"When you're **fundraising**, it's **AI**
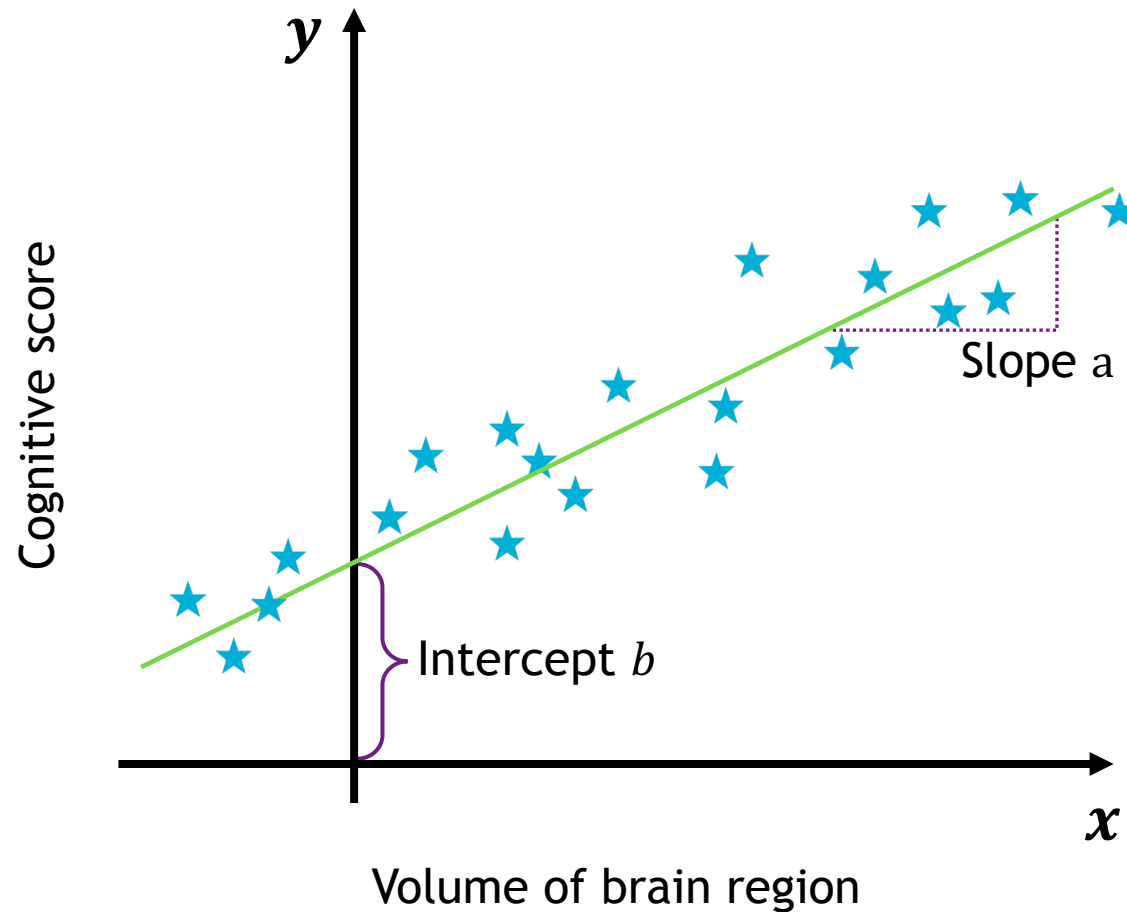
When you're **hiring**, it's **ML**

When you're **implementing**, it's **linear regression**"

— Baron Schwartz

# The basics

## What is linear regression?
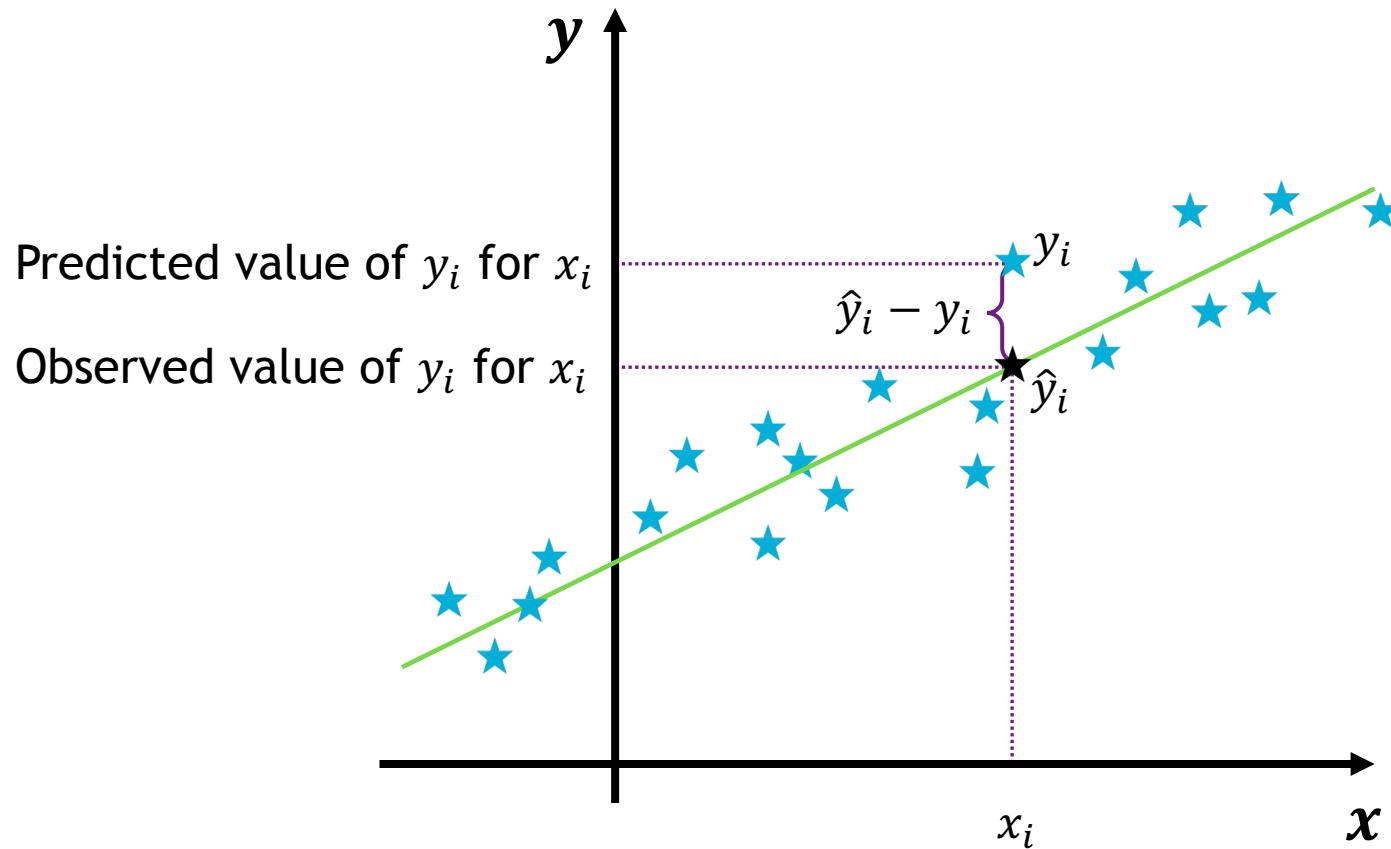


$$y = f(x) = ax + b$$

## Cost function



Predicted value of $y_i$ for $x_i$

Observed value of $y_i$ for $x_i$

$\hat{y}_i - y_i$

$y_i$

$\hat{y}_i$

$x_i$

Mean squared error:

$$J = \frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2$$

## Learning



Predicted value of $y_i$ for $x_i$

Observed value of $y_i$ for $x_i$

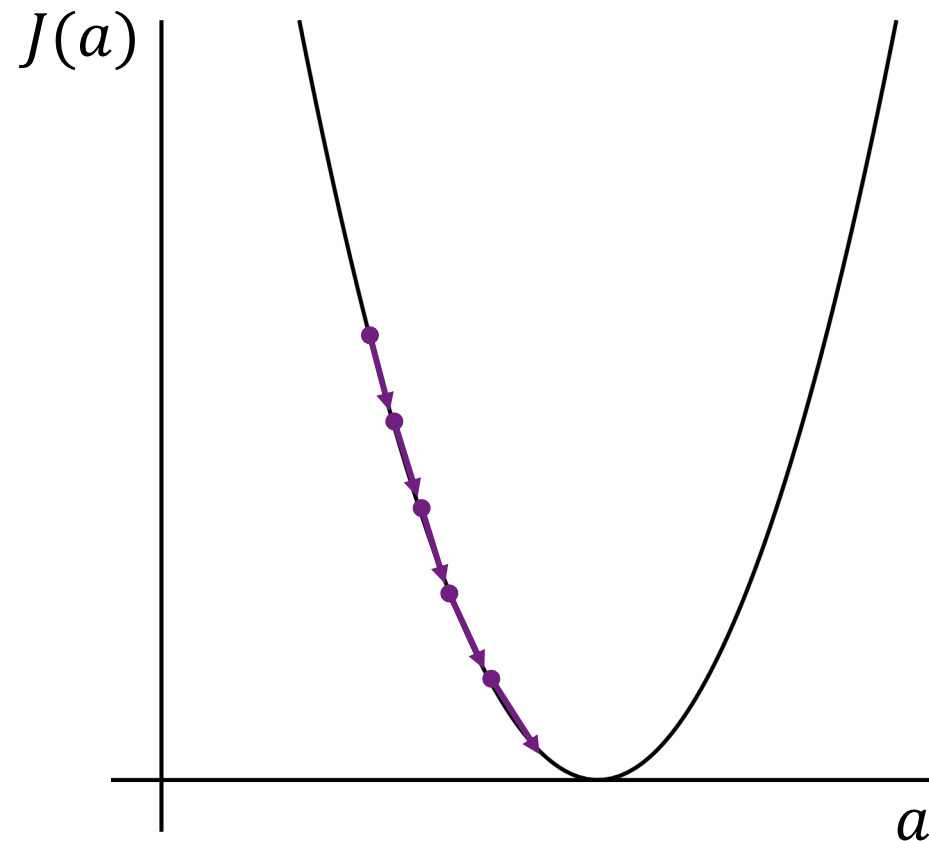$\hat{y}_i - y_i$

$y_i$

$\hat{y}_i$

$x_i$

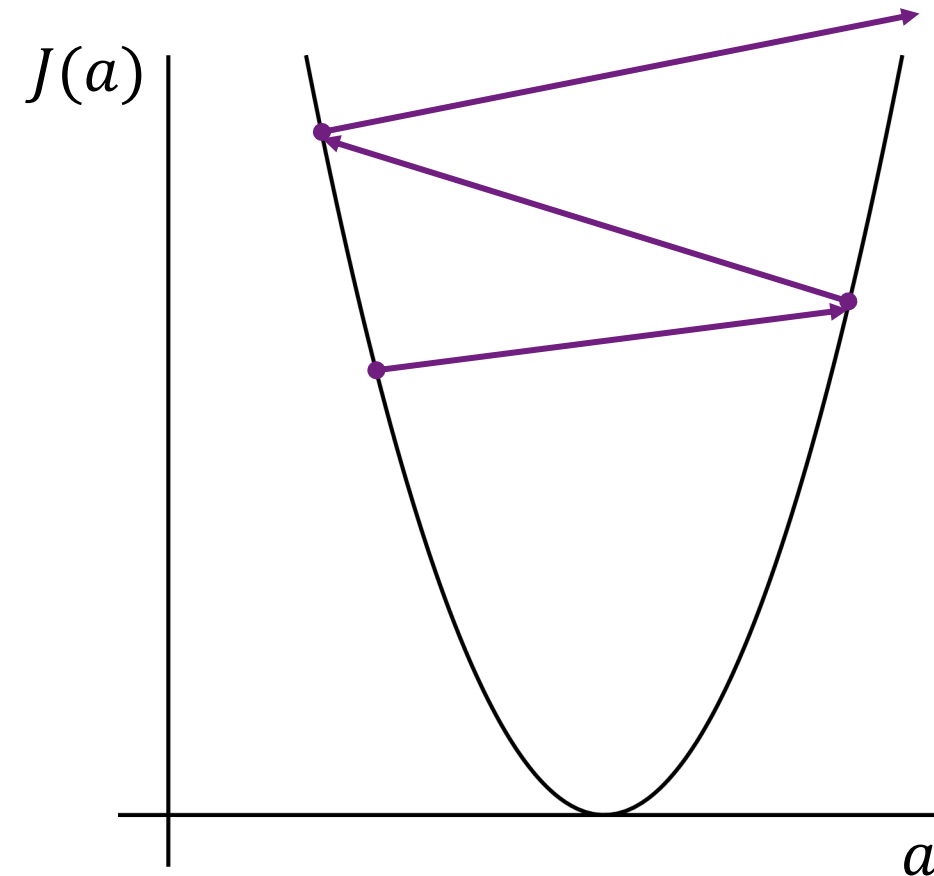Find the model with the minimal error:

$$\hat{f} = arg \min_{f \in \mathcal{F}} J(f)$$

# Linear regression

## Gradient descent

Small learning rate

Large learning rate

## Gradient descent

**Cost function:**

$$J = \frac{1}{n}\sum_{i=1}^{n}(\hat{y}_i - y_i)^2 = \frac{1}{n}\sum_{i=1}^{n}(ax_i + b - y_i)^2$$

**Partial derivates:**

$$\frac{\partial J}{\partial a} = \frac{2}{n}\sum_{i=1}^{n}(ax_i + b - y_i)\cdot x_i = \frac{2}{n}\sum_{i=1}^{n}(\hat{y}_i - y_i)\cdot x_i$$

$$\frac{\partial J}{\partial b} = \frac{2}{n}\sum_{i=1}^{n}(ax_i + b - y_i) = \frac{2}{n}\sum_{i=1}^{n}(\hat{y}_i - y_i)$$
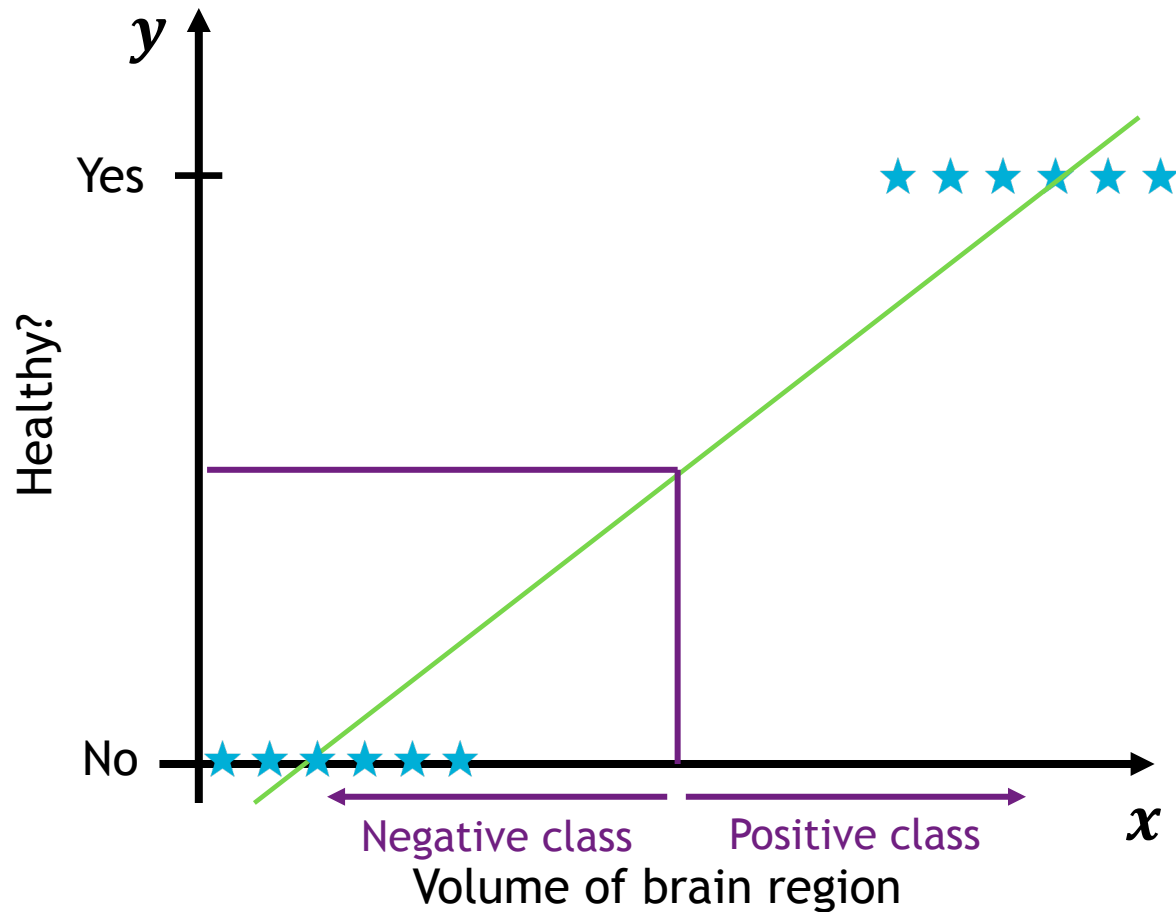
**Update of $a$ and $b$:**

$$a \leftarrow a - \eta\frac{\partial J}{\partial a} = a - \eta\frac{2}{n}\sum_{i=1}^{n}(\hat{y}_i - y_i)$$

$$b \leftarrow b - \eta\frac{\partial J}{\partial b} = b - \eta\frac{2}{n}\sum_{i=1}^{n}(\hat{y}_i - y_i)\cdot x_i$$
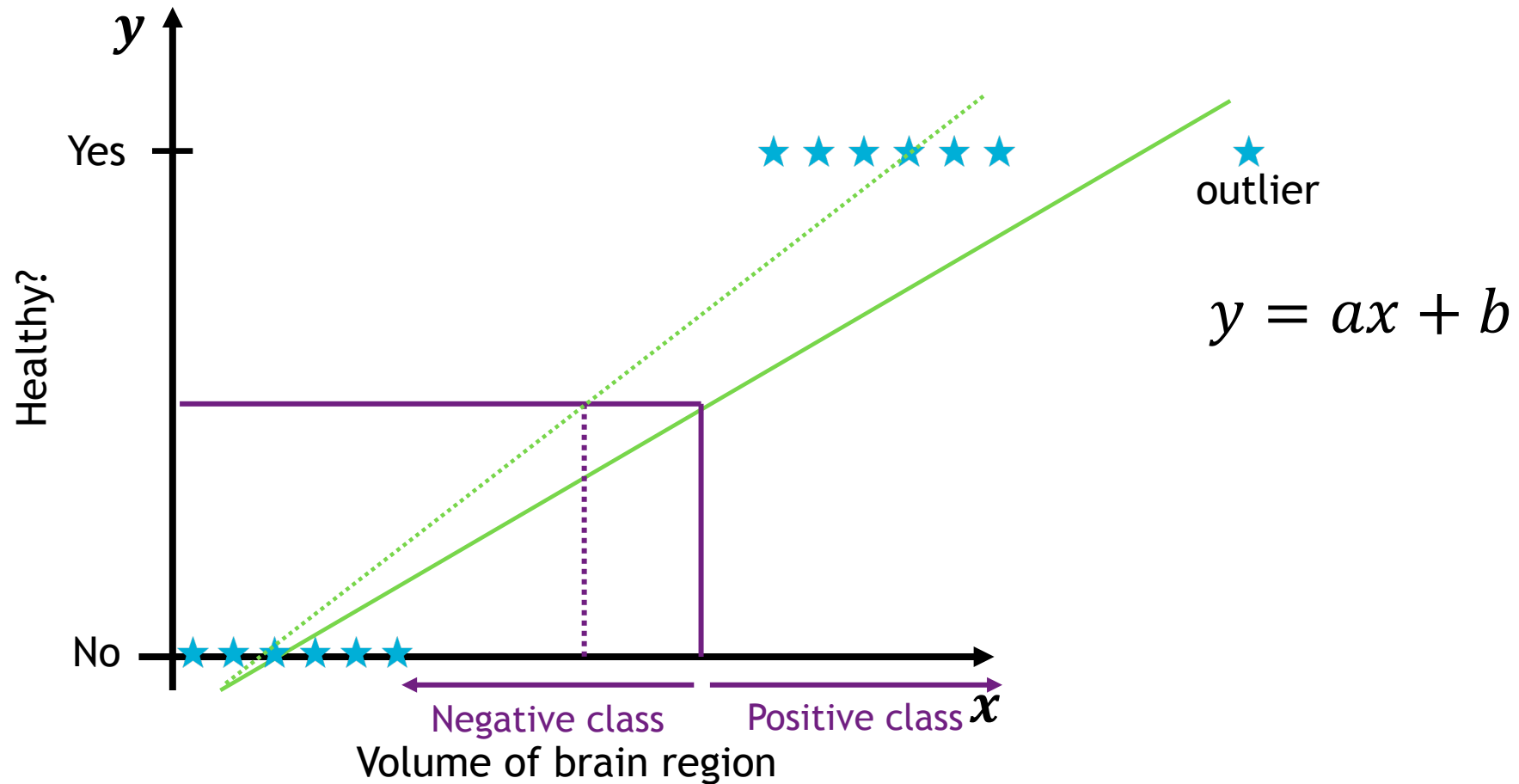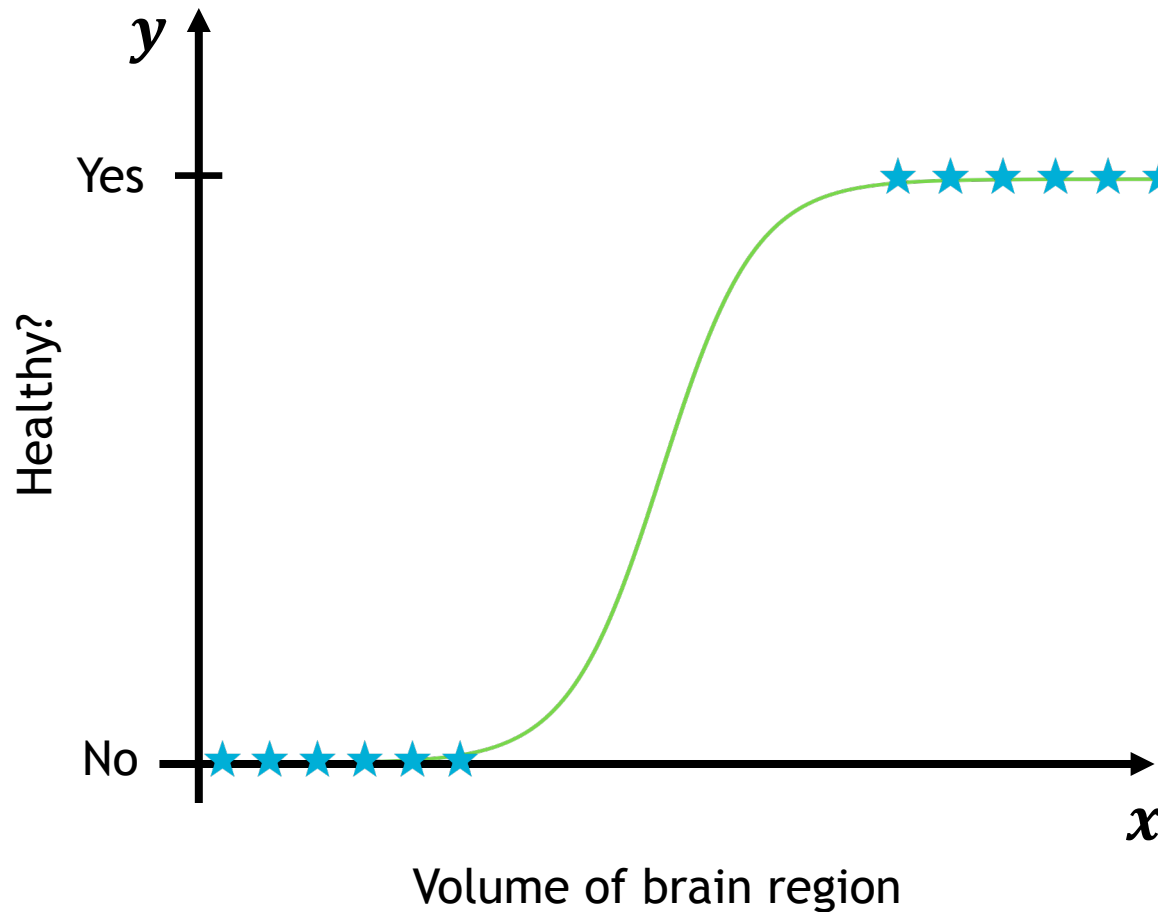
## Linear function



$$y = ax + b$$

## Linear function

## Sigmoid function



$$y = ax + b$$

$$f = \frac{1}{1 + e^{-y}} = \frac{1}{1 + e^{-(ax+b)}}$$

## Multiple input variables

$$z = w_0 + w_1 x_1 + w_2 x_2 + \cdots + w_m x_m = w_0 + \sum_{j=1}^{m} x_j w_j$$

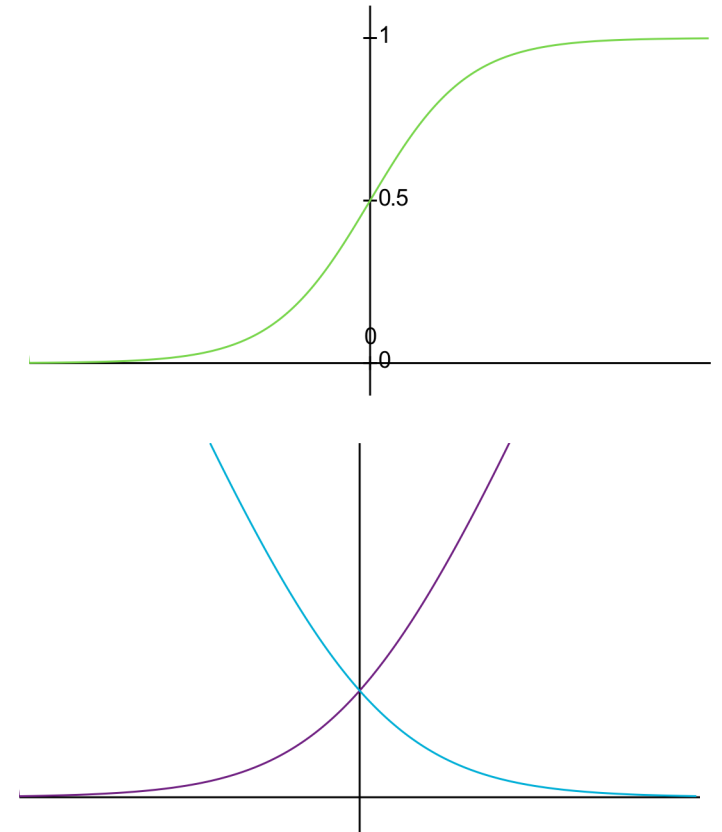$$f = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-\left(w_0 + \sum_{j=1}^{m} x_j w_j\right)}}$$

# Logistic regression

## Loss function

- **Low if prediction is correct**

- **High if prediction is wrong**

$$\begin{cases} \text{if } y = 1, l(f(\boldsymbol{x}), y) \text{ should be low when } f(\boldsymbol{x}) \text{ is high} \\ \text{if } y = 0, l(f(\boldsymbol{x}), y) \text{ should be low when } f(\boldsymbol{x}) \text{ is low} \end{cases}$$

$$l(f(\boldsymbol{x}), y) = \begin{cases} -\log(f(\boldsymbol{x})), & \text{if } y = 1 \\ -\log(1 - f(\boldsymbol{x})), & \text{if } y = 0 \end{cases}$$

$$l(f(\boldsymbol{x}), y) = -y \log(f(\boldsymbol{x})) - (1 - y) \log(1 - f(\boldsymbol{x}))$$

## Cost function

$$J(f) = \frac{1}{n} \sum_{i=1}^{n} l(f(\boldsymbol{x}_i), y_i)$$

$$J(f) = -\frac{1}{n} \sum_{i=1}^{n} \left(y_i \log\bigl(f(\boldsymbol{x}_i)\bigr) + (1 - y_i) \log(1 - f(\boldsymbol{x}_i))\right)$$

## Learning

Find the model with the minimal error:

$$\hat{f} = arg \min_{f \in \mathcal{F}} J(f)$$

# Logistic regression

## Gradient descent

**Cost function:**

$$J(f) = -\frac{1}{n}\sum_{i=1}^{n}\left(y_i\log\bigl(f(\boldsymbol{x}_i)\bigr) + (1 - y_i)\log(1 - f(\boldsymbol{x}_i))\right)$$

**Partial derivates:**

$$\frac{\partial J}{\partial w_j} = -\frac{1}{n}\sum_{i=1}^{n}\left(y_{i,j} - f(\boldsymbol{x}_i)\right)\cdot x_{i,j}$$

**Update of $a$ and $b$:**

$$w_j \leftarrow w_j - \eta\frac{\partial J}{\partial w_j}$$

Input variables: $x = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$

Input features

Output: $y$

Model: $f, y = f(x)$

The "artificial intelligence"

Loss: $l(f(x), y)$

Quantifies how much the prediction is far from the true output

Cost function: $J(f) = \dfrac{1}{n} \sum_{i=1}^{n} l(f(x_i), y_i)$

Quantifies how much the prediction is far from the true output across all training examples

Learning: $\hat{f} = arg \min_{f \in \mathcal{F}} J(f)$

Find the model with the minimal error

Gradient descent: $w \leftarrow w - \eta \dfrac{\partial J}{\partial w}$
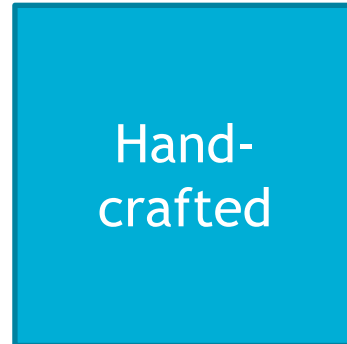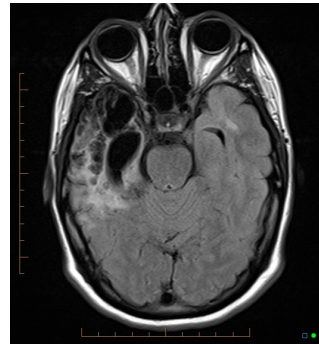
Method to find the model with the minimal error

ML vs DL : learning features

Case courtesy of Dr Chris O'Donnell, Radiopaedia.org, rID: 27433

# Neural networks

## Artificial neuron



$$\hat{y} = h\left(\sum_{j=1}^{m} x_j w_j\right)$$

Output

Linear combination of inputs

Non-linear activation function

Inputs    Weights    Sum    Non-Linearity    Output

MIT Introduction to Deep Learning (introtodeeplearning.com)

## Artificial neuron



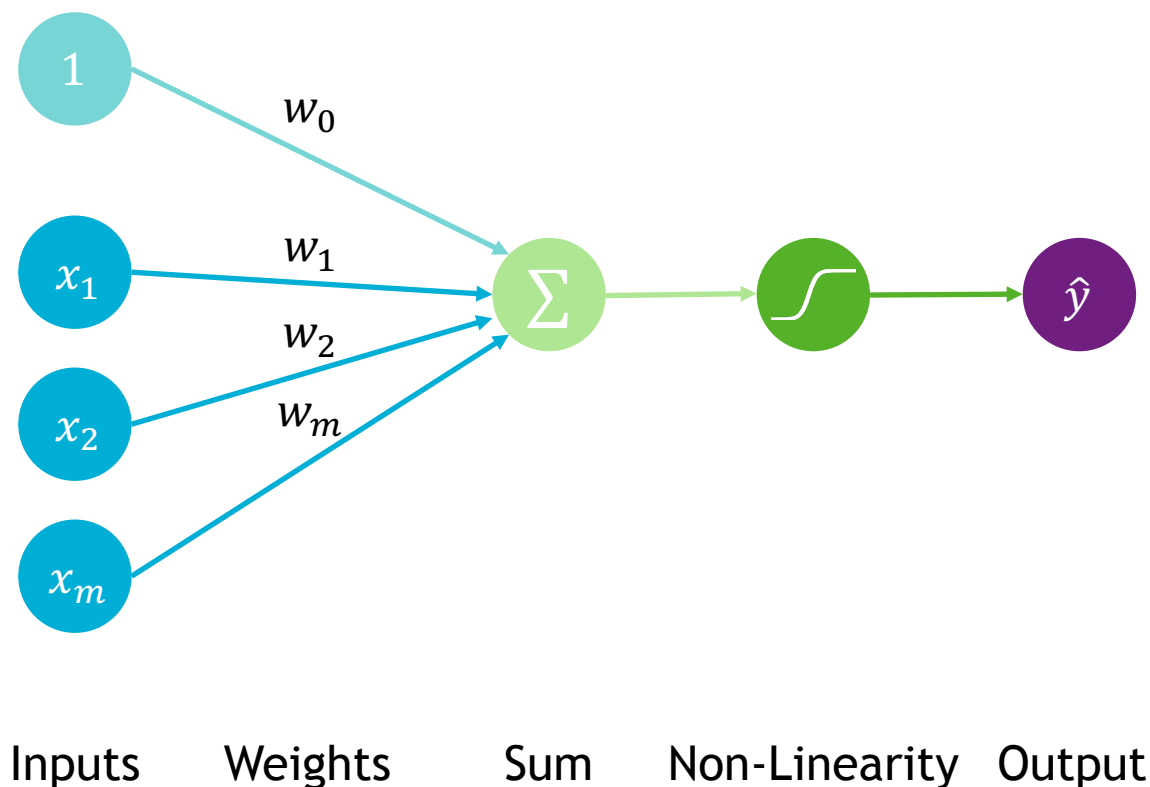Output     Bias     Linear combination of inputs

$$\hat{y} = h\left(w_0 + \sum_{j=1}^{m} x_j w_j\right)$$

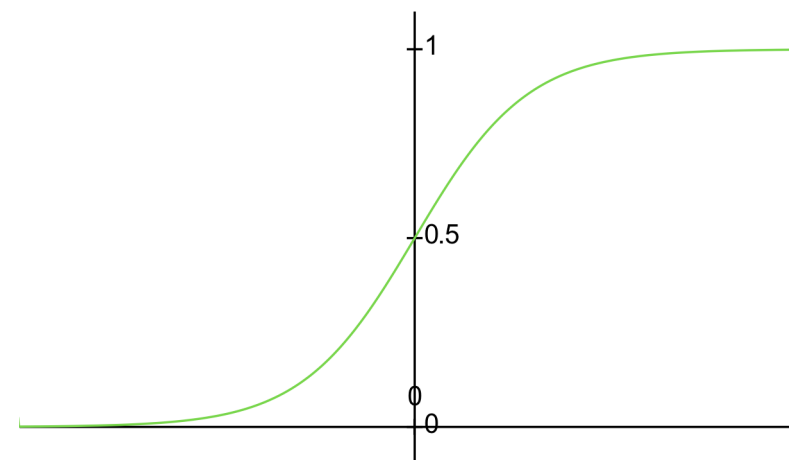Non-linear activation function

$$\hat{y} = h(w_0 + \boldsymbol{X}^T \boldsymbol{W})$$

where $\boldsymbol{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ and $\boldsymbol{W} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$

Inputs    Weights    Sum    Non-Linearity    Output

MIT Introduction to Deep Learning (introtodeeplearning.com)

# Perceptron

## Activation function

$$\hat{y} = h(w_0 + \boldsymbol{X}^T\boldsymbol{W})$$

Sigmoid function: $h(z) = \dfrac{1}{1 + e^{-z}}$



Inputs    Weights    Sum    Non-Linearity    Output

→ Logistic regression

MIT Introduction to Deep Learning (introtodeeplearning.com)

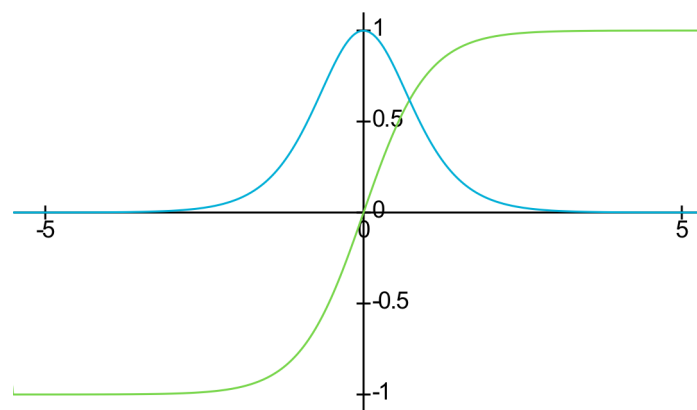## Examples of activation functions

### Sigmoid

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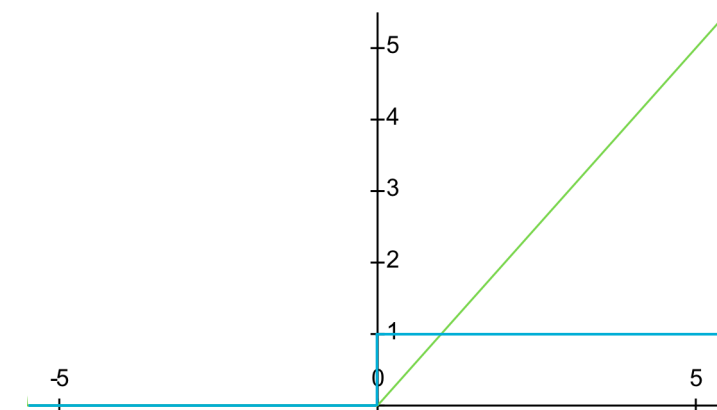
$$h'(z) = h(z)(1 - h(z))$$

### Hyperbolic tangent

$$h(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$
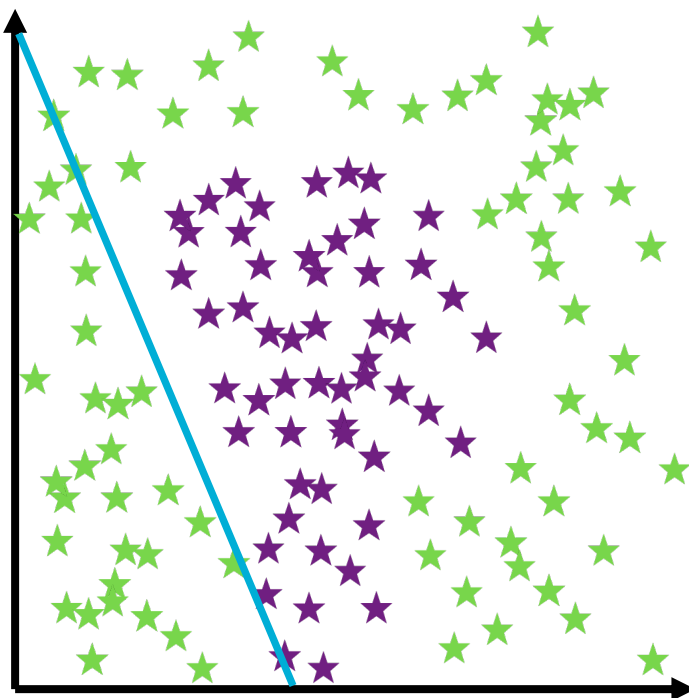
$$h'(z) = 1 - h(z)^2$$
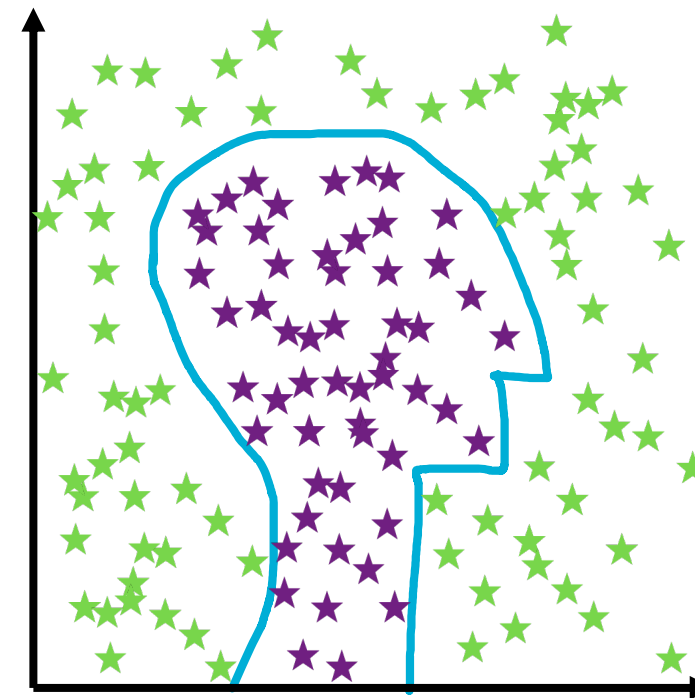
### Rectified linear unit (ReLU)

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

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

MIT Introduction to Deep Learning (introtodeeplearning.com)
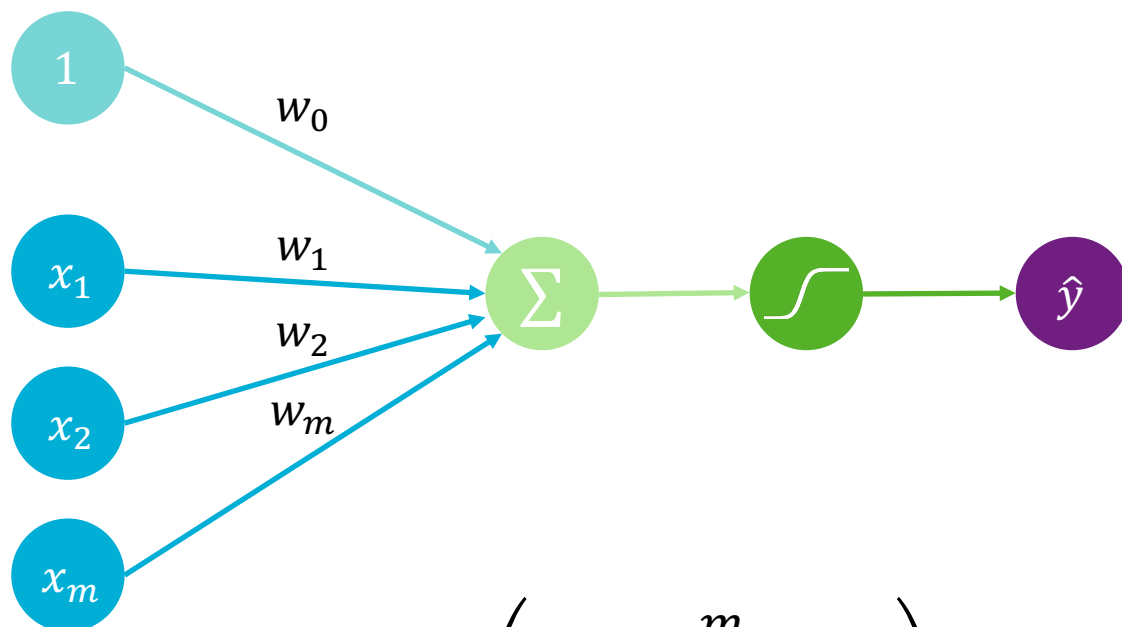
## Importance of non-linear activation functions
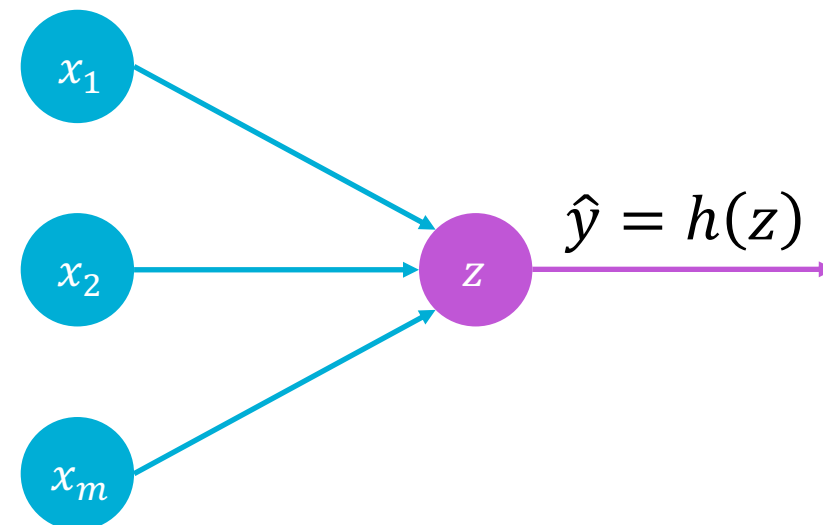


Linear activation functions
→ linear decisions

Non-linear activation functions
→ arbitrarily complex decisions
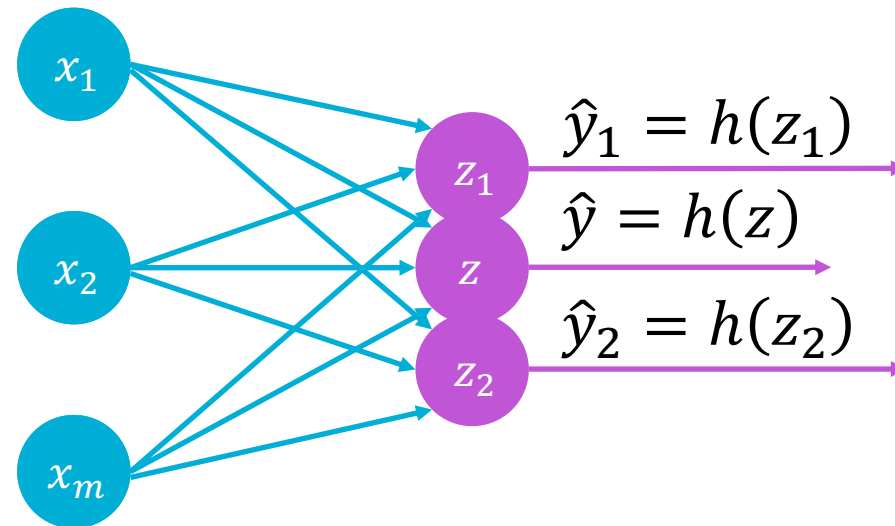
## Simplified notation



$$\hat{y} = h\left(w_0 + \sum_{j=1}^{m} x_j w_j\right)$$
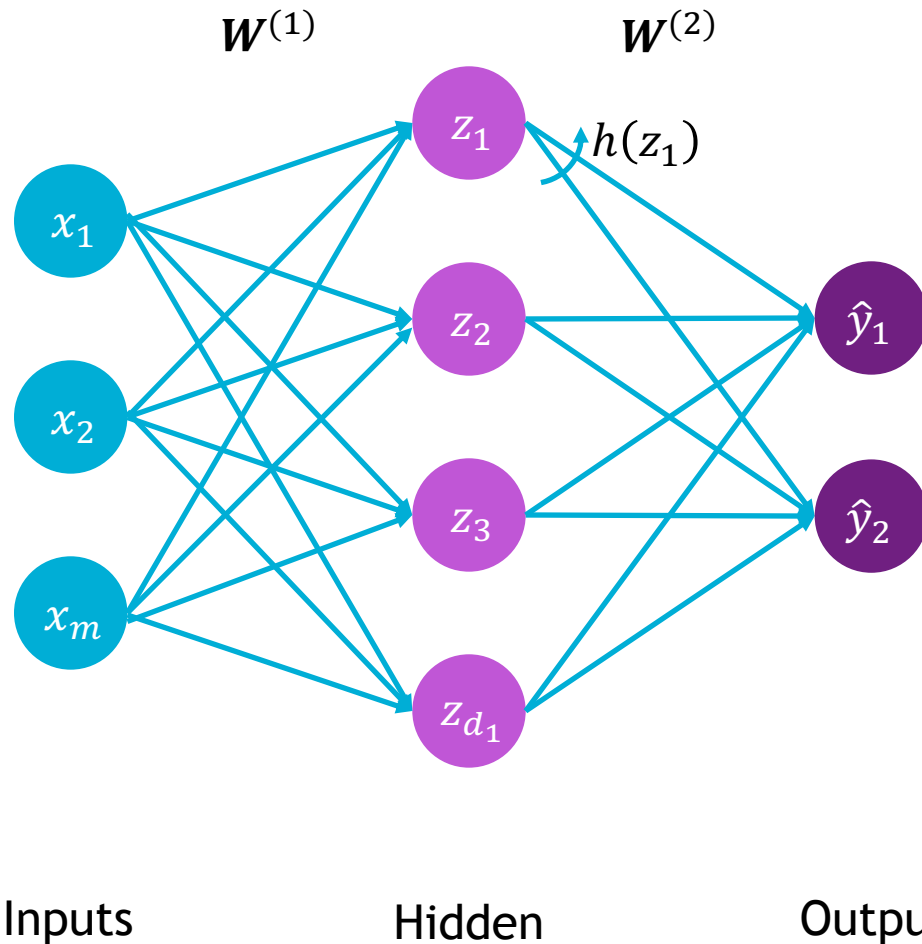
$$z = w_0 + \sum_{j=1}^{m} x_j w_j$$

$$\hat{y} = h(z)$$

## Multi output neural network with dense layers



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

MIT Introduction to Deep Learning (introtodeeplearning.com)

$W^{(1)}$  $W^{(2)}$

$h(z_1)$

$x_1$  $x_2$  $x_m$
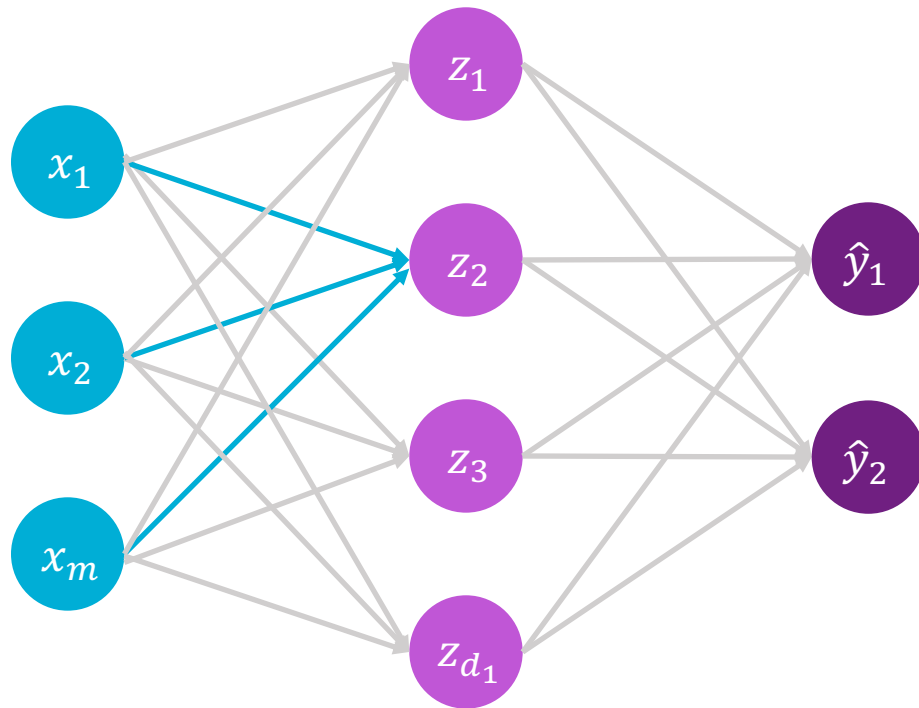
$z_1$  $z_2$  $z_3$  $z_{d_1}$

$\hat{y}_1$  $\hat{y}_2$

Inputs  Hidden  Output

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

$$\hat{y}_i = h\left(w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j w_{j,i}^{(2)}\right)$$

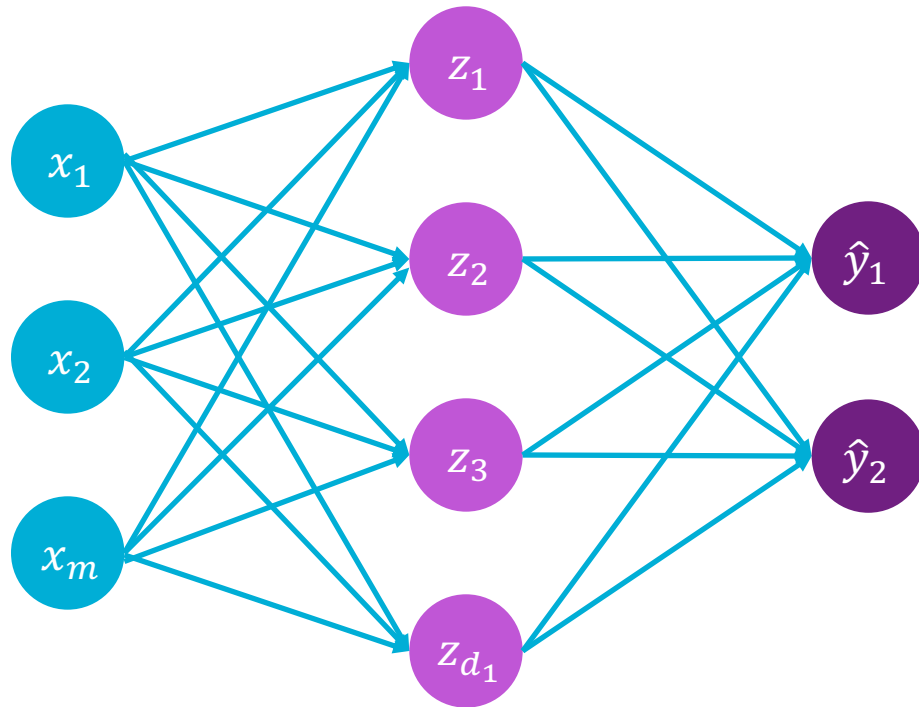MIT Introduction to Deep Learning (introtodeeplearning.com)

$$z_2 = w_{0,2}^{(1)} + \sum_{j=1}^{m} x_j w_{j,2}^{(1)} \mathsf{x}$$

$$= w_{0,2}^{(1)} + x_1 w_{1,2}^{(1)} + x_2 w_{2,2}^{(1)} + x_m w_{m,2}^{(1)}$$

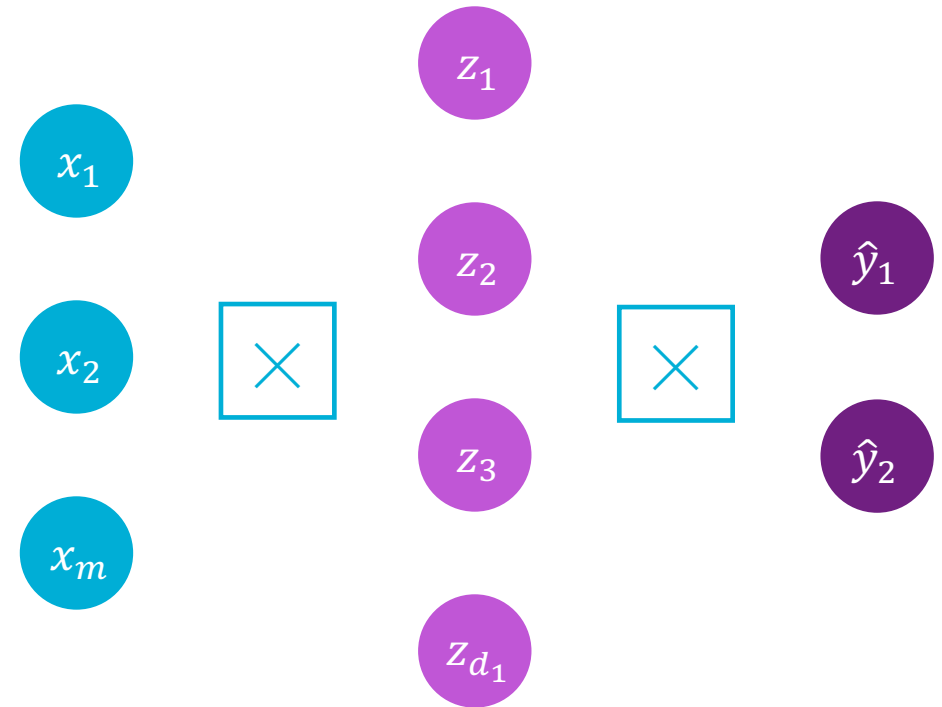Inputs      Hidden      Output

MIT Introduction to Deep Learning (introtodeeplearning.com)

# Single layer neural network

## Simplified notation



Inputs        Hidden        Output        Inputs        Hidden        Output

MIT Introduction to Deep Learning (introtodeeplearning.com)

# Deep neural network

## Stacking layers

$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{d_{k-1}} h(z_{k-1,j}) w_{j,i}^{(k)}$$

Inputs

Hidden

Output

MIT Introduction to Deep Learning (introtodeeplearning.com)

# Applying neural networks

## Diagnosis of dementia based on imaging biomarkers



Is this subject healthy?

$x_2$ = brain metabolism

$x_1$ = volume of fluid in the brain

Healthy

Dementia

MIT Introduction to Deep Learning (introtodeeplearning.com)

## Diagnosis of dementia based on imaging biomarkers

Is this subject healthy?



$$x^{(1)} = [1, 1.5]$$

Predicted: 0.1
Actual: 1

Loss: $l\big(f\big(x^{(i)}; \boldsymbol{W}\big), y^{(i)}\big)$

Predicted  Actual

MIT Introduction to Deep Learning (introtodeeplearning.com)

## Diagnosis of dementia based on imaging biomarkers

Is this subject healthy?



$$X = \begin{bmatrix} 1 & 1.5 \\ 3.5 & 1.2 \\ 5 & 0.9 \\ \vdots & \vdots \end{bmatrix}$$

$f(x)$ $\quad$ $y$

$$\begin{bmatrix} 0.1 \\ 0.4 \\ 0.6 \\ \vdots \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \end{bmatrix}$$
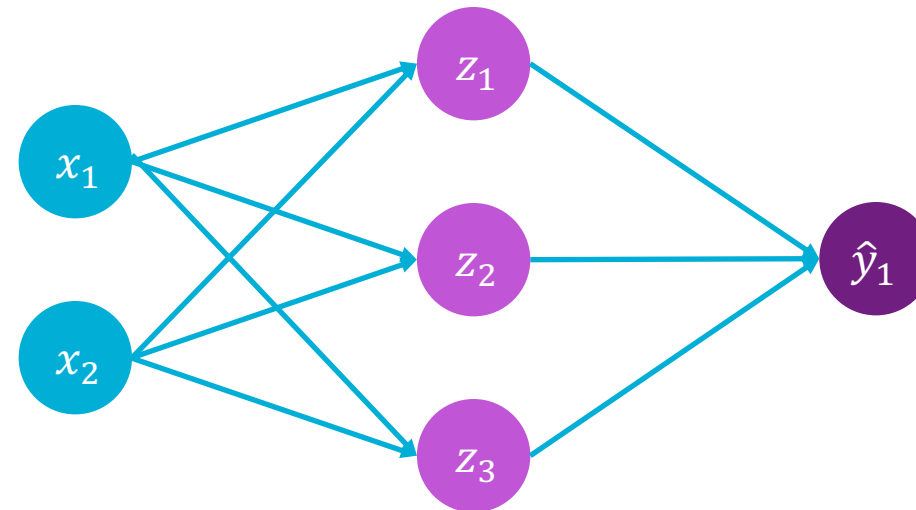
Cost function: $\quad J(\boldsymbol{W}) = \dfrac{1}{n}\sum_{i=1}^{n} l\big(\underline{f\big(x^{(i)}; \boldsymbol{W}\big)}, \underline{y^{(i)}}\big)$

Predicted  Actual

MIT Introduction to Deep Learning (introtodeeplearning.com)

## Diagnosis of dementia based on imaging biomarkers

Is this subject healthy?



$$X = \begin{bmatrix} 1 & 1.5 \\ 3.5 & 1.2 \\ 5 & 0.9 \\ \vdots & \vdots \end{bmatrix}$$

$f(x)$ $\quad y$

$$\begin{bmatrix} 0.1 \\ 0.4 \\ 0.6 \\ \vdots \end{bmatrix} \quad \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \end{bmatrix}$$

Cost function with cross entropy loss:

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

Actual  Predicted  Actual  Predicted

## Diagnosis of dementia based on imaging biomarkers

What is the dementia severity?



$$X = \begin{bmatrix} 1 & 1.5 \\ 3.5 & 1.2 \\ 5 & 0.9 \\ \vdots & \vdots \end{bmatrix}$$

$$f(x) \quad\quad y$$

$$\begin{bmatrix} 3 \\ 1 \\ 8 \\ \vdots \end{bmatrix} \quad\quad \begin{bmatrix} 0 \\ 2 \\ 7 \\ \vdots \end{bmatrix}$$

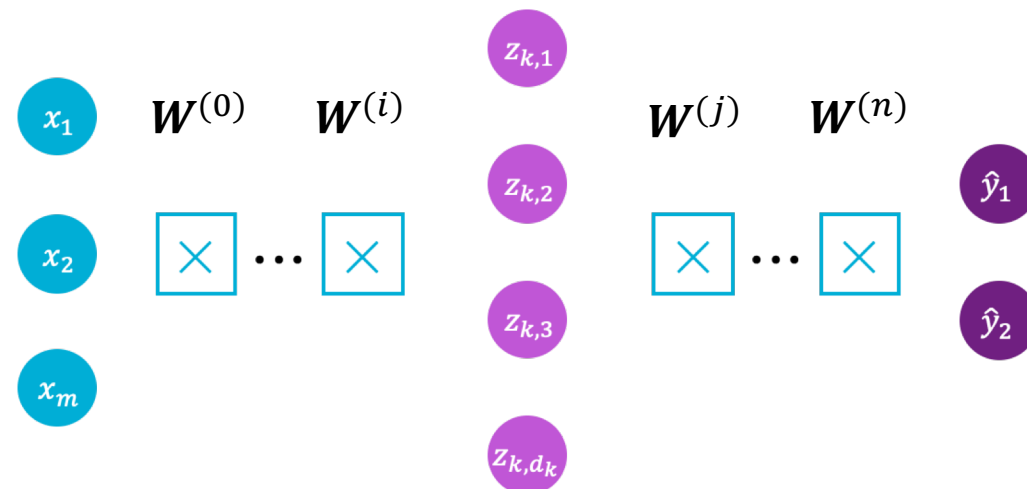Cost function with mean squared error loss:

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

Actual          Predicted

## Loss optimisation

- **Find the network weights that achieve the lowest loss**

$$W^* = \operatorname*{argmin}_{W} J(W) = \operatorname*{argmin}_{W} \frac{1}{n} \sum_{i=1}^{n} l\big(f\big(x^{(i)}; W\big), y^{(i)}\big)$$

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



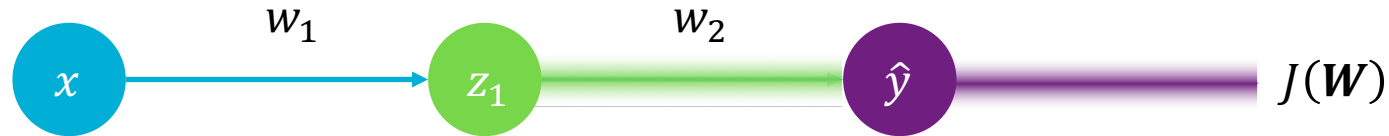MIT Introduction to Deep Learning (introtodeeplearning.com)

## Gradient descent

### Algorithm

1. Initialise weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence

    a. Compute gradient $\frac{\partial J(\boldsymbol{W})}{\partial \boldsymbol{W}}$

    b. Update weights $\boldsymbol{W} \leftarrow \boldsymbol{W} - \eta \frac{\partial J(\boldsymbol{W})}{\partial \boldsymbol{W}}$

3. Return weights

## Computing gradients: backpropagation
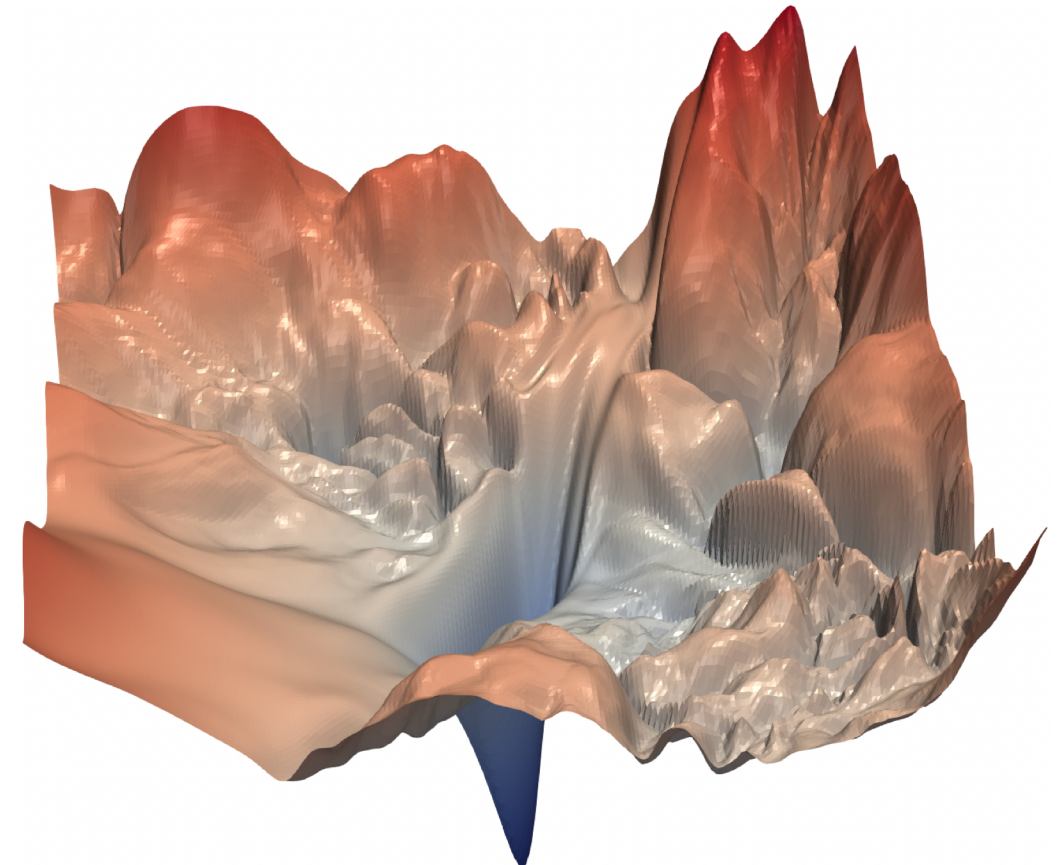
## Computing gradients: backpropagation



$$\frac{\partial J(\boldsymbol{W})}{\partial w_2} = \frac{\partial J(\boldsymbol{W})}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_2}$$

## Computing gradients: backpropagation



$$\frac{\partial J(\boldsymbol{W})}{\partial w_2} = \frac{\partial J(\boldsymbol{W})}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_2}$$

$$\frac{\partial J(\boldsymbol{W})}{\partial w_1} = \frac{\partial J(\boldsymbol{W})}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_1} = \frac{\partial J(\boldsymbol{W})}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

MIT Introduction to Deep Learning (introtodeeplearning.com)

ARAMIS
LAB
BRAIN DATA SCIENCE

## Gradient descent in practice

### Algorithm

1. Initialise weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence

   a. Compute gradient $\frac{\partial J(\boldsymbol{W})}{\partial \boldsymbol{W}}$

   b. Update weights $\boldsymbol{W} \leftarrow \boldsymbol{W} - \eta \frac{\partial J(\boldsymbol{W})}{\partial \boldsymbol{W}}$

3. Return weights

Li et al., Visualizing the Loss Landscape of Neural Nets, NIPS, 2018

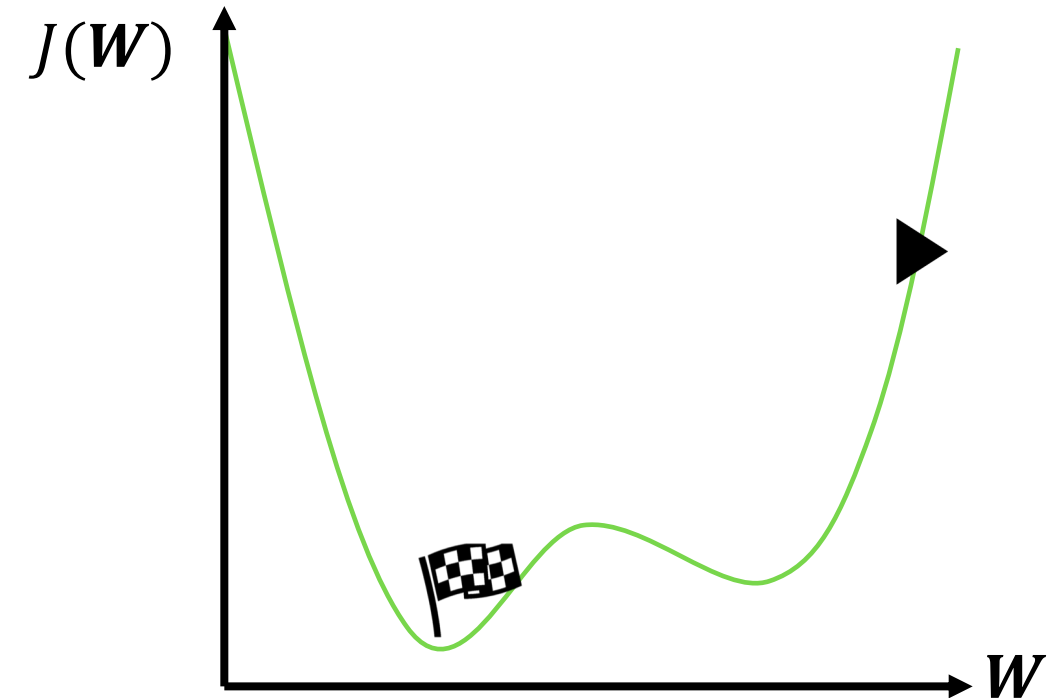MIT Introduction to Deep Learning (introtodeeplearning.com)

## Learning rate

$$W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$$

- If $\eta$ is too small: slow to converge, may be trapped in local minima

## Learning rate

$$W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$$

- If $\eta$ is too small: slow to converge, may be trapped in local minima

- If $\eta$ is too large: may diverge

## Learning rate

$$W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$$

- If $\eta$ is too small: slow to converge, may be trapped in local minima

- If $\eta$ is too large: may diverge

→ Adaptative learning rate



$J(W)$

$W$

## Adaptative learning rate

- Learning rates are no longer fixed

- Can be made larger or smaller depending on:

  - how large the gradient is

  - how fast learning is happening

  - the size of particular weights

  - etc.

- Algorithms:

  - Adam [Kingma et al., Adam: A Method for Stochastic Optimization, 2014]

  - Adadelta [Zeiler et al., ADADELTA: An Adaptive Learning Rate Method, 2012]

  - Adagrad [Duchi et al., Adaptive Subgradient Methods for Online Learning and Stochastic Optimization, 2011]

  - RMSProp [Hinton, Neural Networks for Machine Learning]

## Standard gradient descent

### Algorithm

1. Initialise weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence

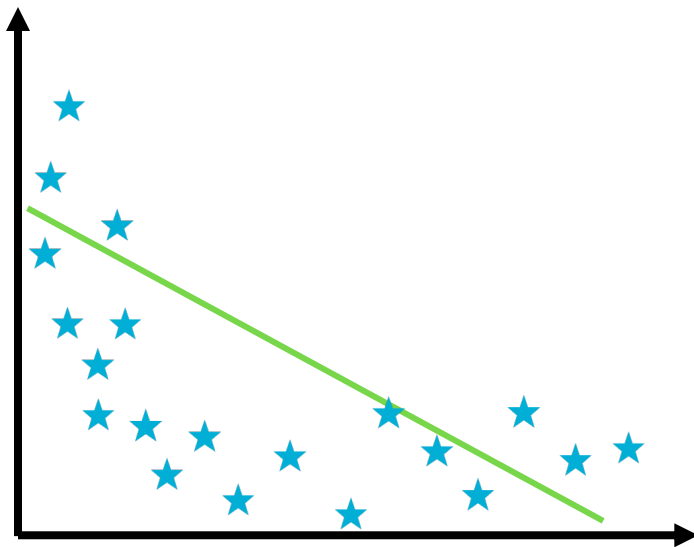   a. Compute gradient $\dfrac{\partial J(\boldsymbol{W})}{\partial \boldsymbol{W}}$    Can be expensive to compute

   b. Update weights $\boldsymbol{W} \leftarrow \boldsymbol{W} - \eta \dfrac{\partial J(\boldsymbol{W})}{\partial \boldsymbol{W}}$

3. Return weights

## Stochastic gradient descent with one sample

### Algorithm

1. Initialise weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence

   a. Pick single data point $i$

   b. Compute gradient $\dfrac{\partial J_i(\boldsymbol{W})}{\partial \boldsymbol{W}}$    Easy to compute but **very noisy**

   c. Update weights $\boldsymbol{W} \leftarrow \boldsymbol{W} - \eta \dfrac{\partial J(\boldsymbol{W})}{\partial \boldsymbol{W}}$

3. Return weights

## Stochastic gradient descent with several samples (mini-batch)

**Algorithm**

1. Initialise weights randomly $\sim \mathcal{N}(0, \sigma^2)$

2. Loop until convergence

   a. Pick batch of $B$ data points

   *Fast to compute and good estimate of the true gradient*

   b. Compute gradient $\boxed{\dfrac{\partial J(\boldsymbol{W})}{\partial \boldsymbol{W}} = \dfrac{1}{B} \sum_{k=1}^{B} \dfrac{\partial J_k(\boldsymbol{W})}{\partial \boldsymbol{W}}}$

   c. Update weights $\boldsymbol{W} \leftarrow \boldsymbol{W} - \eta \dfrac{\partial J(\boldsymbol{W})}{\partial \boldsymbol{W}}$

3. Return weights

## Overfitting



Underfitting            Ideal fit            Overfitting

## Regularisation: Dropout



MIT Introduction to Deep Learning (introtodeeplearning.com)

## Regularisation: Early stopping



Loss

Underfitting

Overfitting

Testing

Training

Training iterations

# Summary

## The Perceptron

- Structural building blocks
- Nonlinear activation functions

## Neural Networks

- Stacking Perceptrons to form neural networks
- Optimisation through backpropagation

## Training in Practice

- Adaptive learning
- Batching
- Regularisation



MIT Introduction to Deep Learning (introtodeeplearning.com)

# **Convolutional neural networks**

## Diagnosis of dementia based on imaging biomarkers

Is this subject healthy?



$x^{(1)} = [1, 1.5]$

volume of fluid in the brain → $x_1$

brain metabolism → $x_2$

$z_1$

$z_2$

$z_3$

$\hat{y}_1$

Predicted: 0.1

MIT Introduction to Deep Learning (introtodeeplearning.com)

## Image = matrix of numbers

# Using an image as input of a neural network

## Fully connected neural network



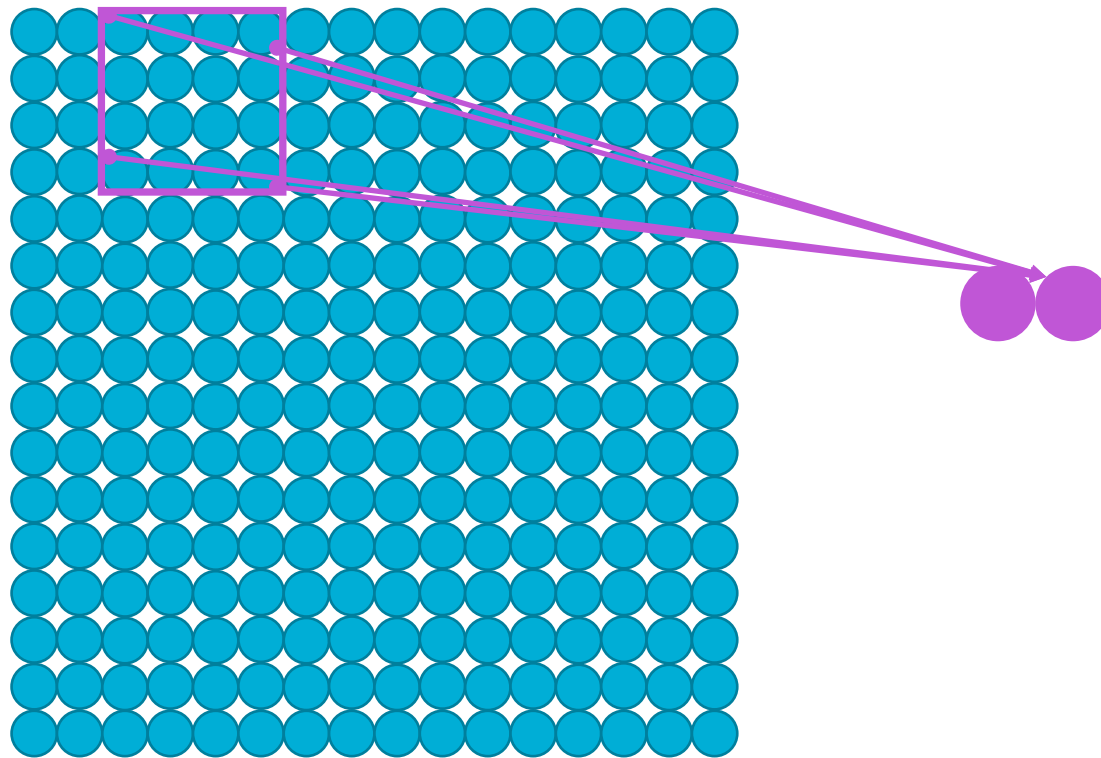No spatial information

Many, many parameters

## Using spatial features



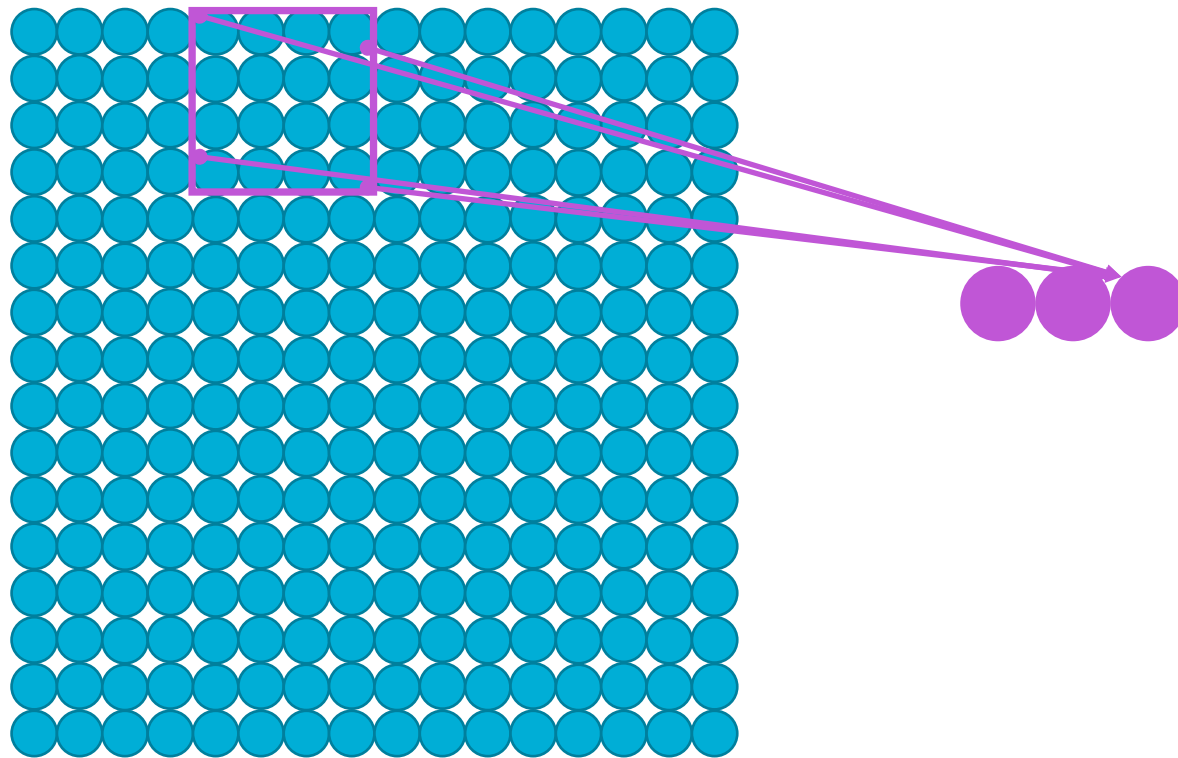**Idea**: connect patches of input to neurons in hidden layer

## Using spatial features

- Slide patch window across input image

- Weight pixels inside the patch
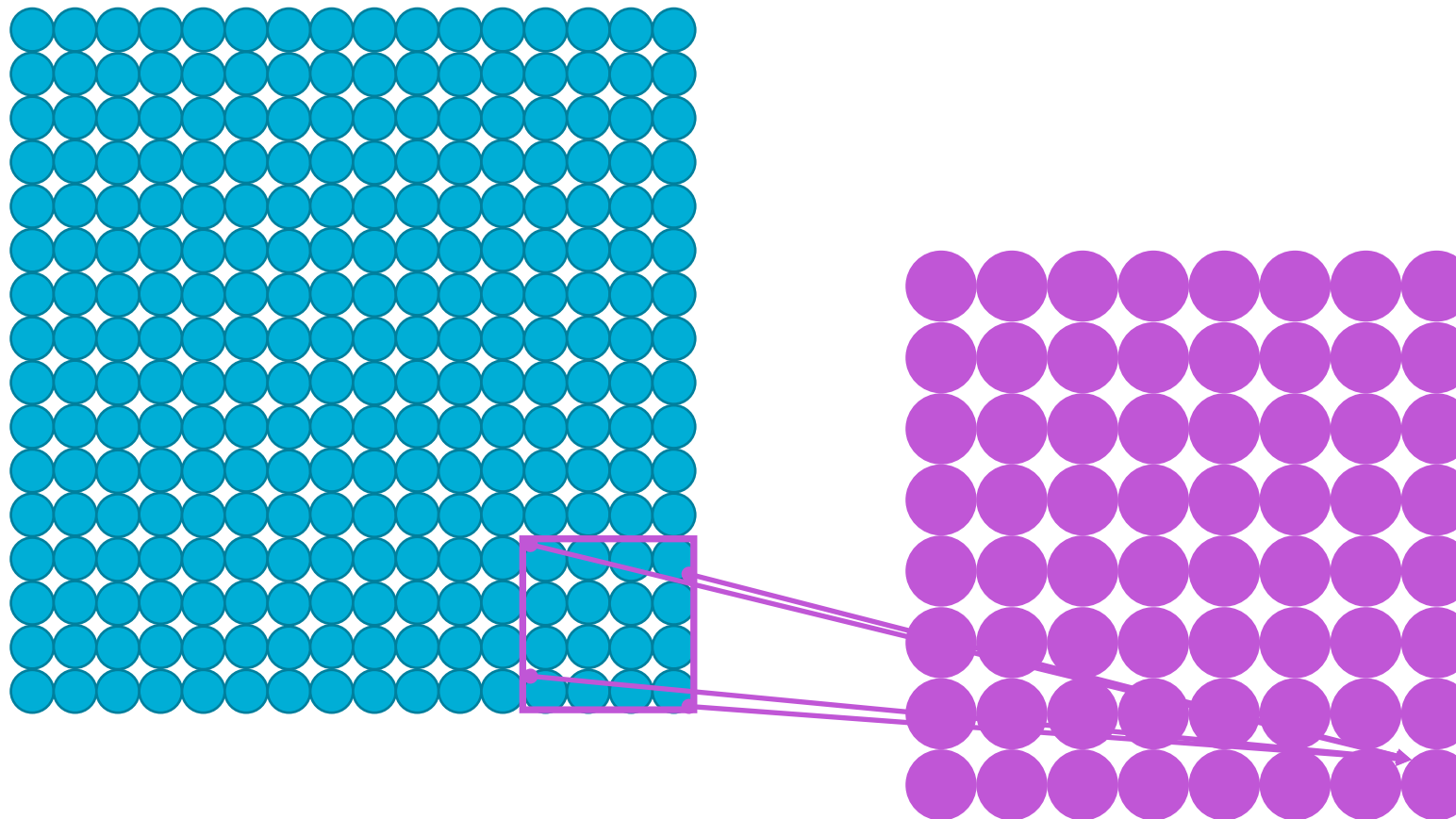
- Apply weighted summation

→ Convolution

## Using spatial features



- Slide patch window across input image

- Weight pixels inside the patch

- Apply weighted summation

→ Convolution

## Using spatial features



- Slide patch window across input image

- Weight pixels inside the patch

- Apply weighted summation

→ Convolution

MIT Introduction to Deep Learning (introtodeeplearning.com)

## The convolution operation

- Slide the filter over the input image

- Element-wise multiply

- Add the outputs

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

$\otimes$

| | | |
|---|---|---|
| 1 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 1 |

Image                                Filter

MIT Introduction to Deep Learning (introtodeeplearning.com)

## The convolution operation

- **Slide the filter over the input image**

- **Element-wise multiply**

- **Add the outputs**

$$1{\times}1 + 1{\times}0 + 1{\times}1$$
$$+\ 0{\times}0 + 1{\times}1 + 1{\times}0$$
$$+\ 0{\times}1 + 0{\times}0 + 1{\times}1$$
$$=\ 4$$



Image         Filter         Feature map

MIT Introduction to Deep Learning (introtodeeplearning.com)

## The convolution operation

- **Slide the filter over the input image**

- **Element-wise multiply**

- **Add the outputs**

$$1 \times 1 + 1 \times 0 + 0 \times 1$$
$$+ 1 \times 0 + 1 \times 1 + 1 \times 0$$
$$+ 0 \times 1 + 1 \times 0 + 1 \times 1$$
$$= 3$$



Image          Filter          Feature map

## The convolution operation

- **Slide the filter over the input image**

- **Element-wise multiply**

- **Add the outputs**

$$1{\times}1 + 1{\times}0 + 1{\times}1$$
$$+ 1{\times}0 + 1{\times}1 + 0{\times}0$$
$$+ 1{\times}1 + 0{\times}0 + 0{\times}1$$
$$= 4$$



Image       Filter       Feature map

MIT Introduction to Deep Learning (introtodeeplearning.com)

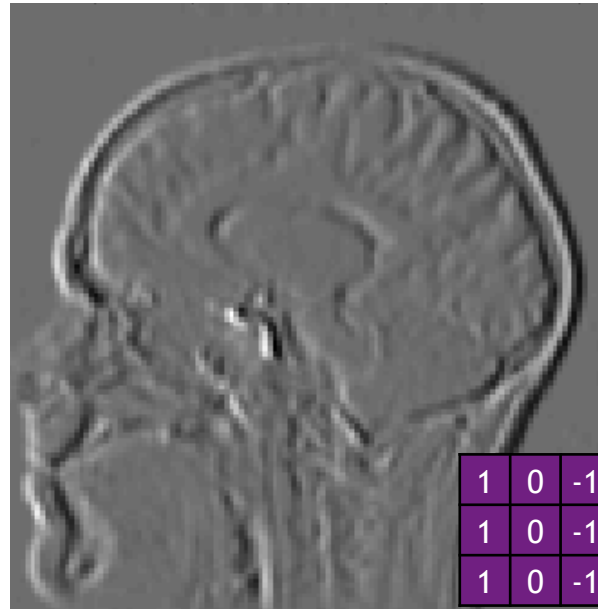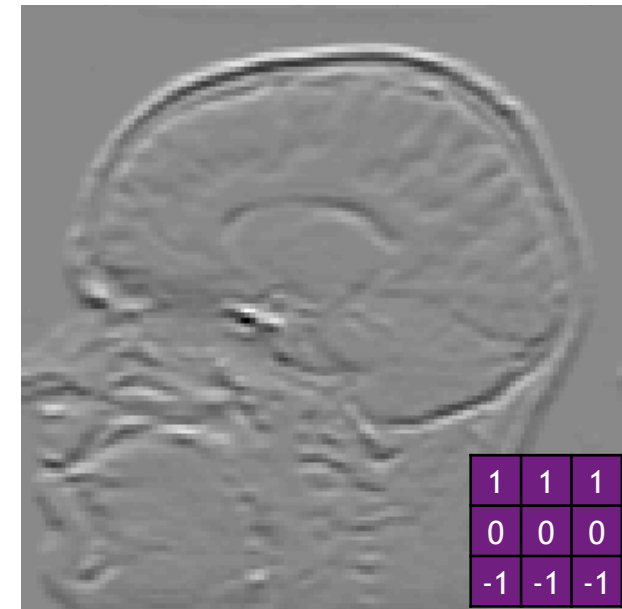## Different filters = different feature maps
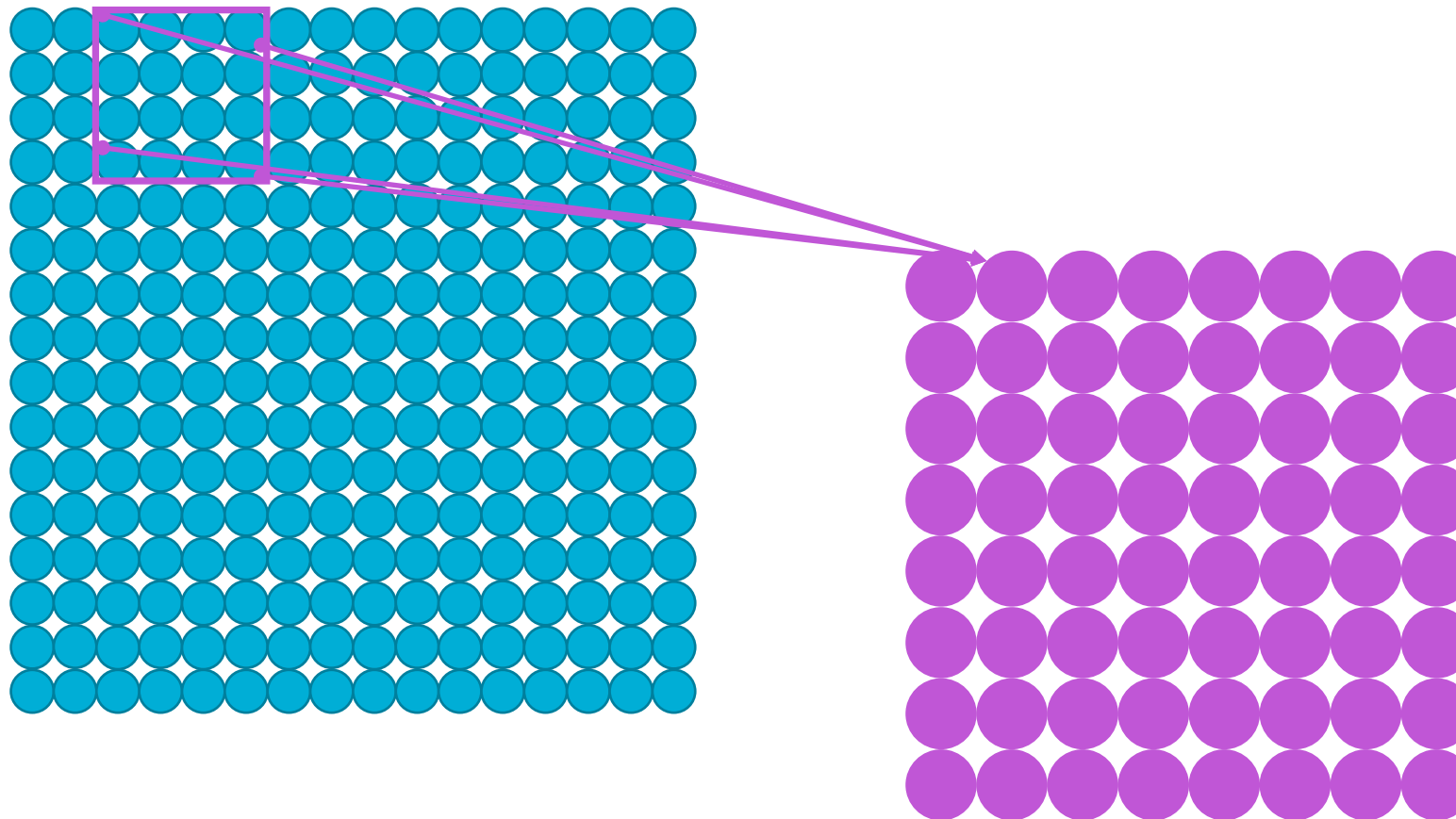


Original image        Vertical edge detection        Horizontal edge detection
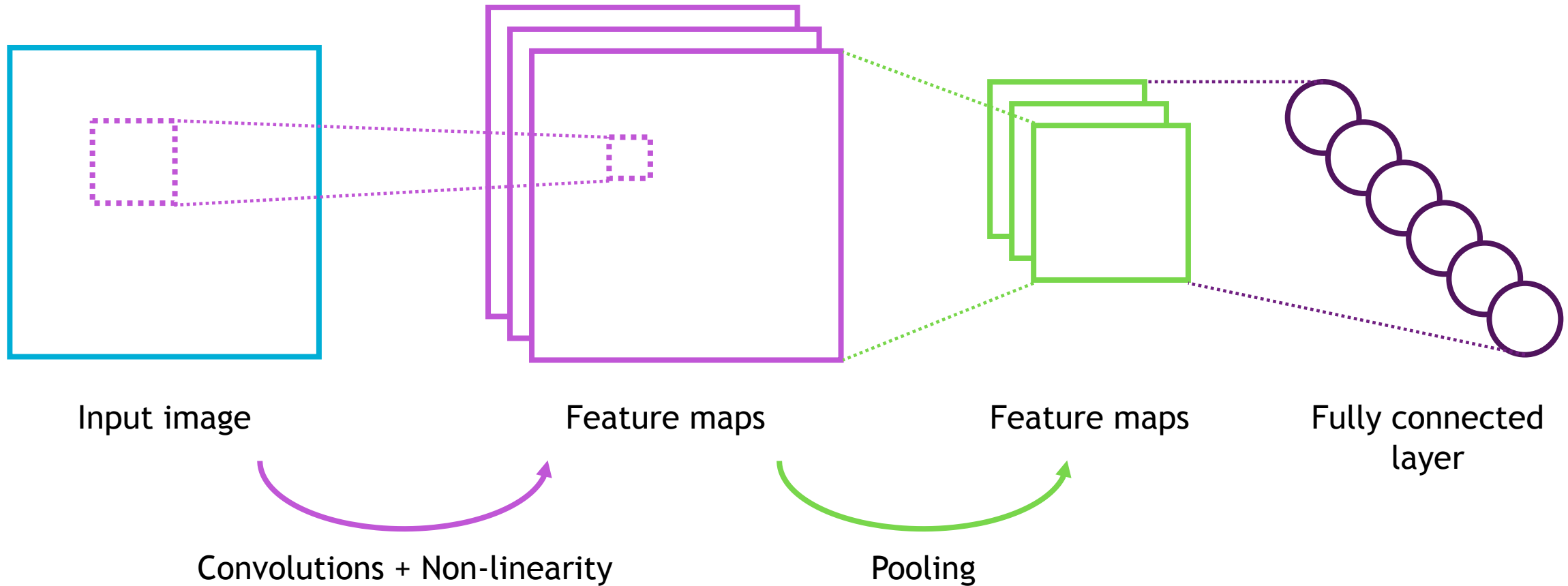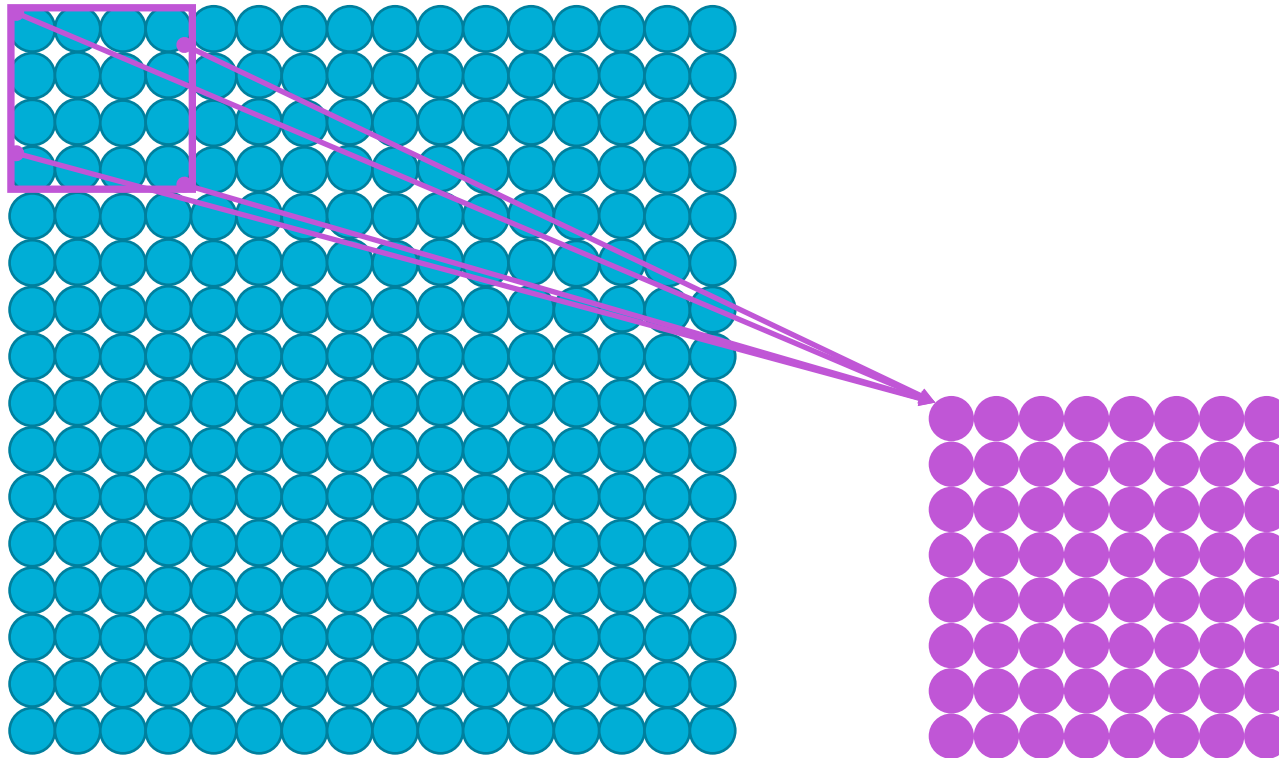
## Using spatial features



- Apply a set of weights – a filter – to extract **local features**

- Use **multiple filters** to extract different features

# Convolutional neural networks

## CNNs for classification



Input image      Feature maps      Feature maps      Fully connected layer

Convolutions + Non-linearity      Pooling
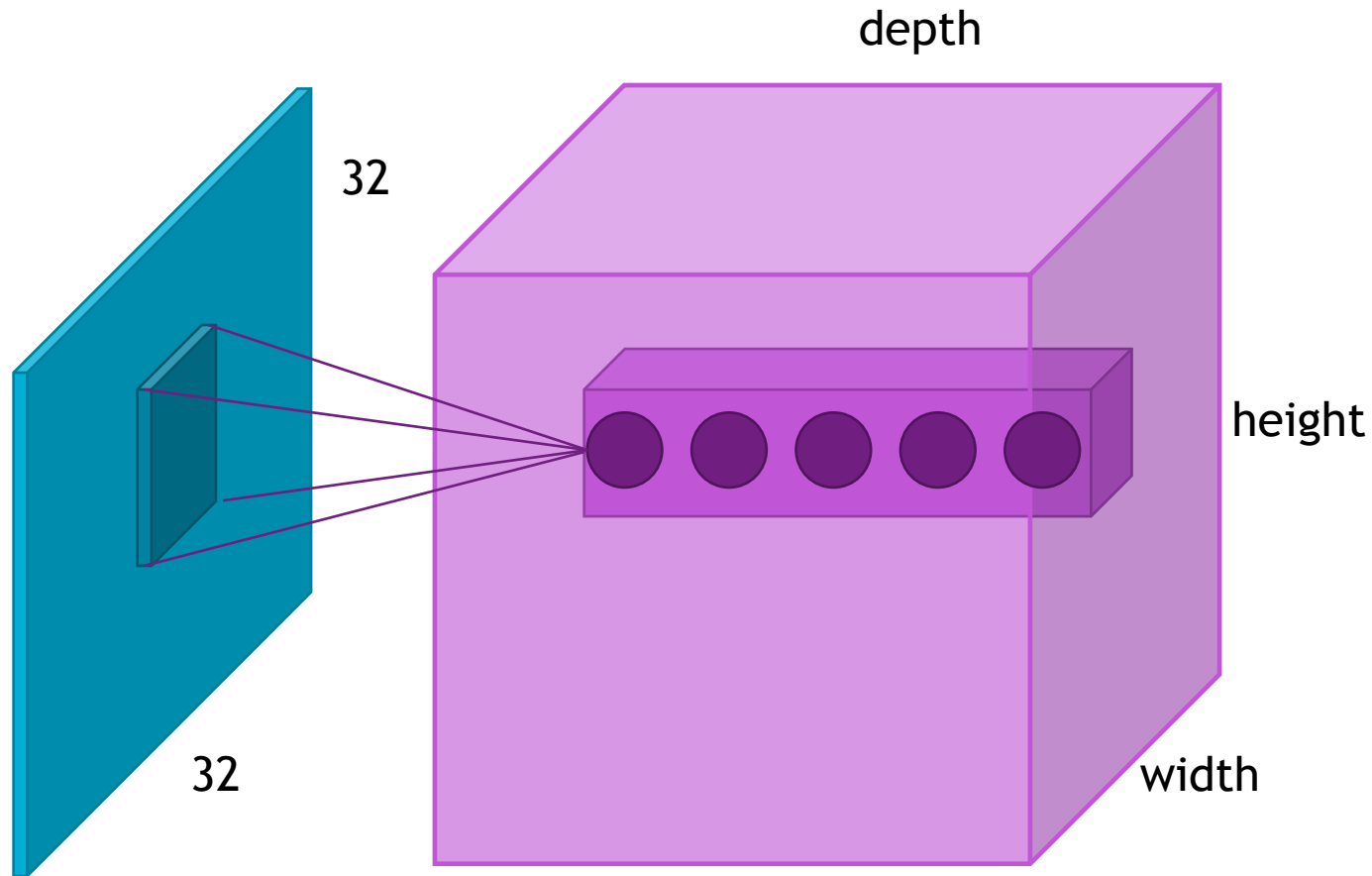
## Convolutional layer



For a neuron in hidden layer:
- Take inputs from patch
- Compute weighted sum
- Apply bias
- Activate with non-linear function

Neuron $(p, q)$ in hidden layer
Filter size 4×4
Weights $w_{i,j}$

$$h\left(\sum_{i=1}^{4}\sum_{j=1}^{4} w_{i,j} x_{i+p,j+q} + b\right)$$

MIT Introduction to Deep Learning (introtodeeplearning.com)

## Spatial arrangement of output volume



depth

32

32

1

height

width

**Layer Dimensions:**
$$h \times w \times d$$
$h$ & $w$ = spatial dimensions
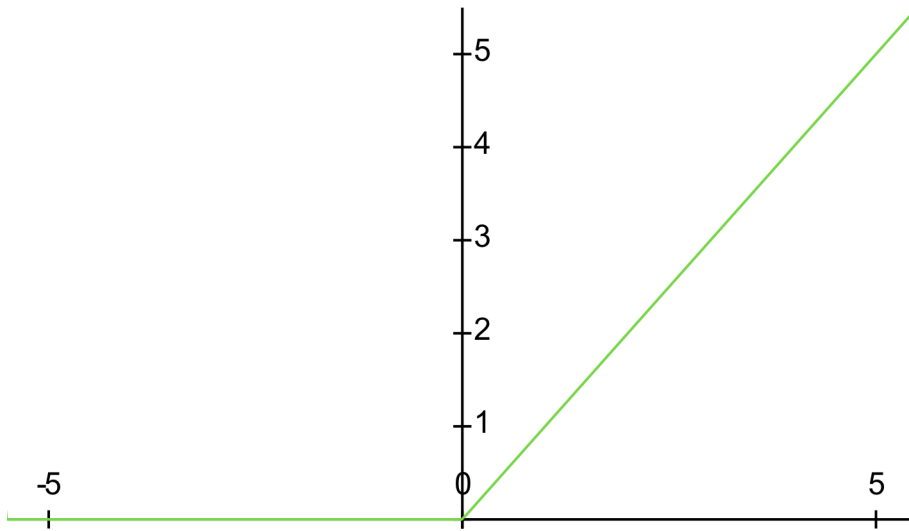$d$ = number of filters

**Stride:**
Filter step size

**Receptive Field:**
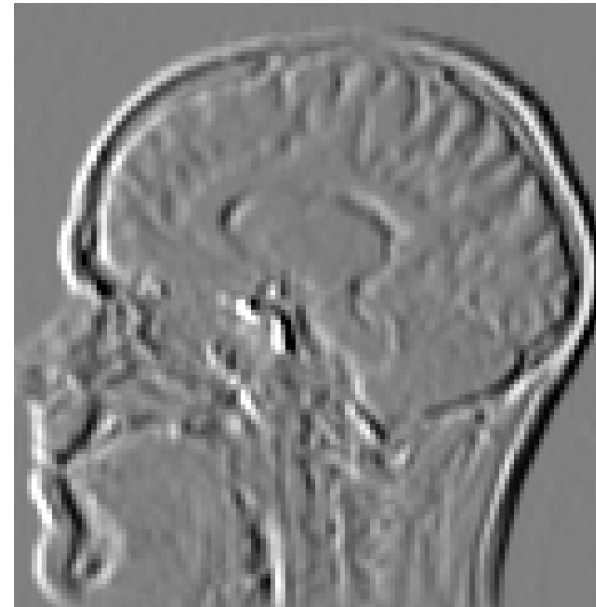Locations in input image that a node is connected to

## Introducing non-linearity
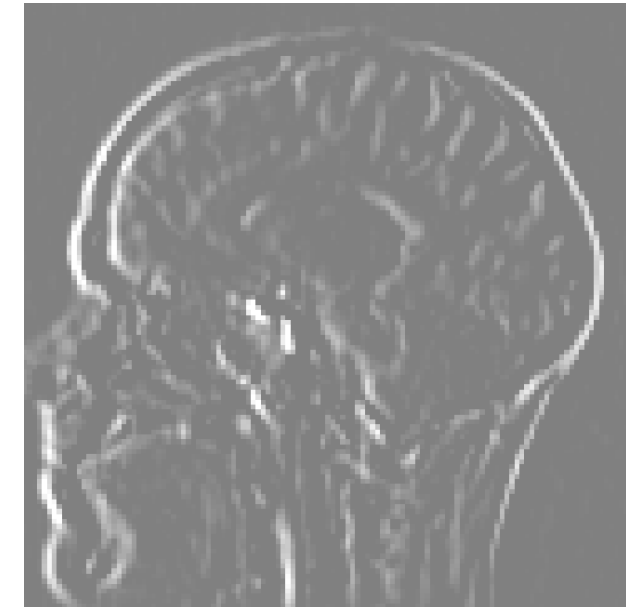
### Rectified linear unit (ReLU)



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

### Input feature map          Rectified feature map



ReLU

Black: negative values - White: positive values

MIT Introduction to Deep Learning (introtodeeplearning.com)

## Pooling

- **Reduce dimensionality while preserving spatial invariance**
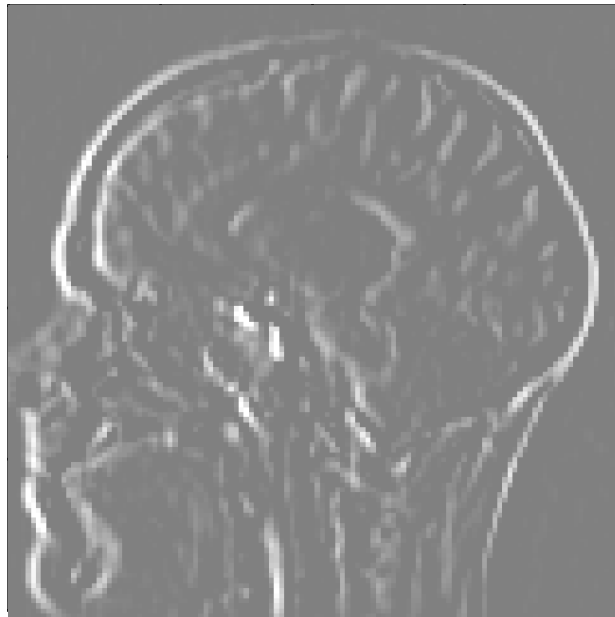


Input feature map

Max pooling with
2×2 filter and stride 2

Pooled feature map

## Pooling

- **Reduce dimensionality while preserving spatial invariance**

Input feature map
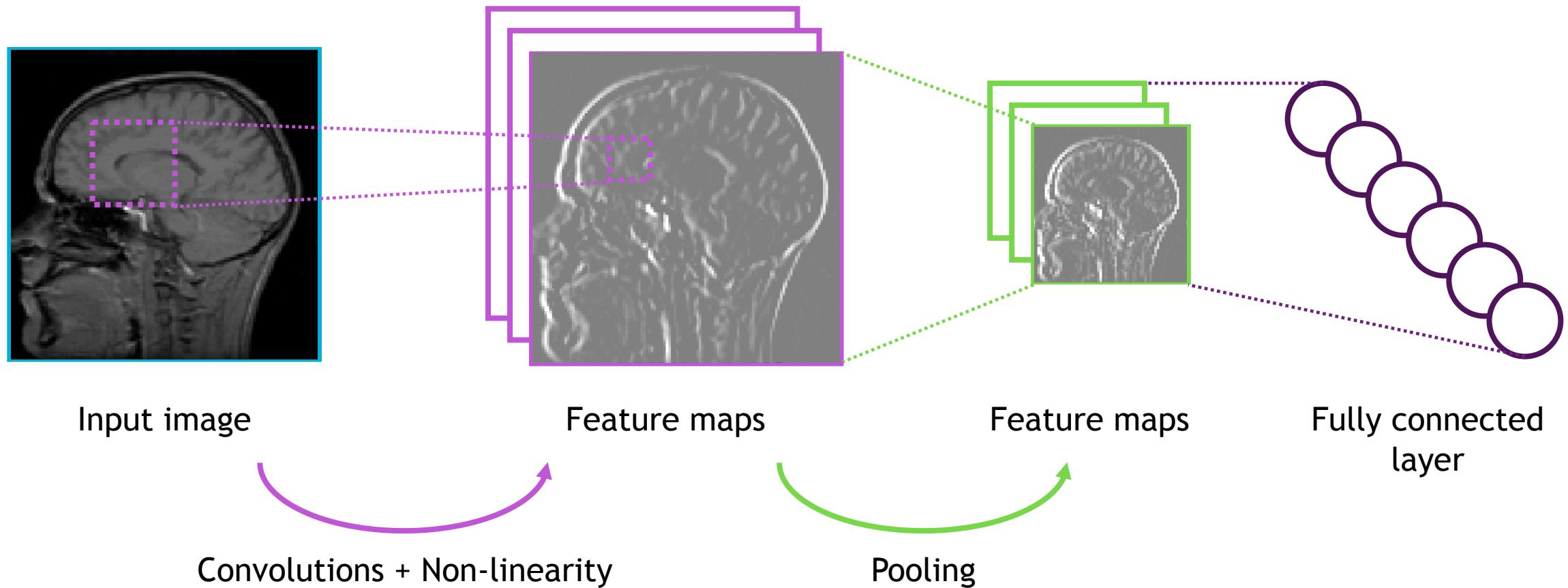
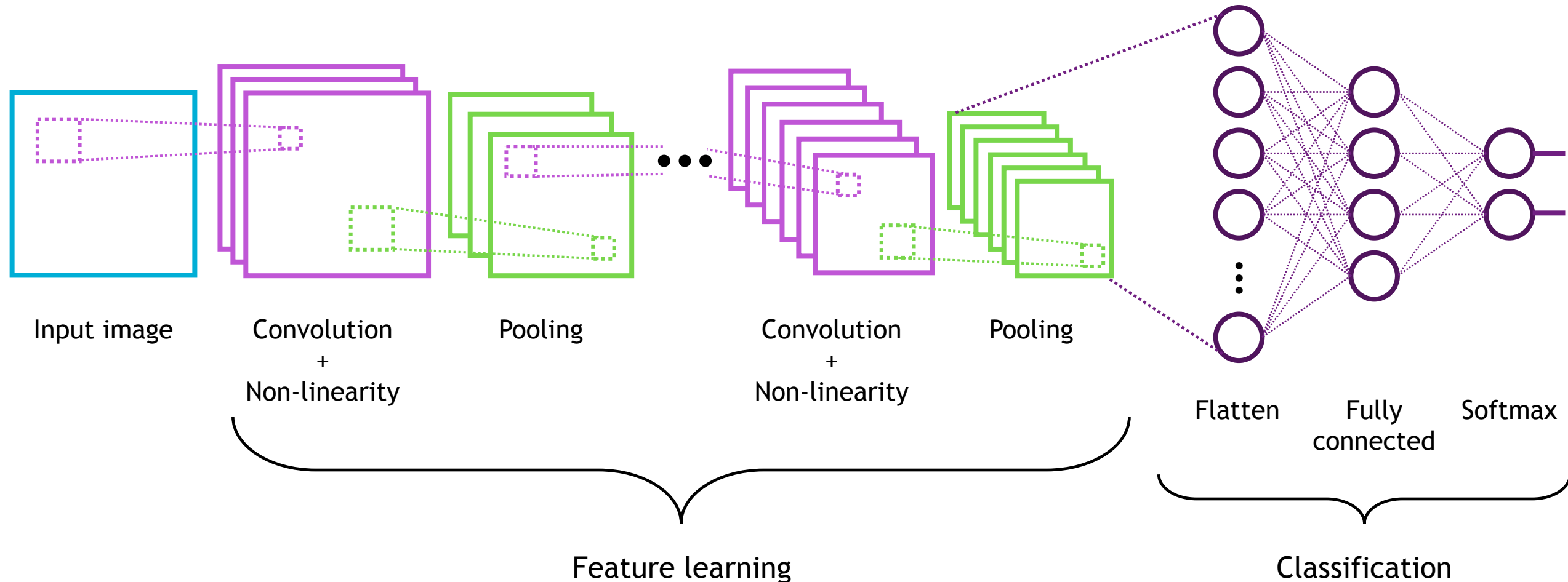Max pooling with
2×2 filter and stride 2

→

Pooled feature map

## CNNs for classification



Input image | Feature maps | Feature maps | Fully connected layer

Convolutions + Non-linearity

Pooling

## CNNs for classification



Input image — Convolution + Non-linearity — Pooling — Convolution + Non-linearity — Pooling — Flatten — Fully connected — Softmax

Feature learning — Classification

## CNNs for many applications



Input image

Convolution
+
Non-linearity

Pooling

Convolution
+
Non-linearity

Pooling

Feature learning

Organ segmentation

Lesion detection

Image denoising

Flatten      Fully        Softmax
             connected
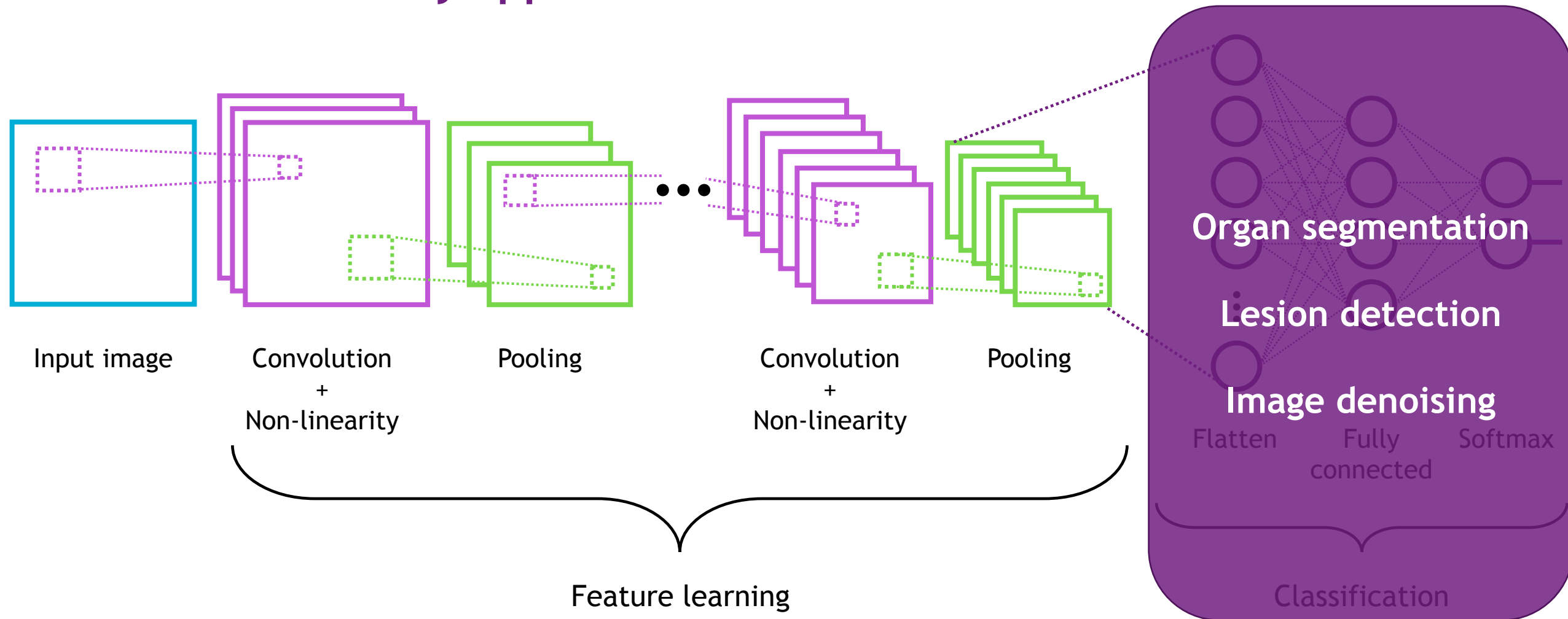
Classification

# Summary

## Images

- Representing images
- Convolutions for feature extraction
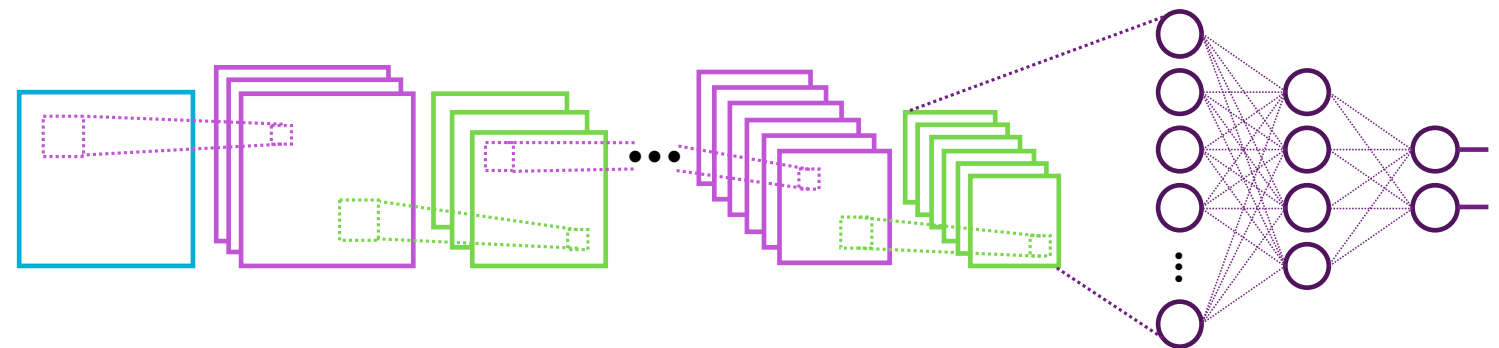
## CNNs
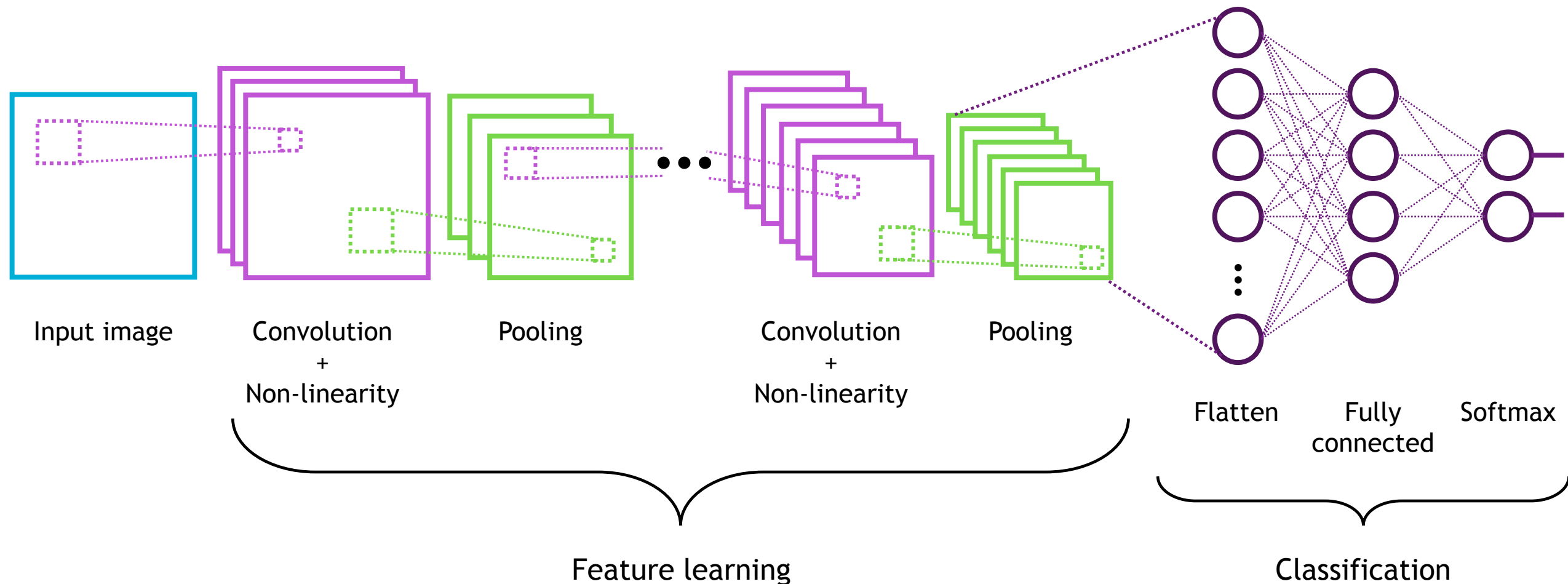
- Convolution → non-linearity → pooling
- Stacking layers

## Applications

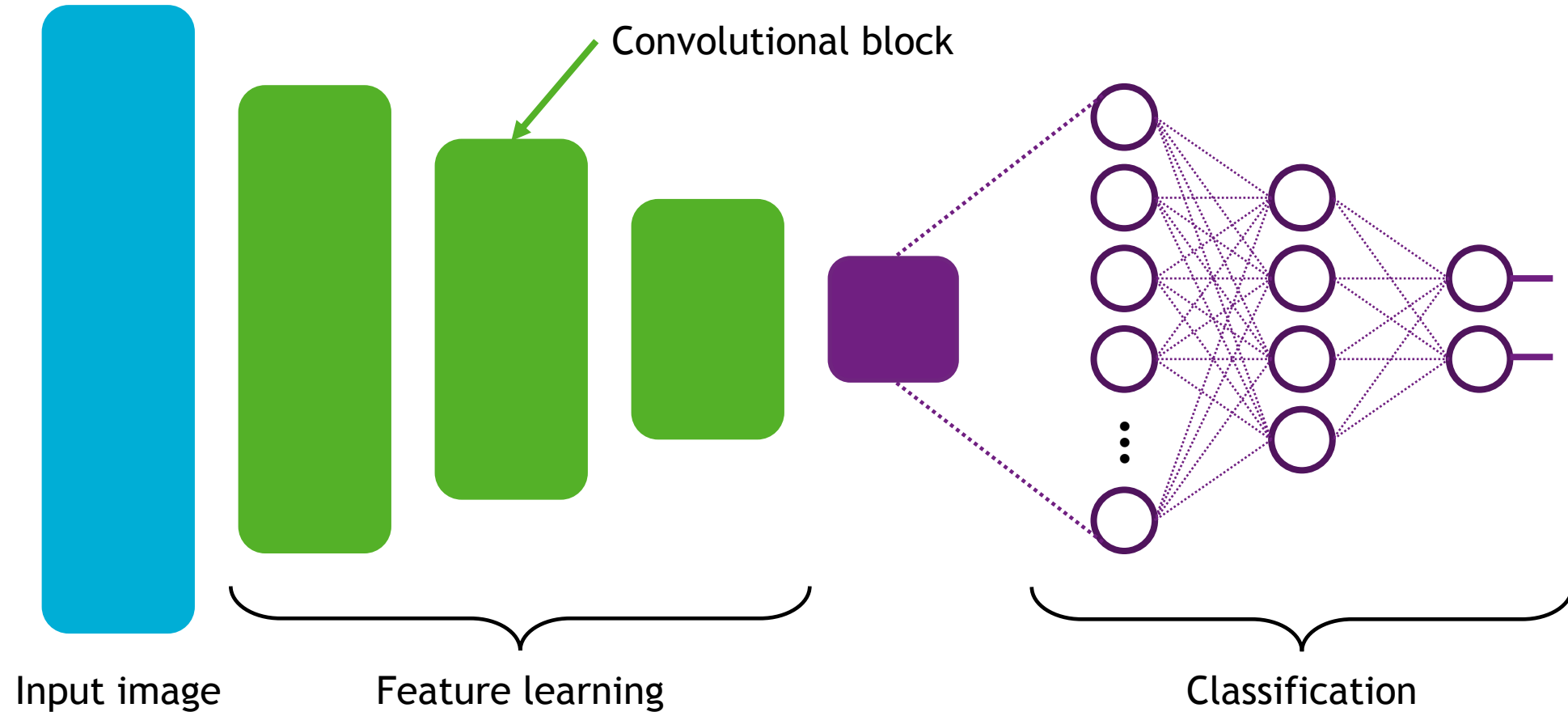- Classification
- Segmentation
- Detection

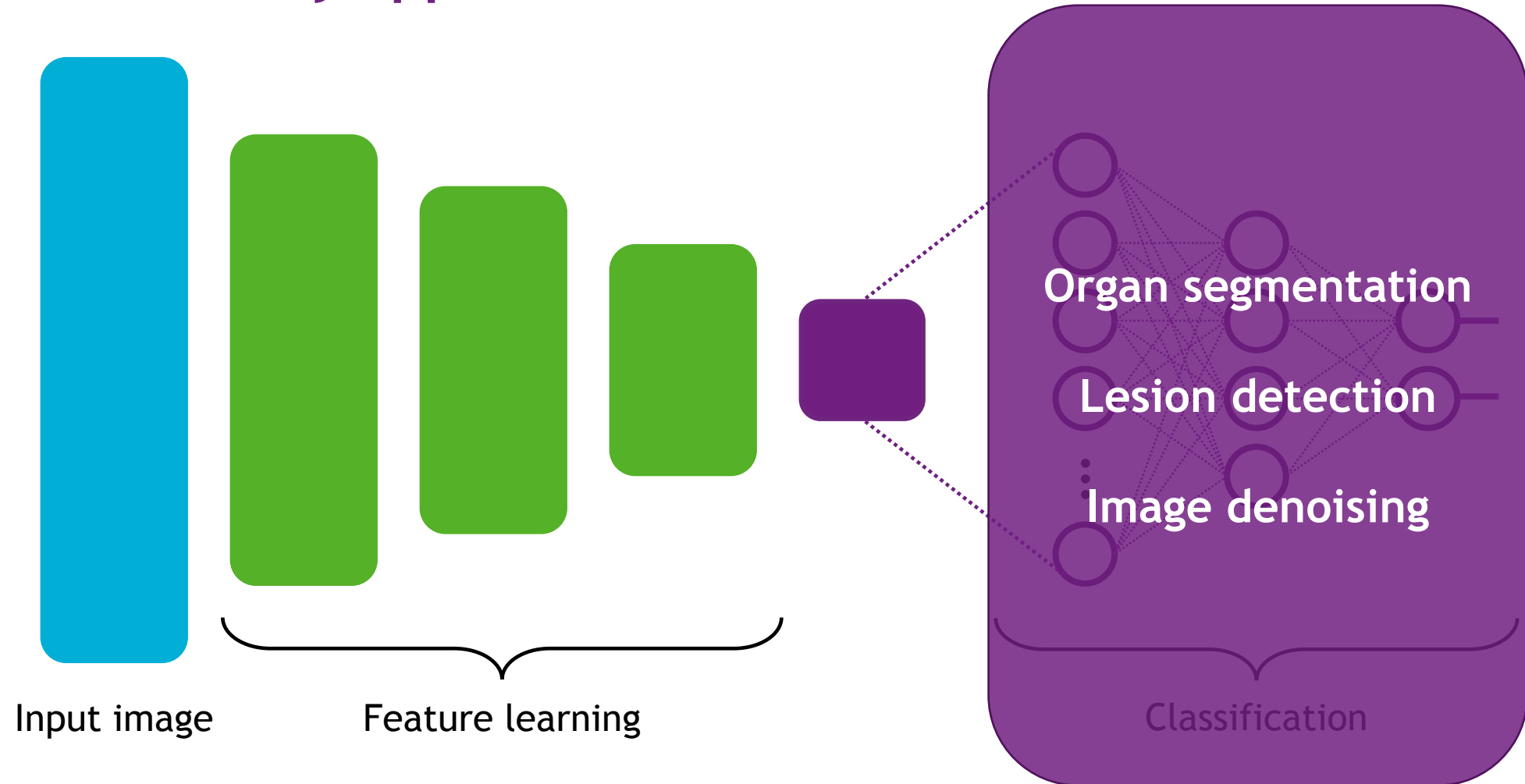| 1 | 0 | -1 |
|---|---|----|
| 1 | 0 | -1 |
| 1 | 0 | -1 |

MIT Introduction to Deep Learning (introtodeeplearning.com)

# Generative Deep Learning

## CNNs for classification



Input image — Convolution + Non-linearity — Pooling — ⋯ — Convolution + Non-linearity — Pooling — Flatten — Fully connected — Softmax

Feature learning

Classification

MIT Introduction to Deep Learning (introtodeeplearning.com)

## CNNs for classification



Convolutional block

Input image — Feature learning — Classification

# Convolutional neural networks

## CNNs for many applications

Input image

Feature learning

**Organ segmentation**

**Lesion detection**

**Image denoising**

Classification

MIT Introduction to Deep Learning (introtodeeplearning.com)

## CNNs for image generation



Latent space

Input image | Feature learning | Reconstruction | Reconstructed image

Observed variable

Latent variable

# Autoencoders
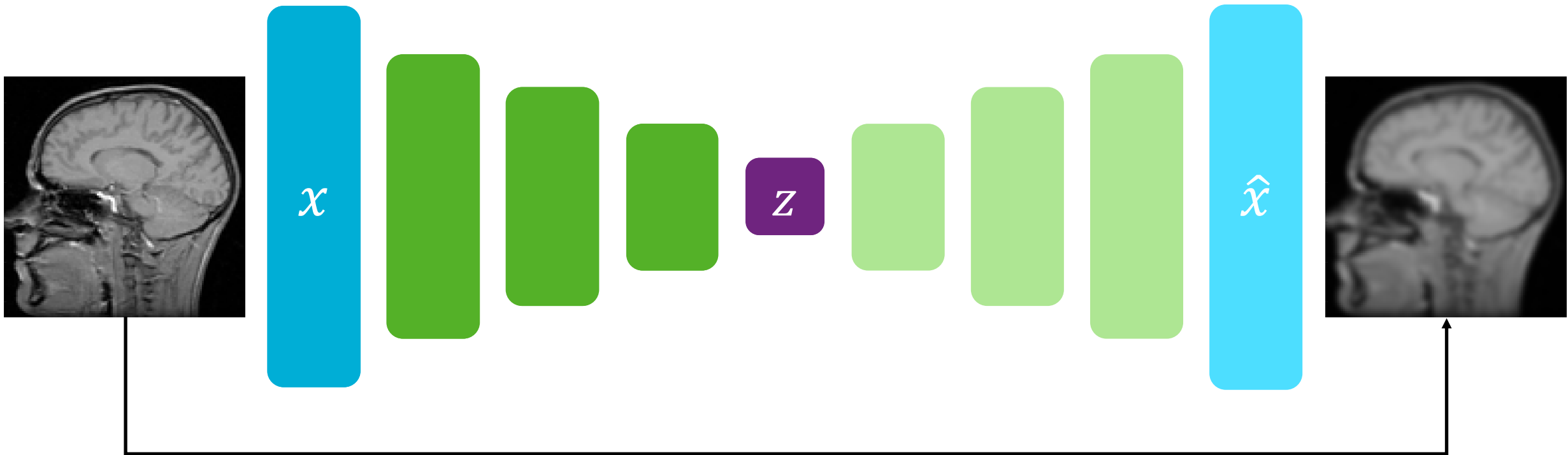
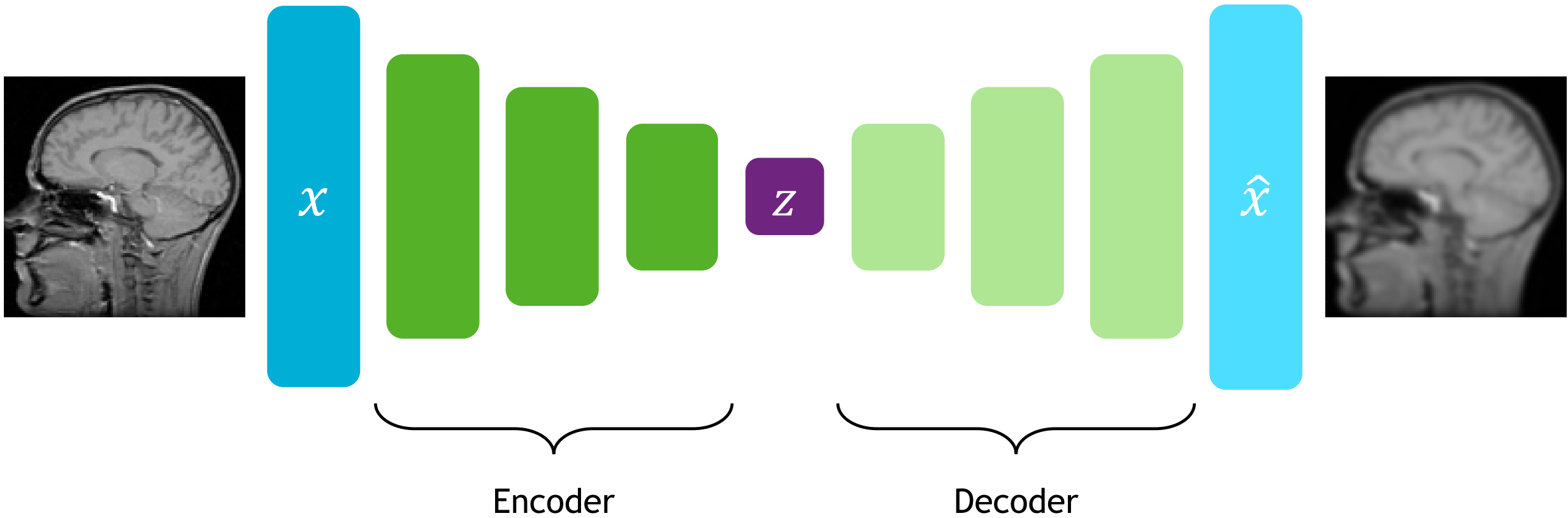## Encoder



Input image
=
observed data

Latent space
=
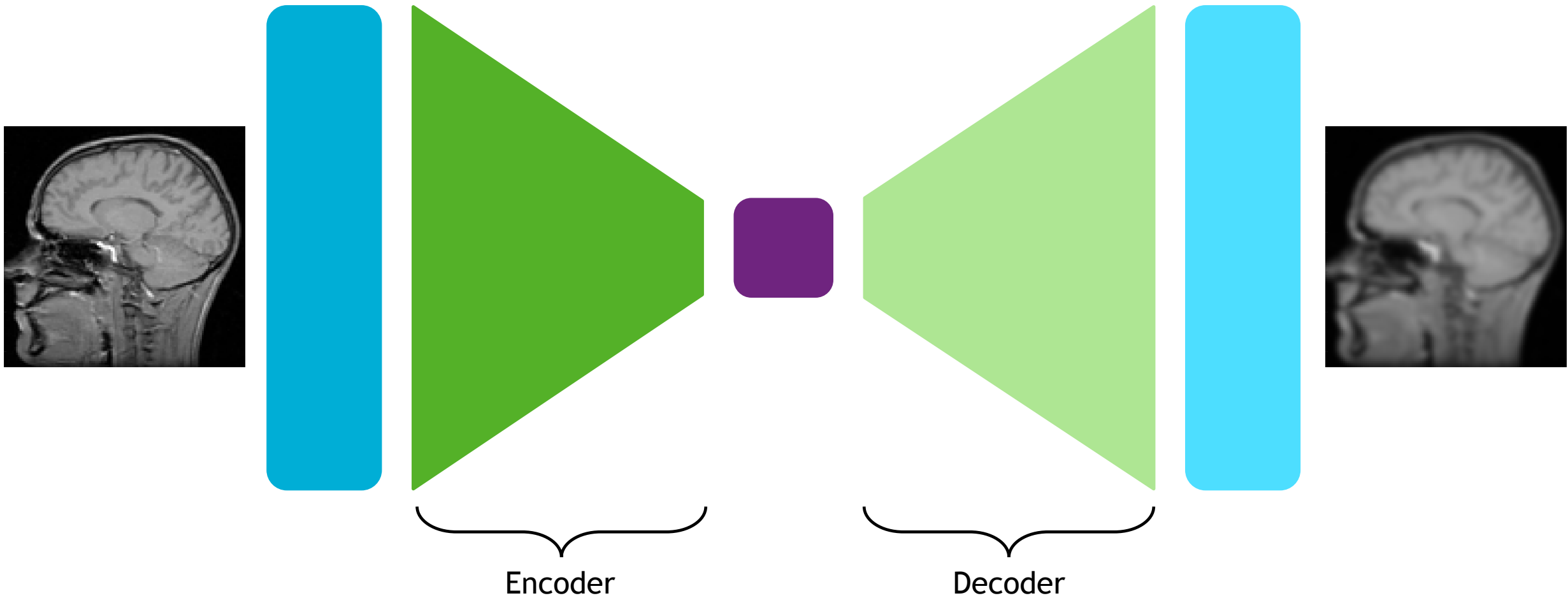low dimensional representation
of the observed data

MIT Introduction to Deep Learning (introtodeeplearning.com)

## Training autoencoders



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

Encoder

Decoder

# Autoencoders

Encoder

Decoder

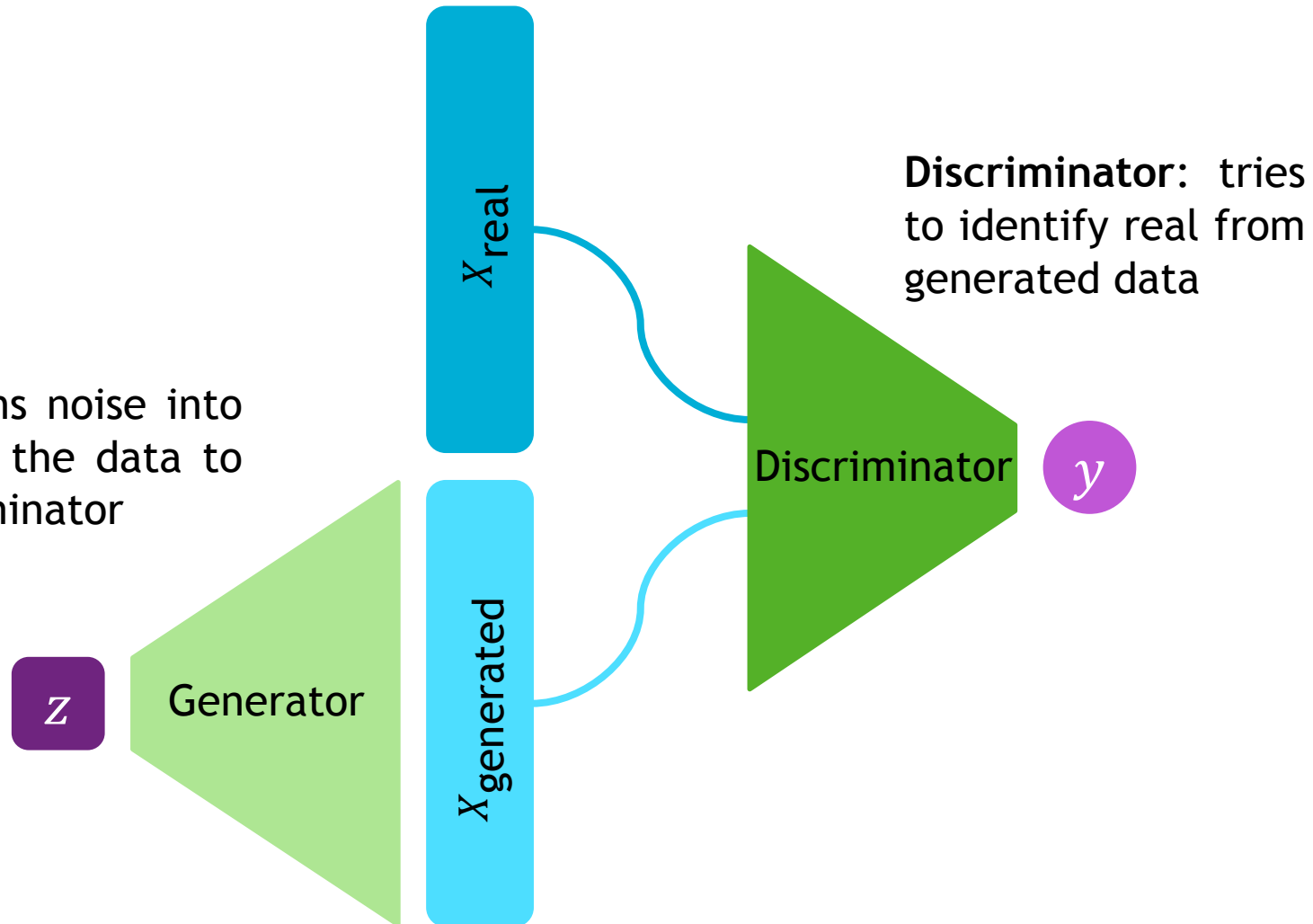MIT Introduction to Deep Learning (introtodeeplearning.com)

## Generating images from scratch



Random noise

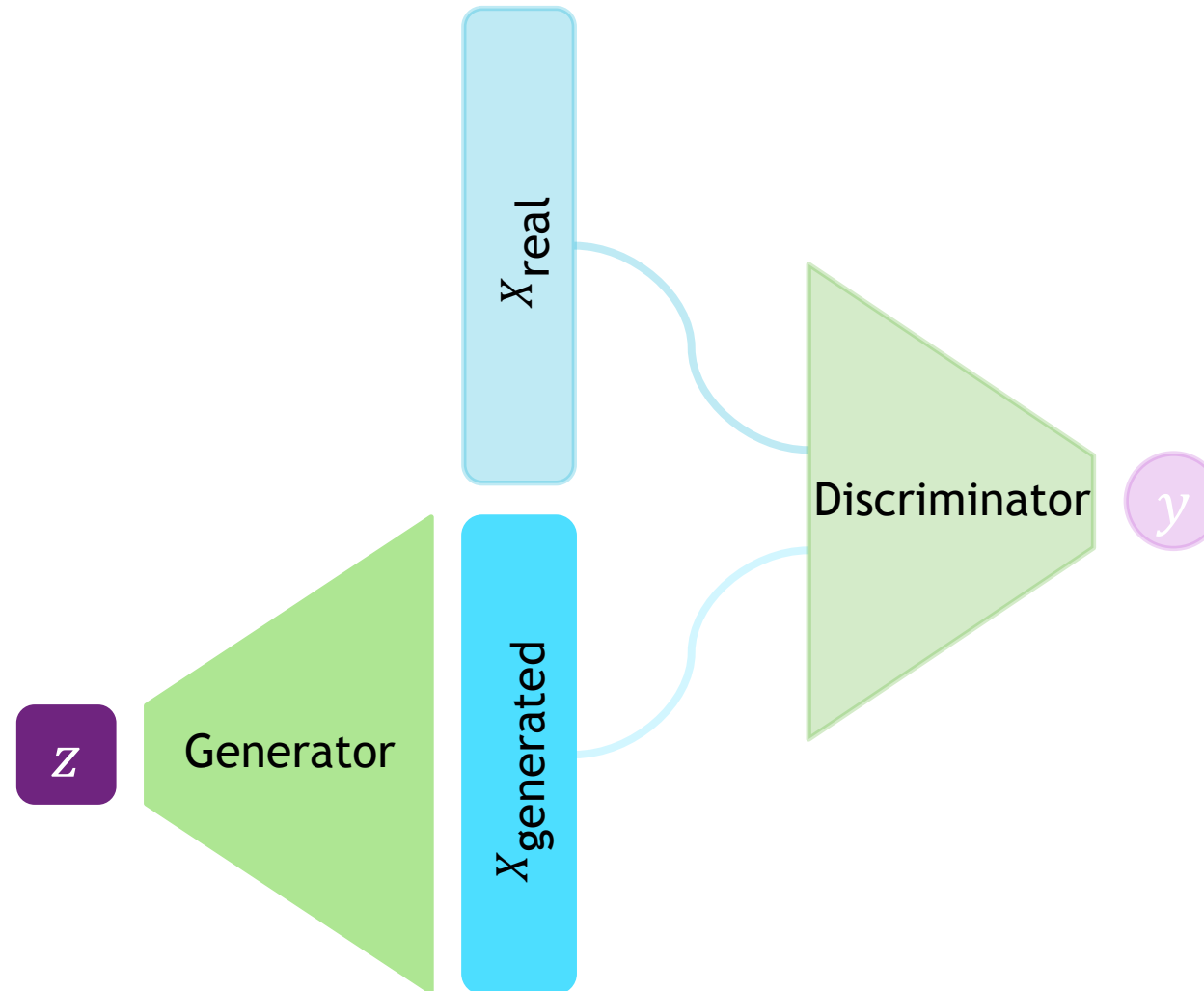## Competing networks



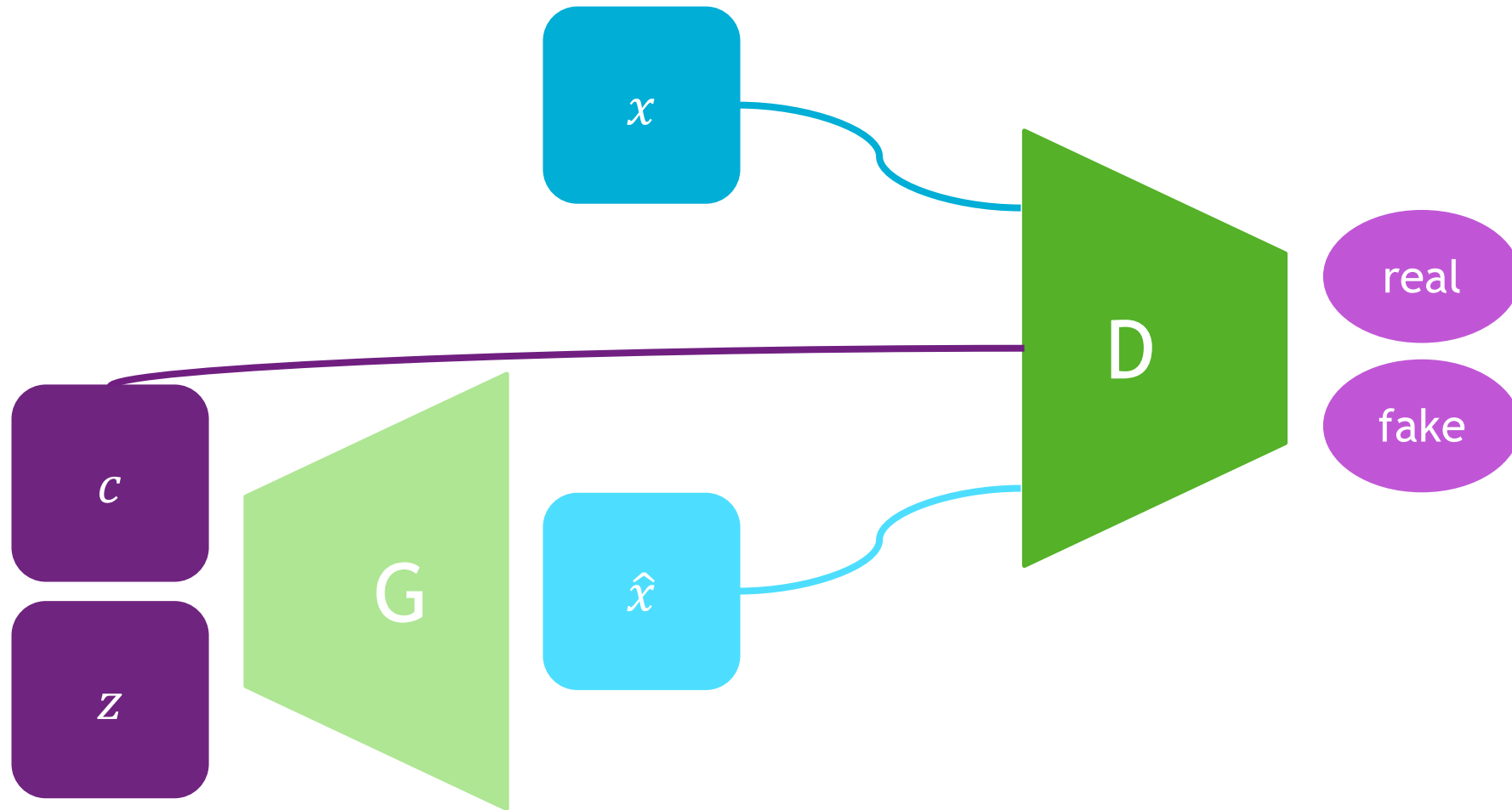**Generator**: turns noise into an imitation of the data to trick the discriminator

**Discriminator**: tries to identify real from generated data

$X_{real}$

$X_{generated}$

$z$   Generator

Discriminator

$y$

MIT Introduction to Deep Learning (introtodeeplearning.com)

# Generative adversarial networks

## Generating new images



MIT Introduction to Deep Learning (introtodeeplearning.com)
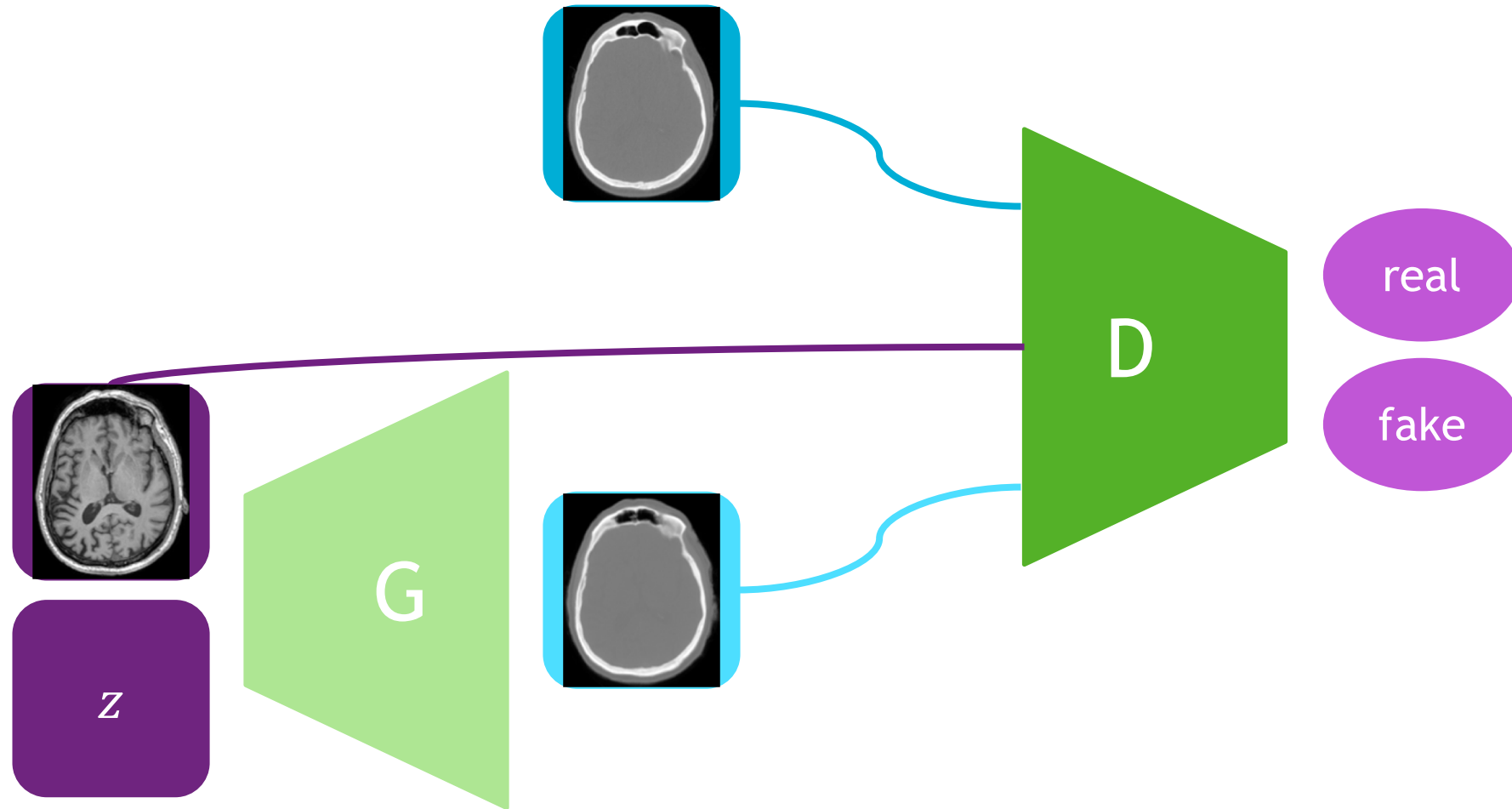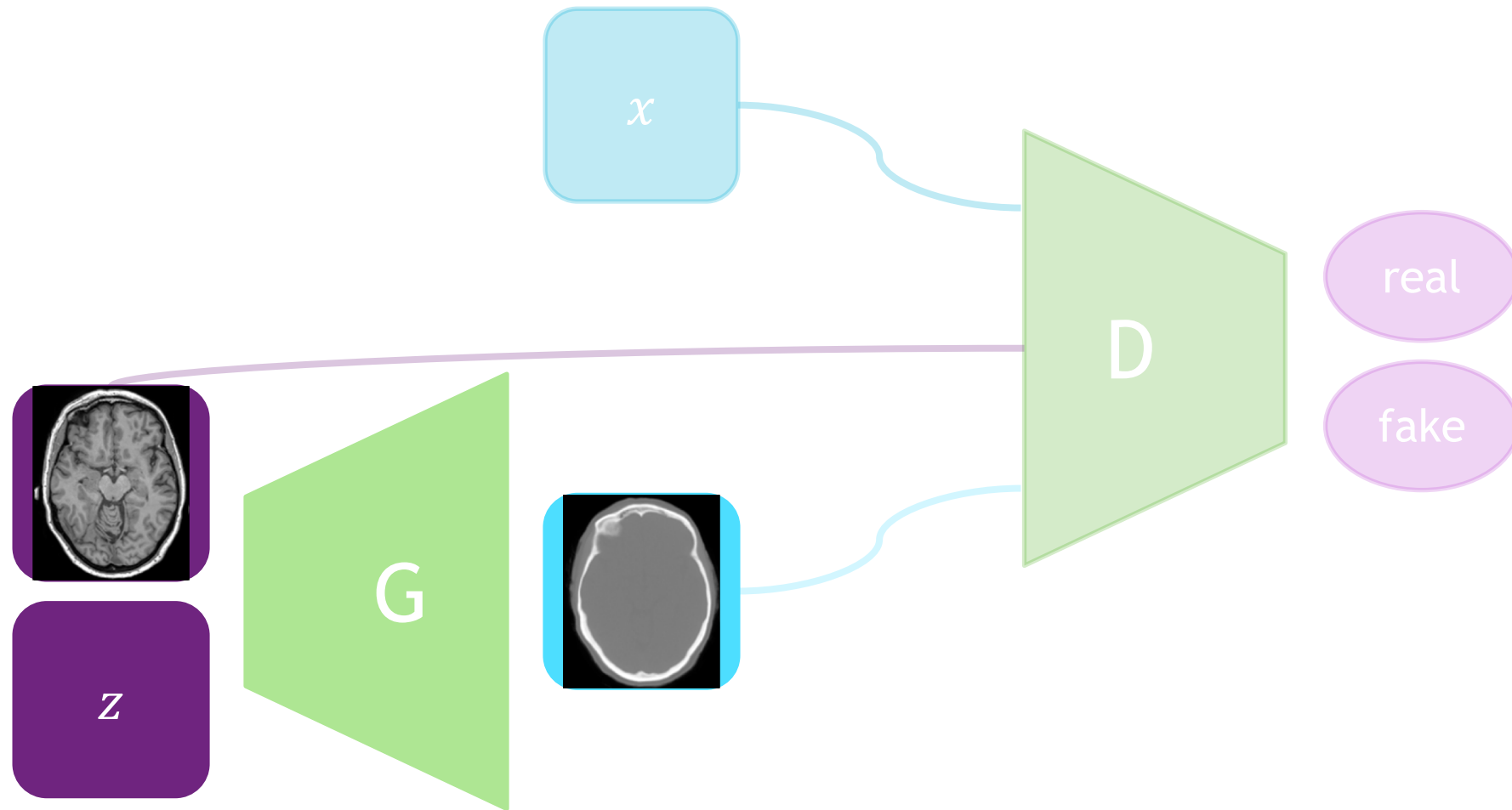
# Generative adversarial networks

## Image translation with conditional GANs

## Image translation with conditional GANs

## Image translation with conditional GANs

# Summary

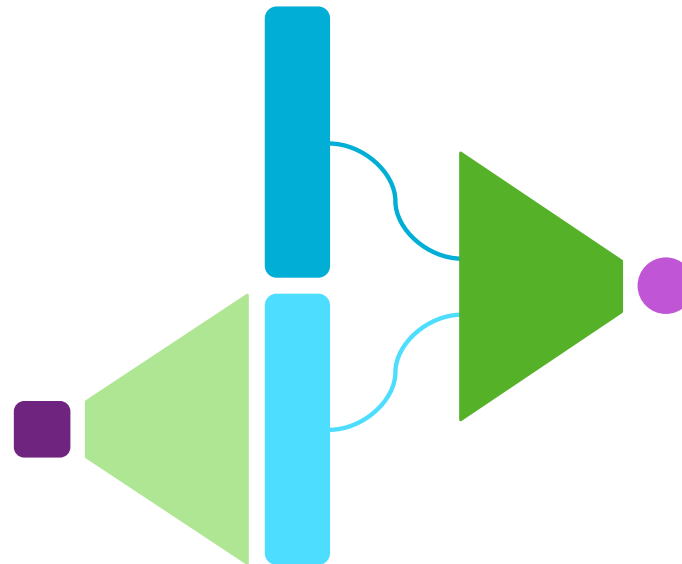| Autoencoders | GANs | Conditional GANs |
|---|---|---|

- Learn low dimensional latent space

- Competing generator and discriminator networks

- For image translation



MIT Introduction to Deep Learning (introtodeeplearning.com)