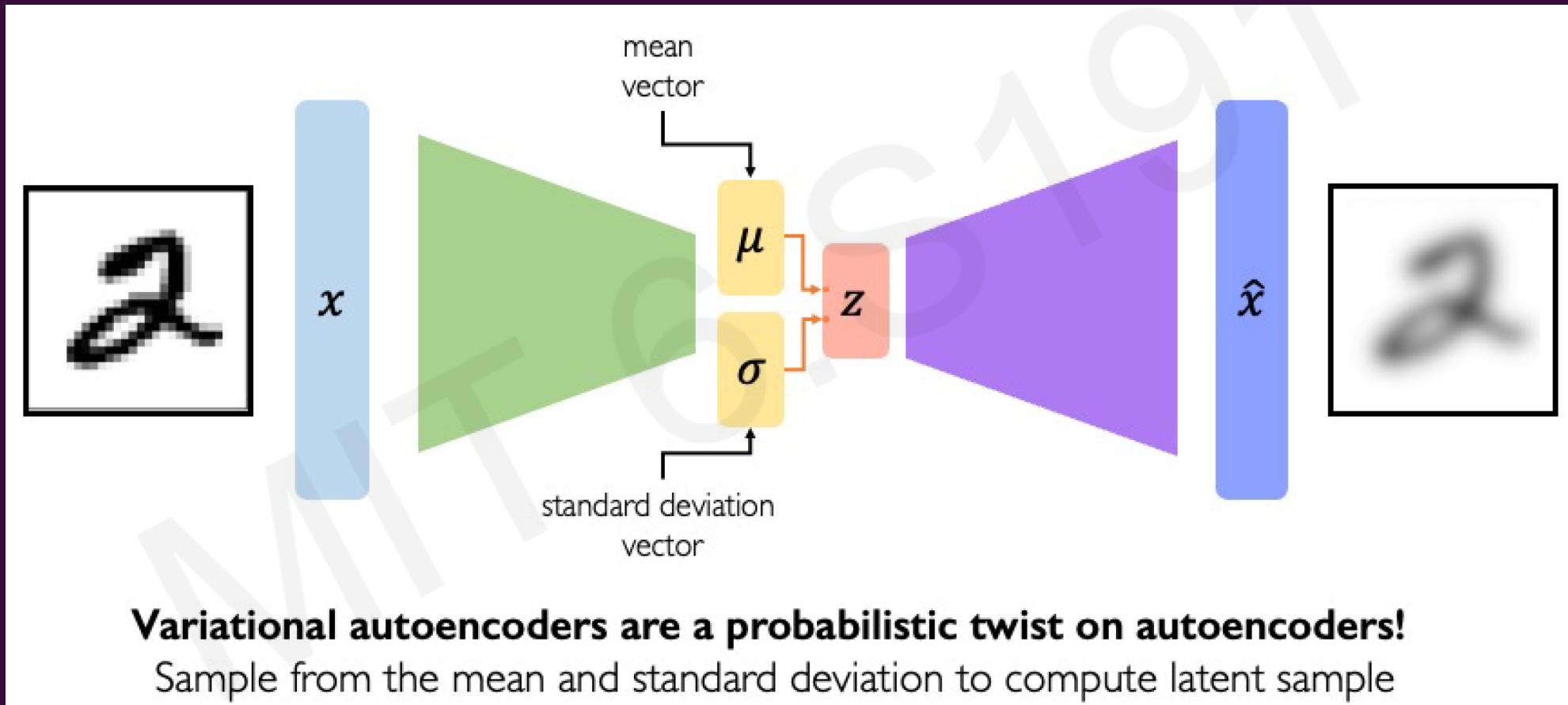


Co-funded by
the European Union



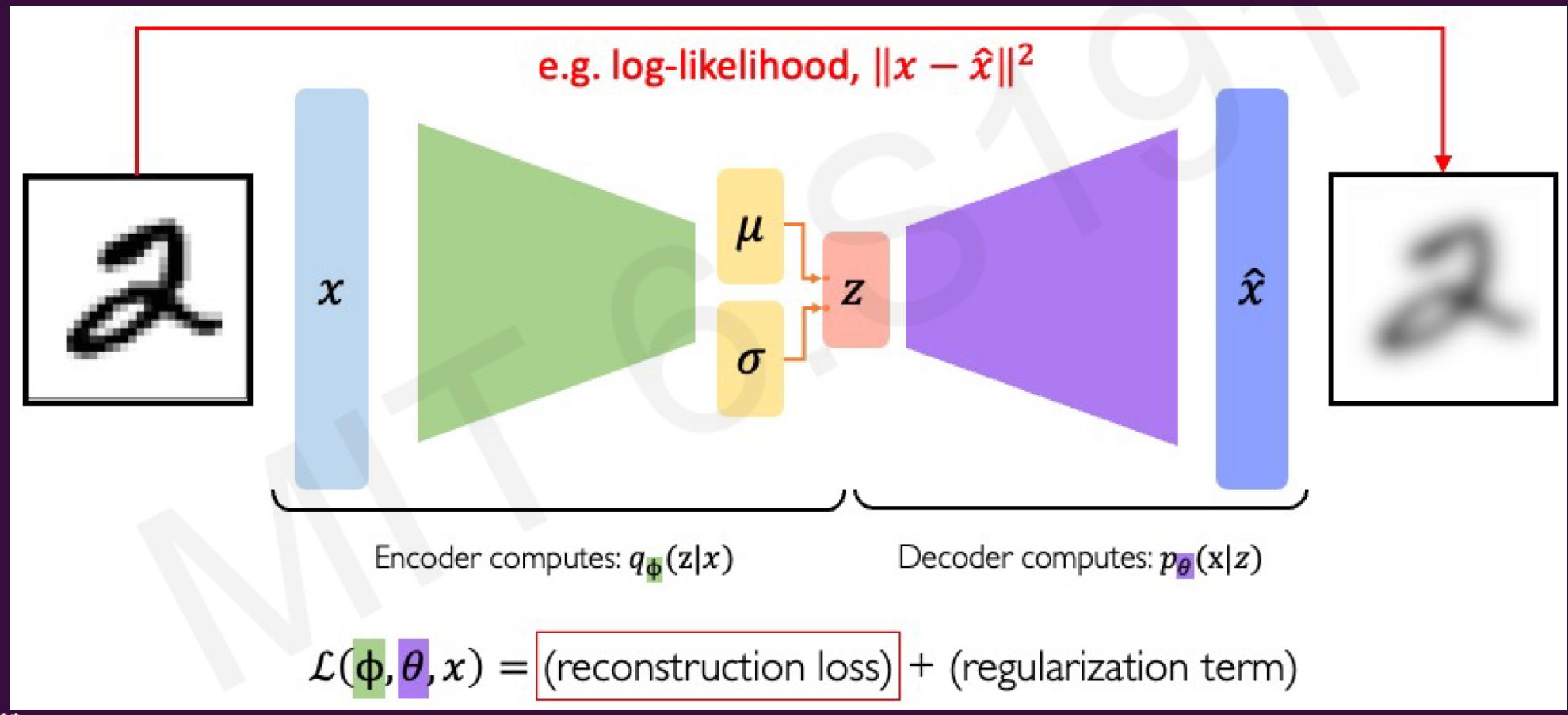
VAEs: Key Difference with Traditional Autoencoder

VAEs extend traditional autoencoders by introducing probabilistic modeling, enabling more structured and meaningful latent representations.



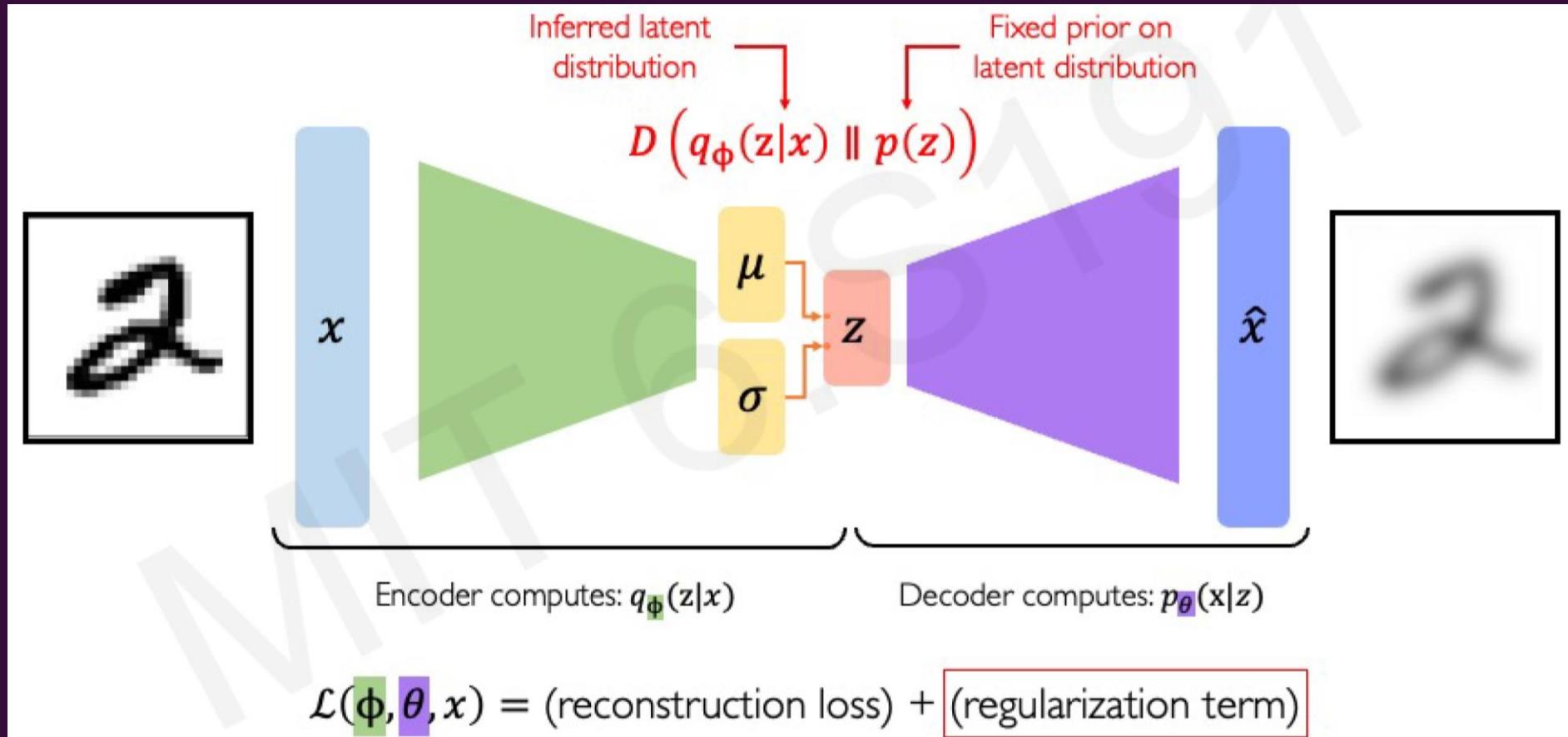
VAE Optimization

They optimize a balance between reconstruction accuracy and latent space regularization, making them useful for data generation.



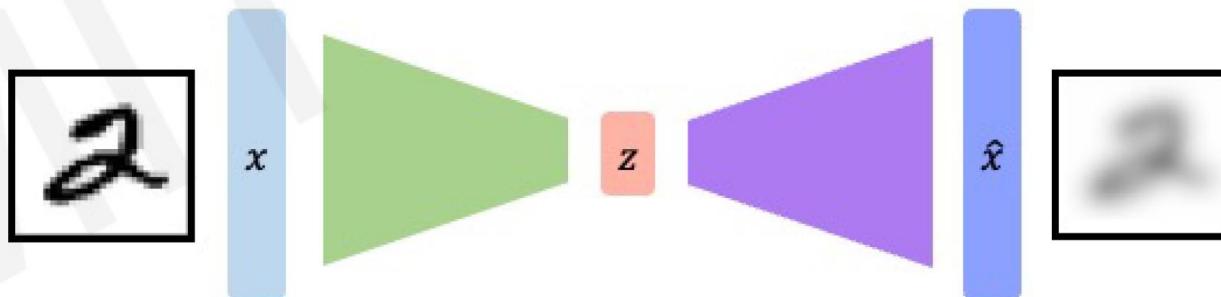
VAE Optimization

- Variational Autoencoders (VAEs) learn compressed representations
- These models enable applications like image synthesis and data augmentation of data



VAE Summary

1. Compress representation of world to something we can use to learn
2. Reconstruction allows for unsupervised learning (no labels!)
3. Reparameterization trick to train end-to-end
4. Interpret hidden latent variables using perturbation
5. Generating new examples



Generative Adversarial Network (GAN)



Co-funded by
the European Union

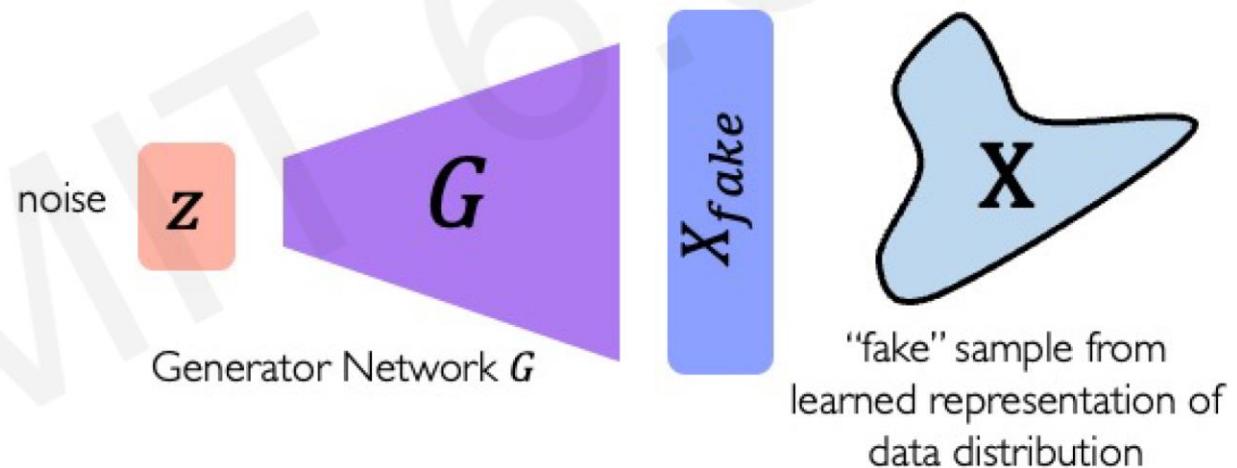


What if we Just Want to Sample?

Idea: don't explicitly model density, and instead just sample to generate new instances.

Problem: want to sample from complex distribution – can't do this directly!

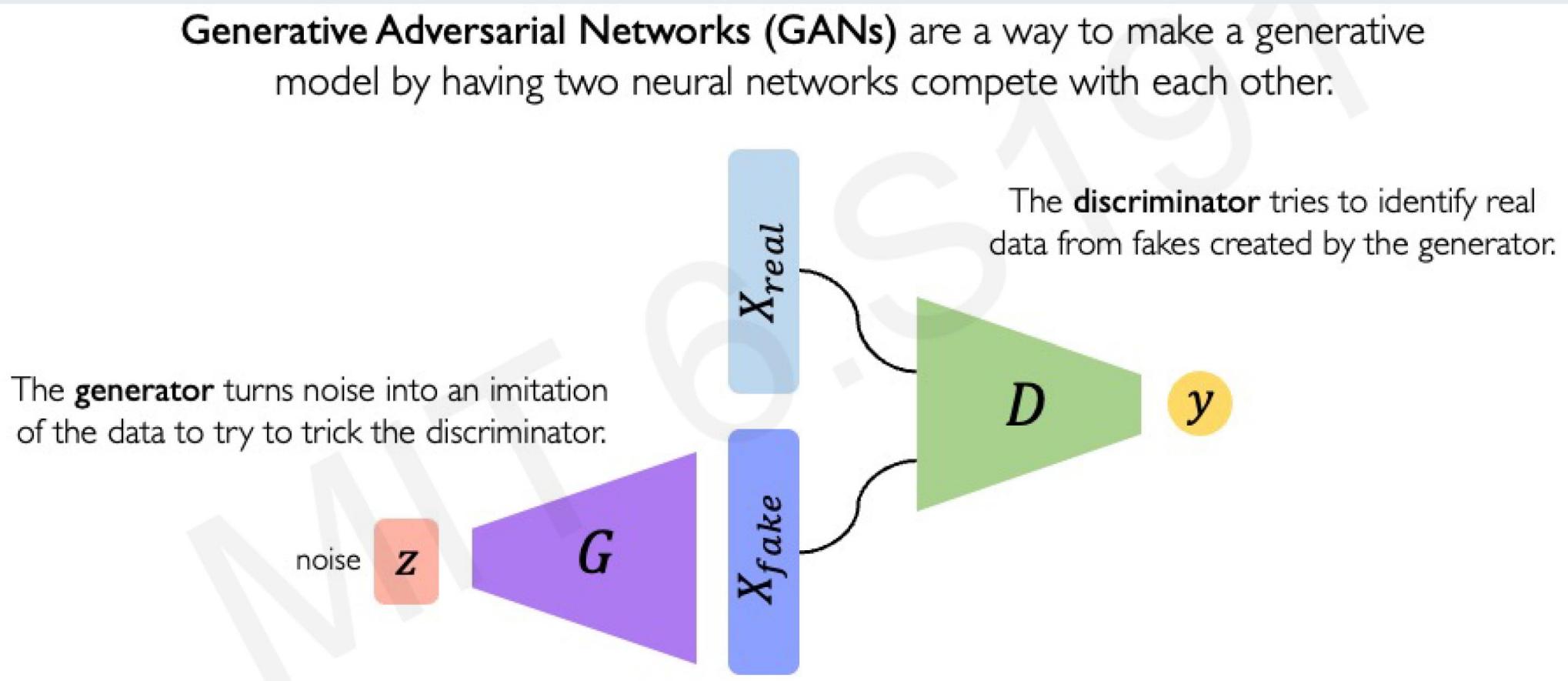
Solution: sample from something simple (e.g., noise), learn a transformation to the data distribution.



Generative Adversarial Network (GAN)

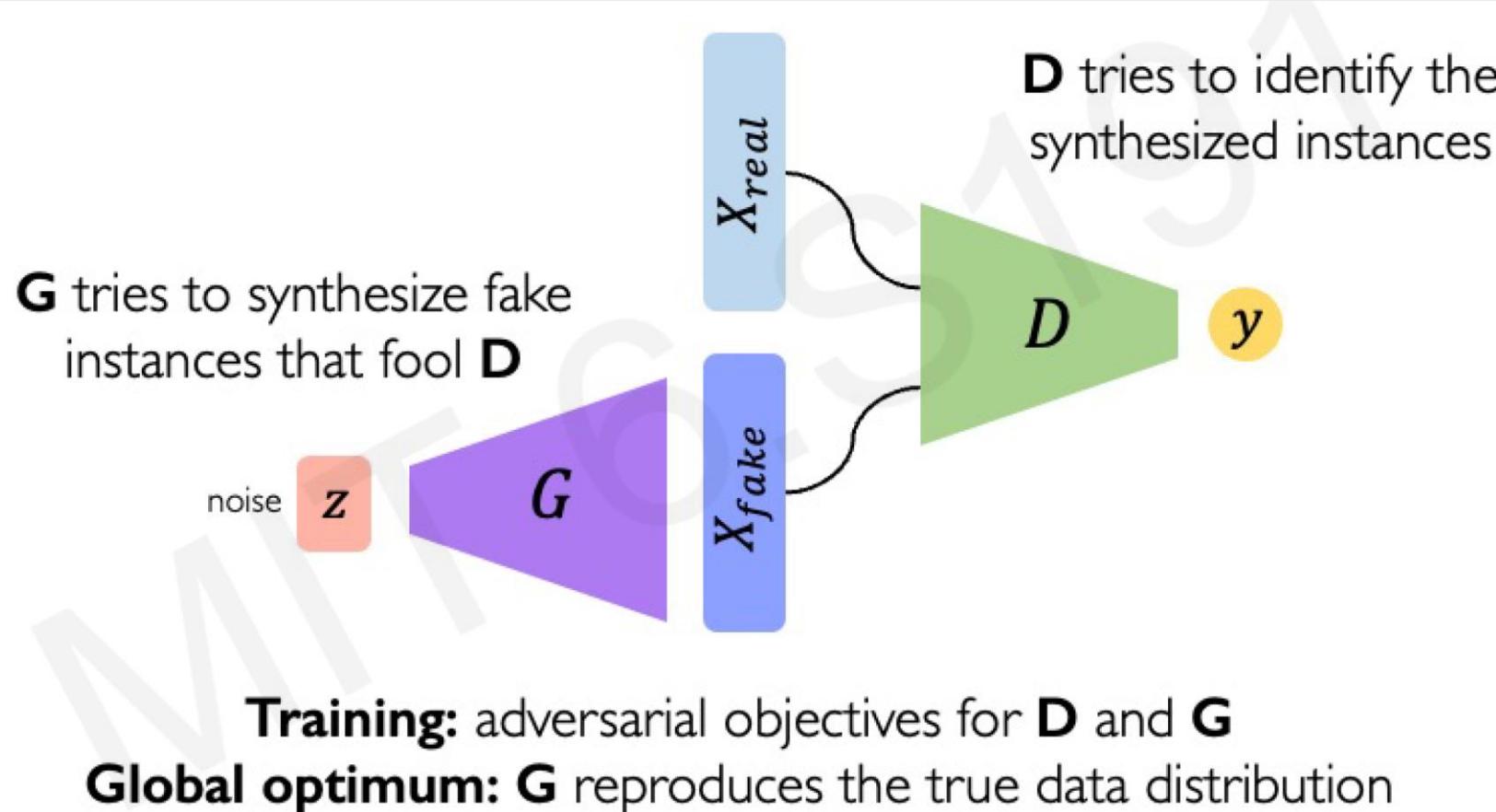
GANs use two competing networks: a generator that creates data and a discriminator that evaluates it

Generative Adversarial Networks (GANs) are a way to make a generative model by having two neural networks compete with each other.



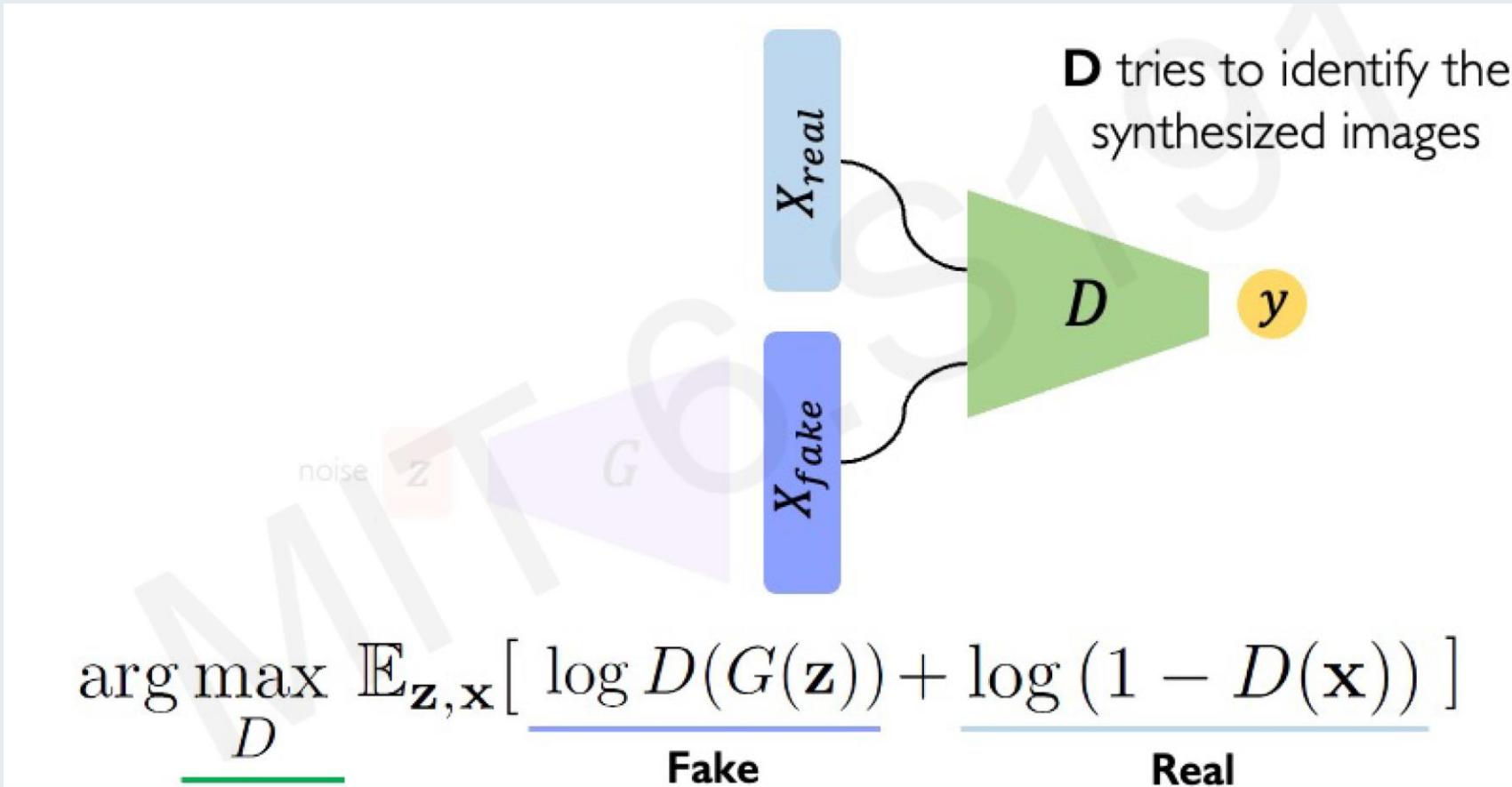
Training GANs

The generator creates fake data, while the discriminator learns to distinguish between real and fake samples, improving both over time.

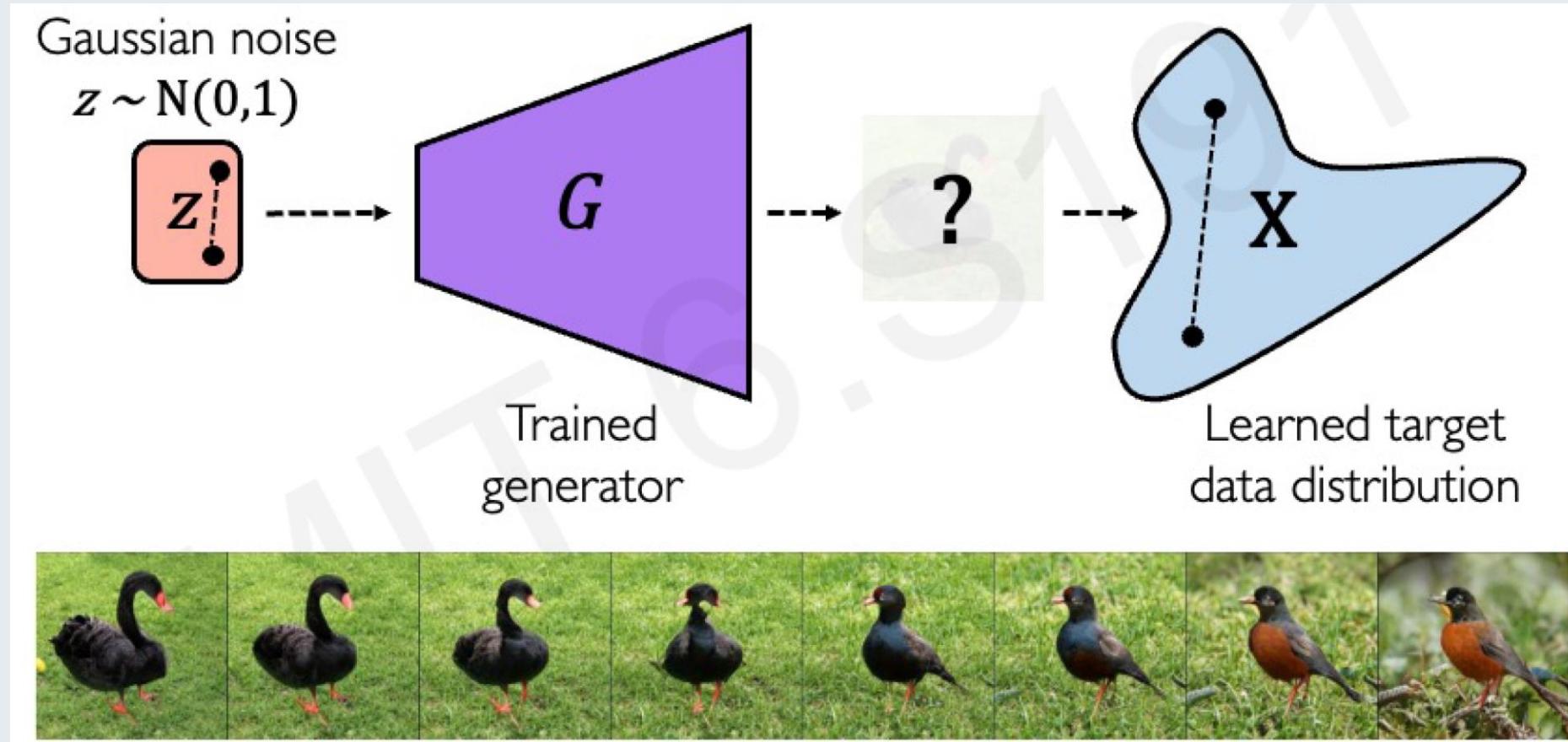


Training GANs: Loss Function

GANs use a min-max game approach where the generator tries to fool the discriminator, and the discriminator aims to correctly classify samples.



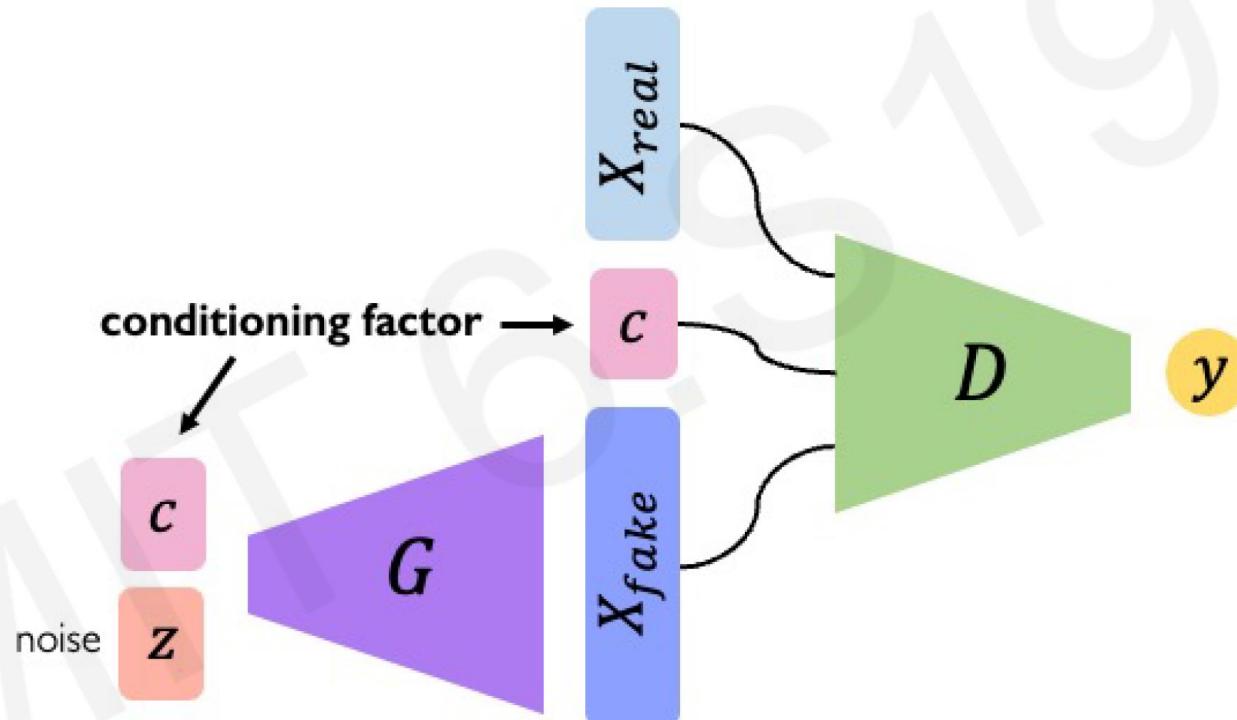
GANs are Distribution Transformers



Co-funded by
the European Union

Conditional GANs

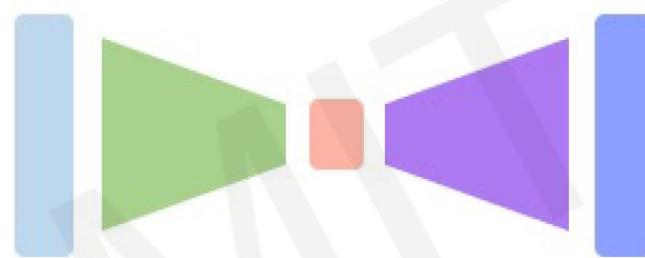
What if we want to control the nature of the output, by **conditioning** on a label?



Deep Generative Modeling Summary

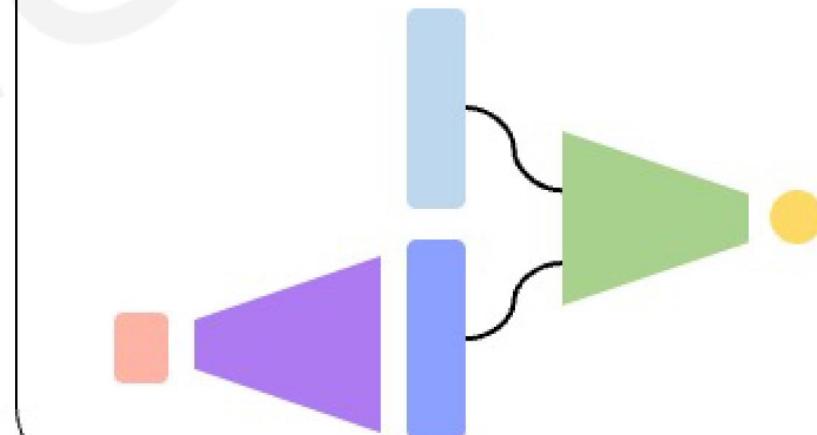
Autoencoders and Variational Autoencoders (VAEs)

Learn lower-dimensional **latent space** and **sample** to generate input reconstructions



Generative Adversarial Networks (GANs)

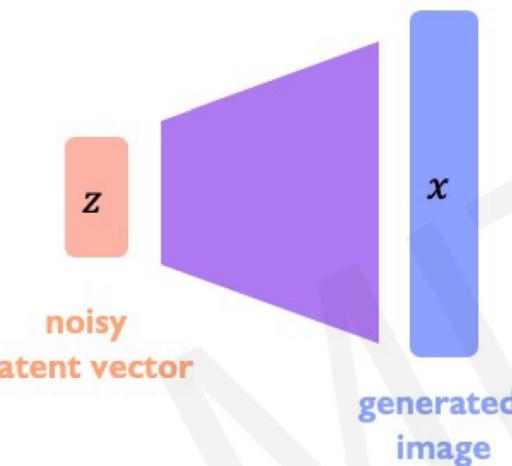
Competing **generator** and **discriminator** networks



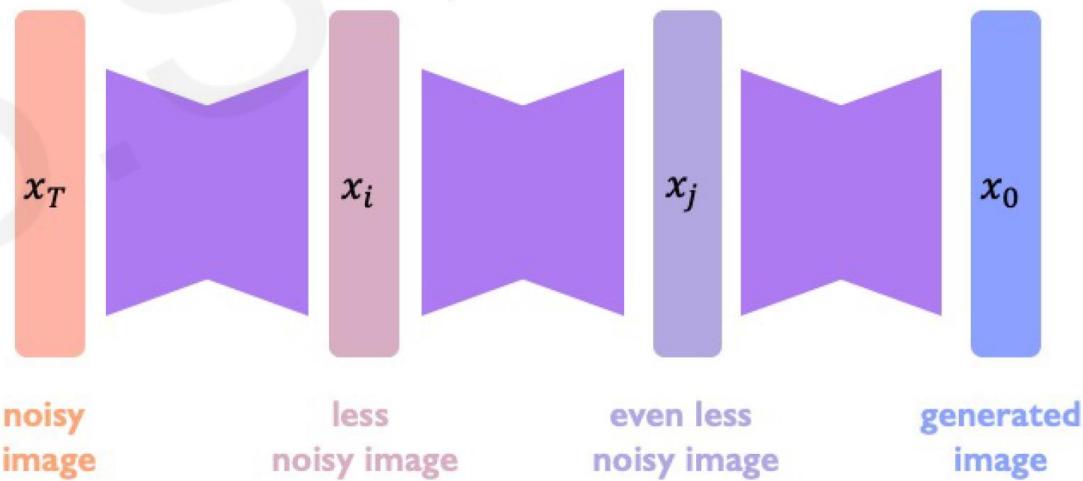
Diffusion Models

Diffusion models iteratively transform noise into structured data, often used for image synthesis and denoising tasks.

VAEs/GANs: Generating images in one-shot directly from low-dimensional latent variables

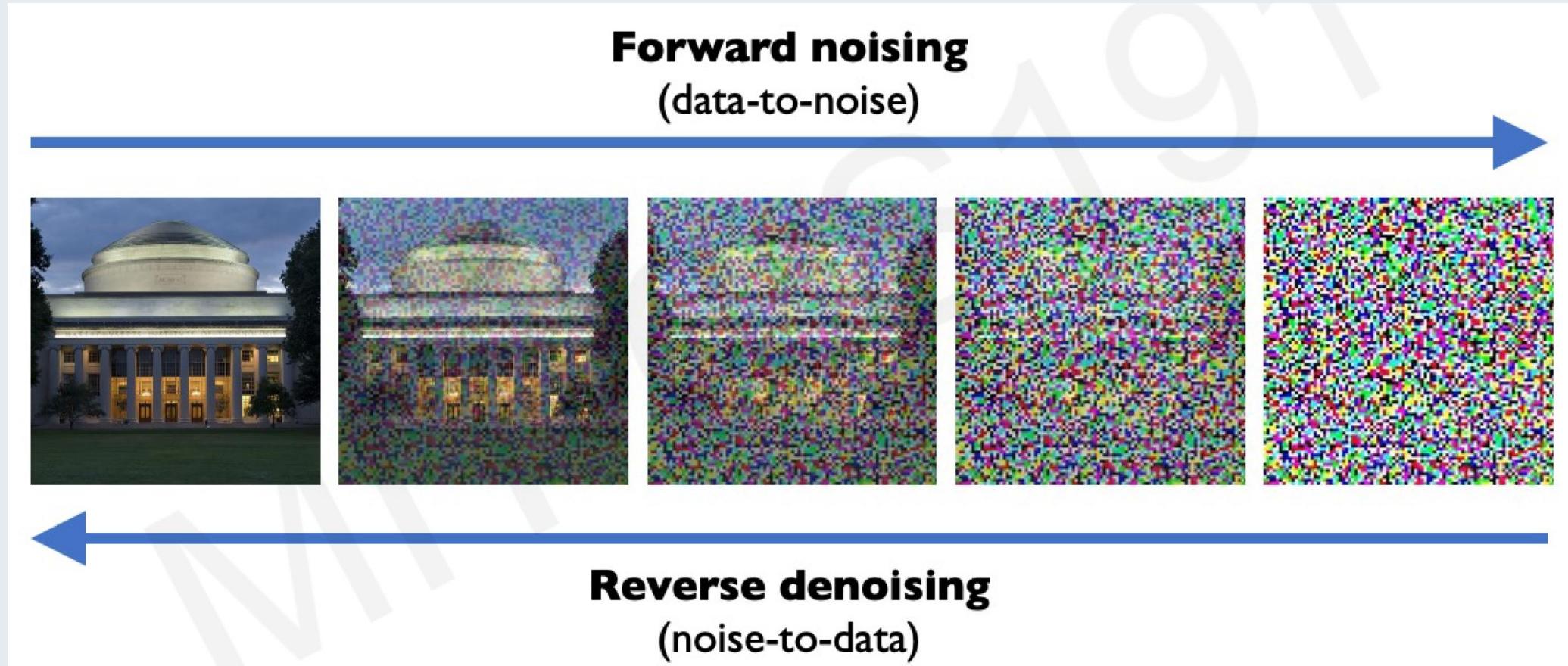


Diffusion: Generating images iteratively by repeatedly refining and removing noise

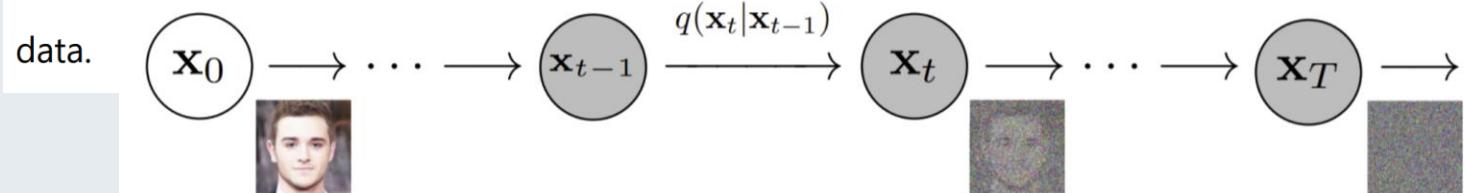


The Diffusion Process

Diffusion models gradually add noise to data and learn to reverse this process

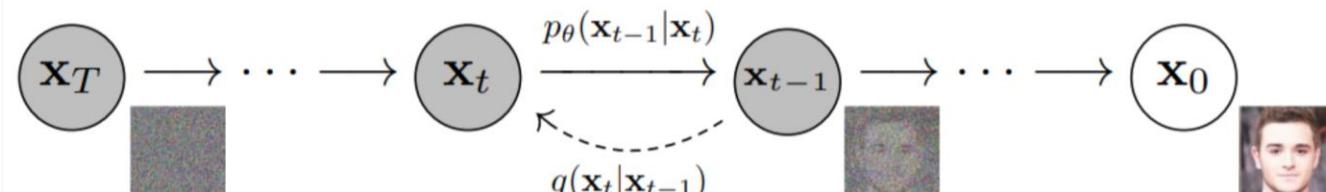


More specifically, a Diffusion Model is a latent variable model which maps to the latent space using a fixed Markov chain. This chain gradually adds noise to the data in order to obtain the approximate posterior $q(\mathbf{x}_1:T|\mathbf{x}_0)$, where $\mathbf{x}_1, \dots, \mathbf{x}_T$ are the latent variables with the same dimensionality as \mathbf{x}_0 . In the figure below, we see such a Markov chain manifested for image data.



(Modified from [source](#))

Ultimately, the image is asymptotically transformed to pure Gaussian noise. The **goal** of training a diffusion model is to learn the **reverse** process - i.e. training $p_\theta(x_{t-1}|x_t)$. By traversing backwards along this chain, we can generate new data.



(Modified from [source](#))

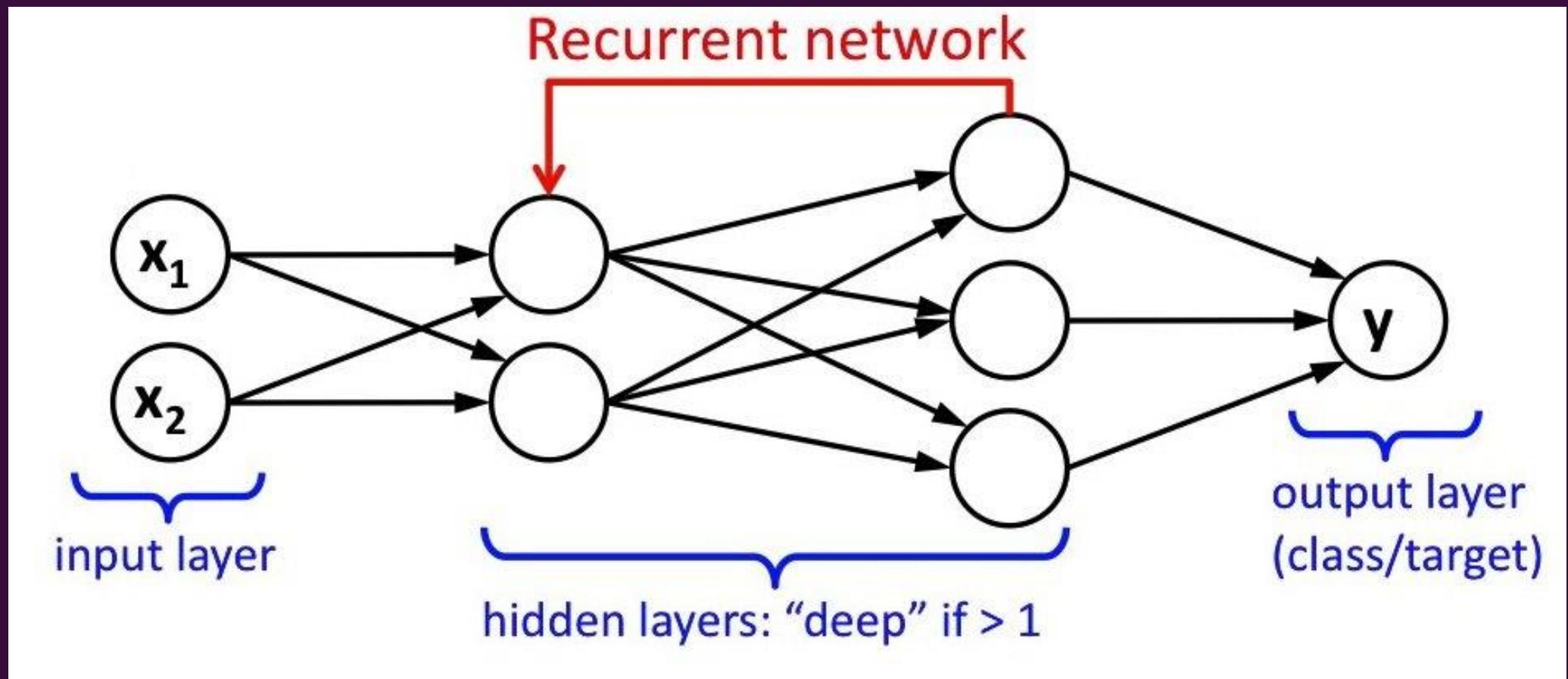


Co-funded by
the European Union

<https://www.assemblyai.com/blog/diffusion-models-for-machine-learning-introduction/>

Recurrent Neural Networks

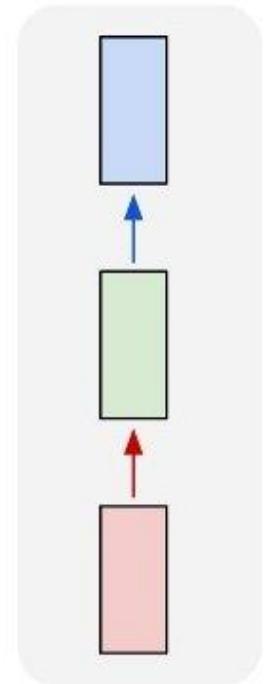
RNNs are designed for sequential data, such as time series and natural language, by maintaining memory across steps.



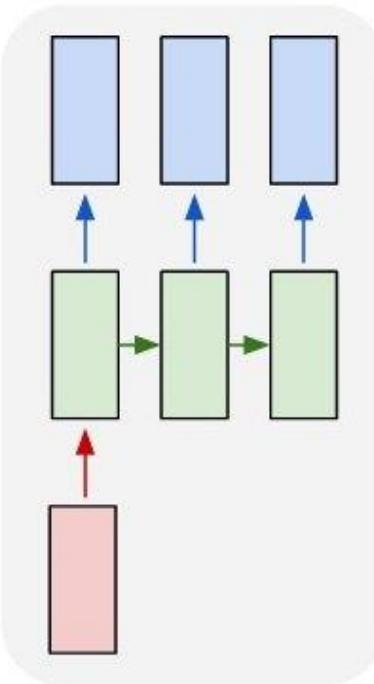
Recurrent Neural Networks

Different RNN architectures suit different types of sequence processing tasks

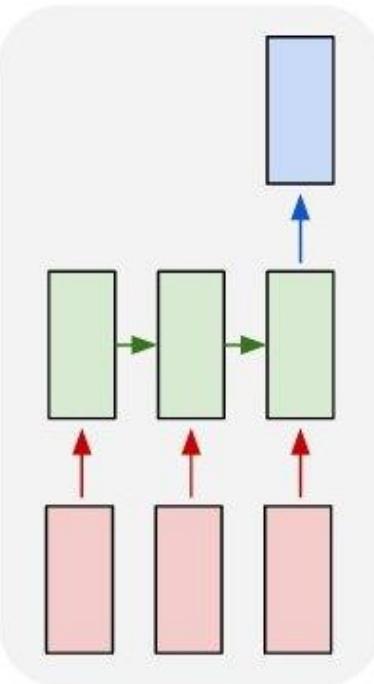
one to one



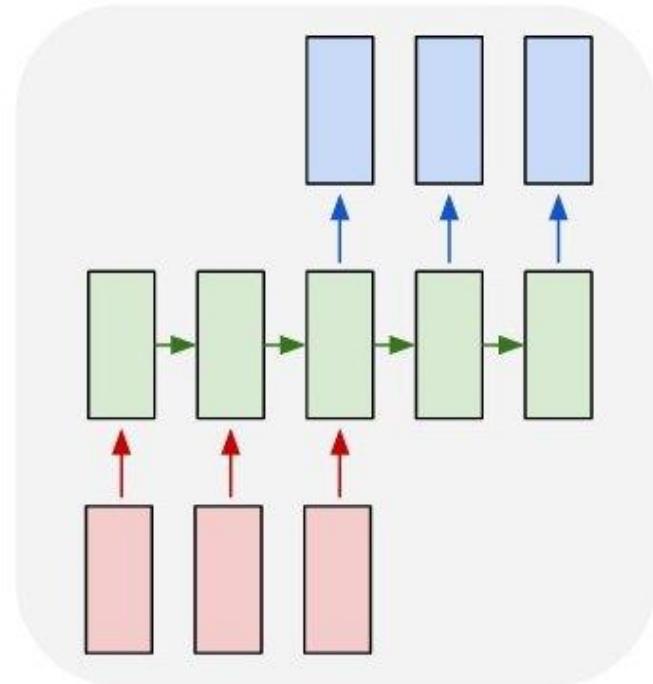
one to many



many to one



many to many



many to many

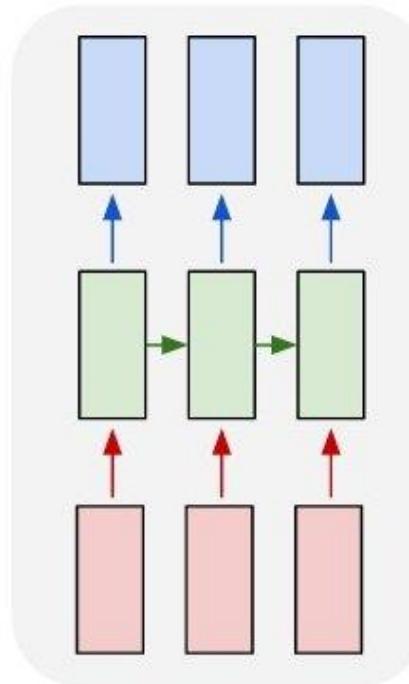


Image in
Label out

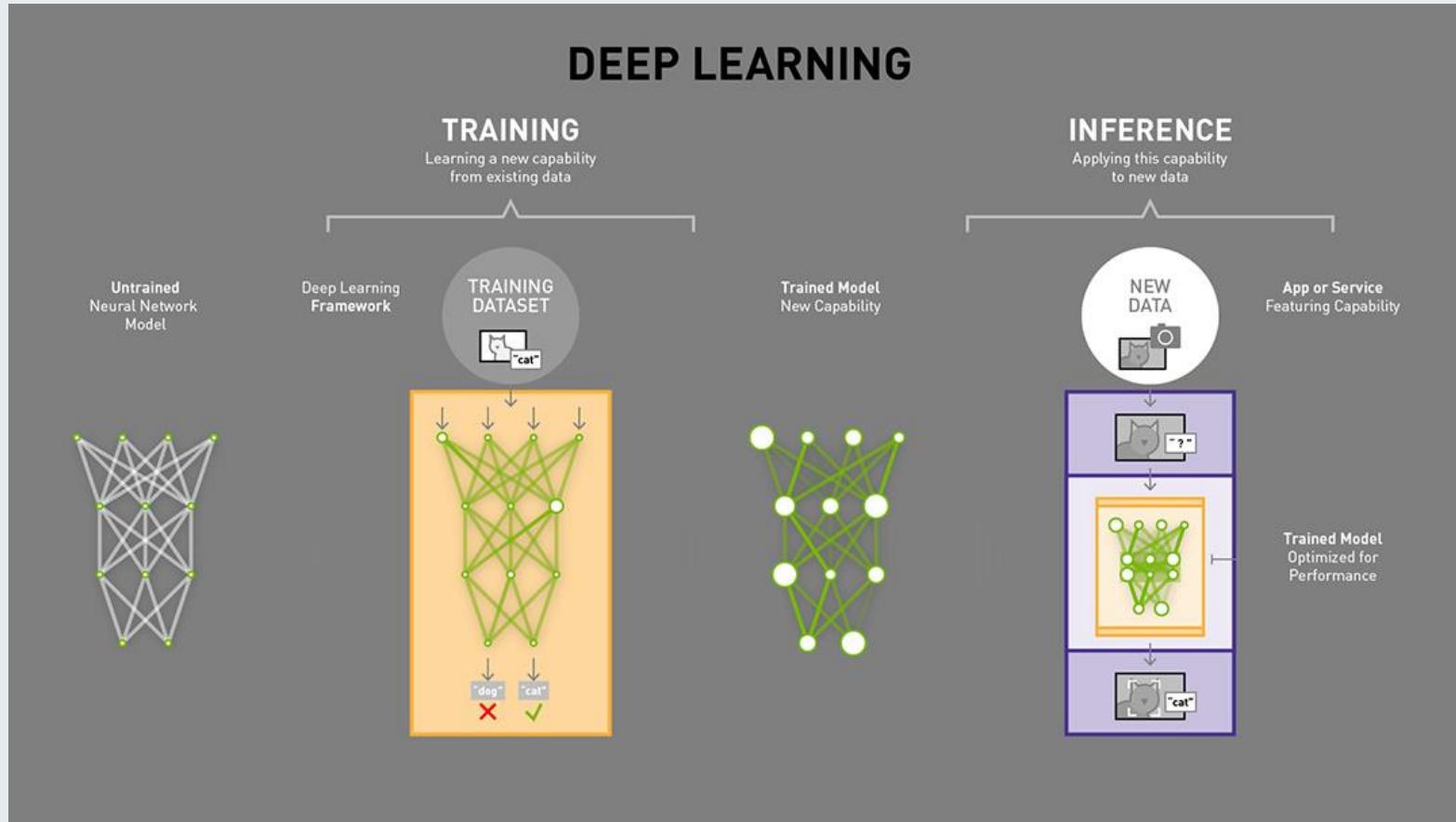
Image in
Words out

Words in
Sentiment out

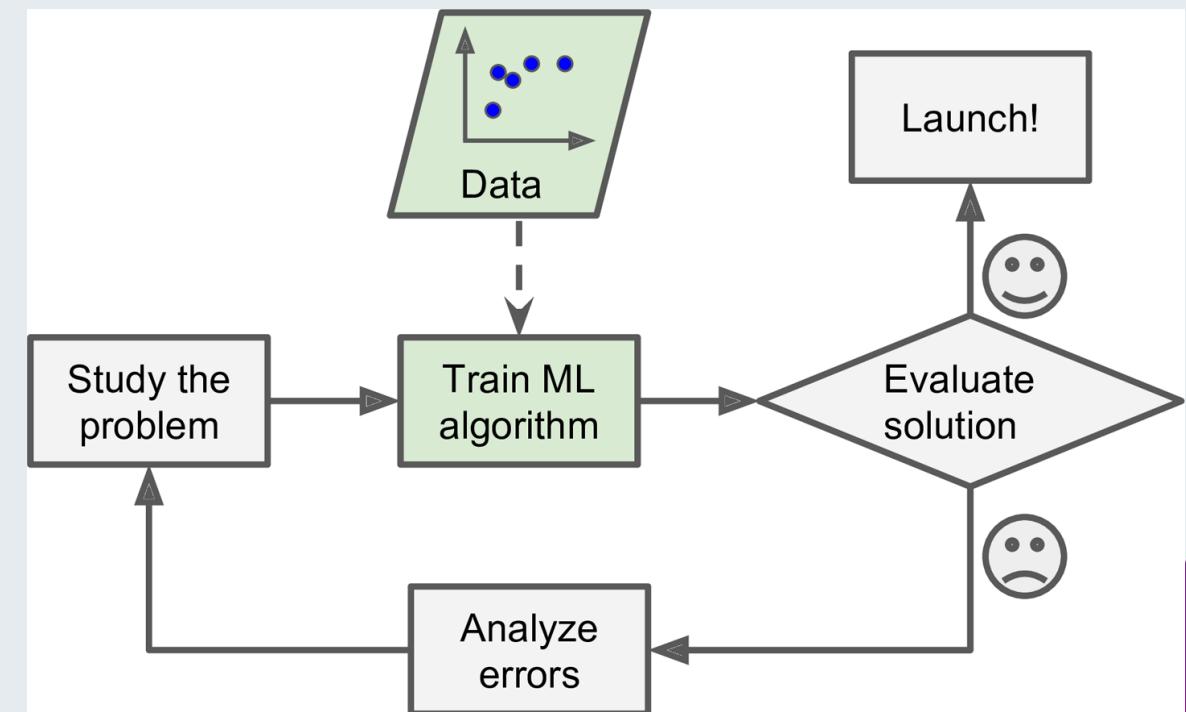
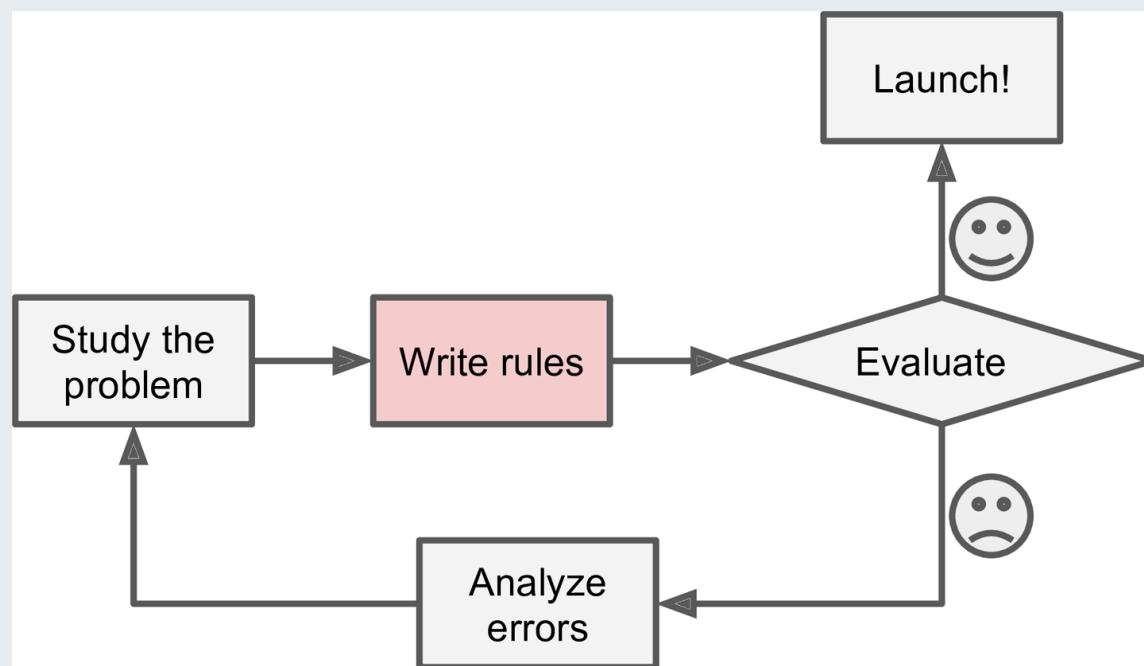
English in
Portuguese out

Video In
Labels out

<https://blogs.nvidia.com/blog/2016/08/22/difference-deep-learning-training-inference-ai/>



Traditional vs ML problem solving



Neural Networks Limitations

- Very **data hungry** (eg. often millions of examples)
- **Computationally intensive** to train and deploy (tractably requires GPUs)
- Easily fooled by **adversarial examples**
- Can be subject to **algorithmic bias**
- Poor at **representing uncertainty** (how do you know what the model knows?)
- Uninterpretable **black boxes**, difficult to trust
- Often require **expert knowledge** to design, fine tune architectures
- Difficult to **encode structure** and prior knowledge during learning
- **Extrapolation**: struggle to go beyond the data