

# VIBOT 2018

## Machine learning with Neural Networks

**Francesco Ciompi**

[francesco.ciompi@radboudumc.nl](mailto:francesco.ciompi@radboudumc.nl)

# Computational Pathology Research Group

The Computational Pathology Group of the [Department of Pathology](#) of Radboudumc works in close collaboration with DIAG. The history of the group reaches back to the early 1980s, while the current group has been active since 2013. We have a strong translational focus, owing to our strong embedding in a clinical environment. We have a leading position in our research field, witnessed for instance by the highly successful [CAMELYON16](#) and [CAMELYON17](#) grand challenges which we organized and by our active role in committees and conference organisation (e.g. the [Computational Pathology Symposium](#), held in conjunction with the European Congress of Pathology).

## Faculty



Jeroen van der Laak [✉](#)  
+31 24 3614367



Geert Litjens [✉](#)  
+31 243614322



Francesco Ciompi [✉](#)  
+31 24 3614322

## Technicians



Karel Gerbrands [✉](#)  
+31 24 3652283



Rob van de Loo [✉](#)  
+31 24 36 93153

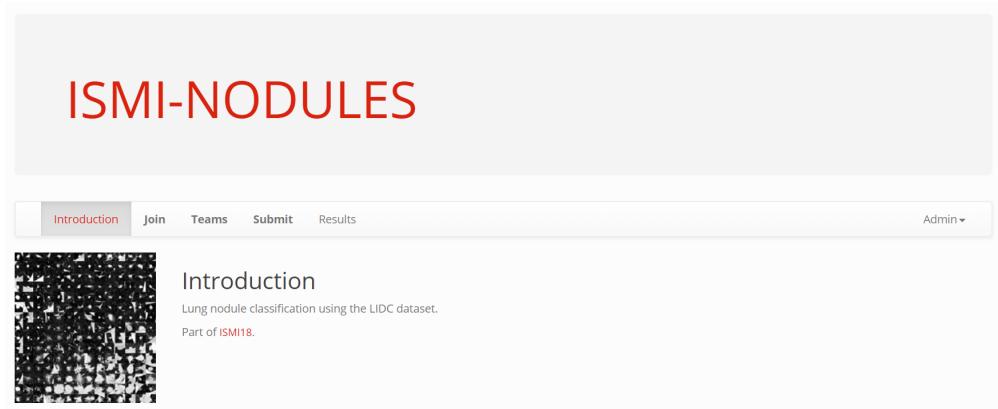
[Home](#)  
[About](#)  
[People](#)  
[Research](#)  
[Publications](#)  
[Thesis Gallery](#)  
[Vacancies](#)  
[Student Projects](#)  
[Contact](#)

# Course overview

	Theory	Workshop
Monday 9 April	10:00 – 12:00	15:00 – 17:00
Tuesday 10 April	12:00 – 14:00	15:00 – 17:00
Wednesday 11 April	--	12:00 – 14:00

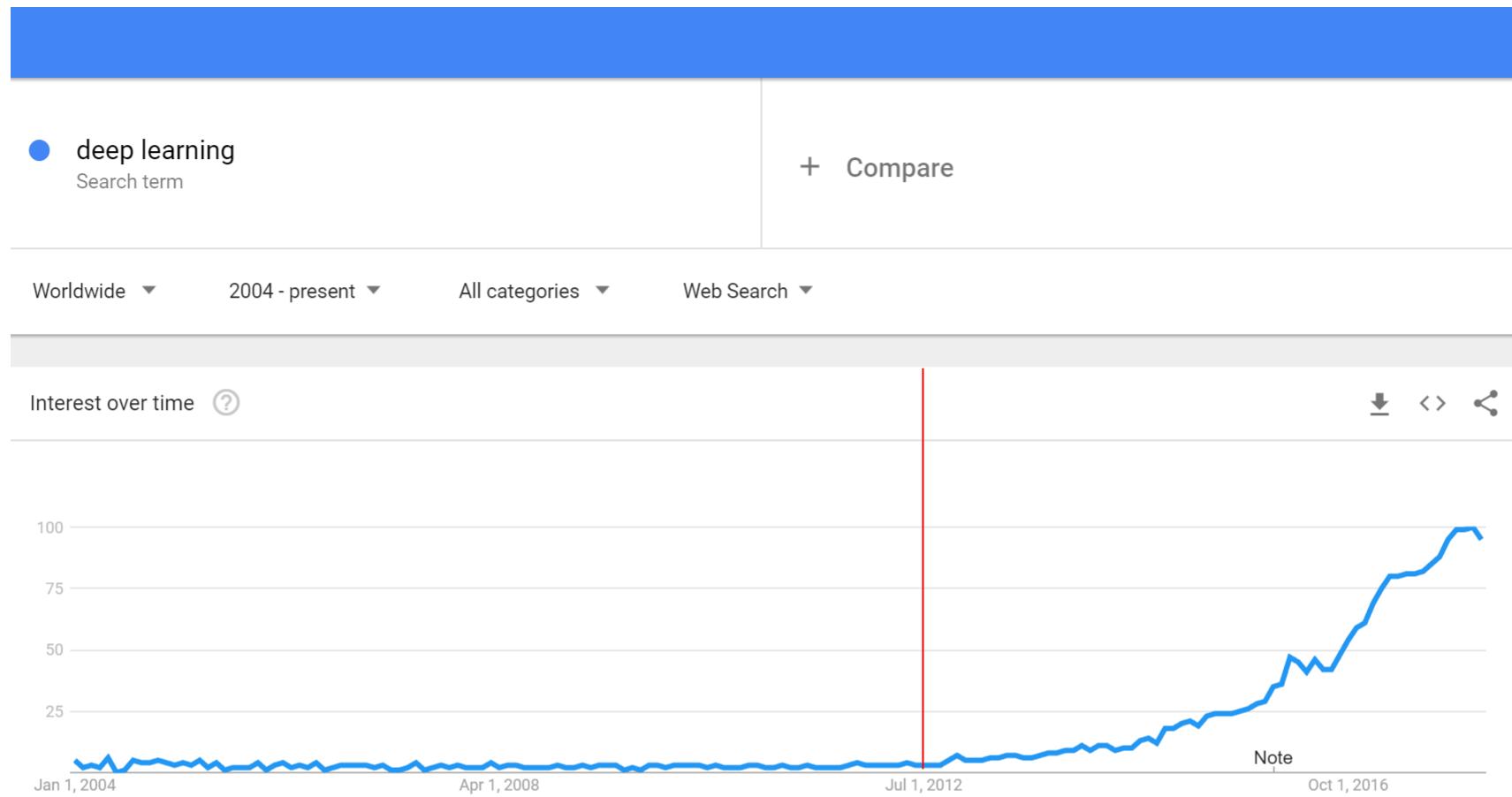
# Course material

- **Lecture slides**
  - Robert sends you the slides
- **Python: Jupyter notebooks**
  - Assignment submission deadline: May 15, 2018
- [www.grand-challenge.org](http://www.grand-challenge.org)

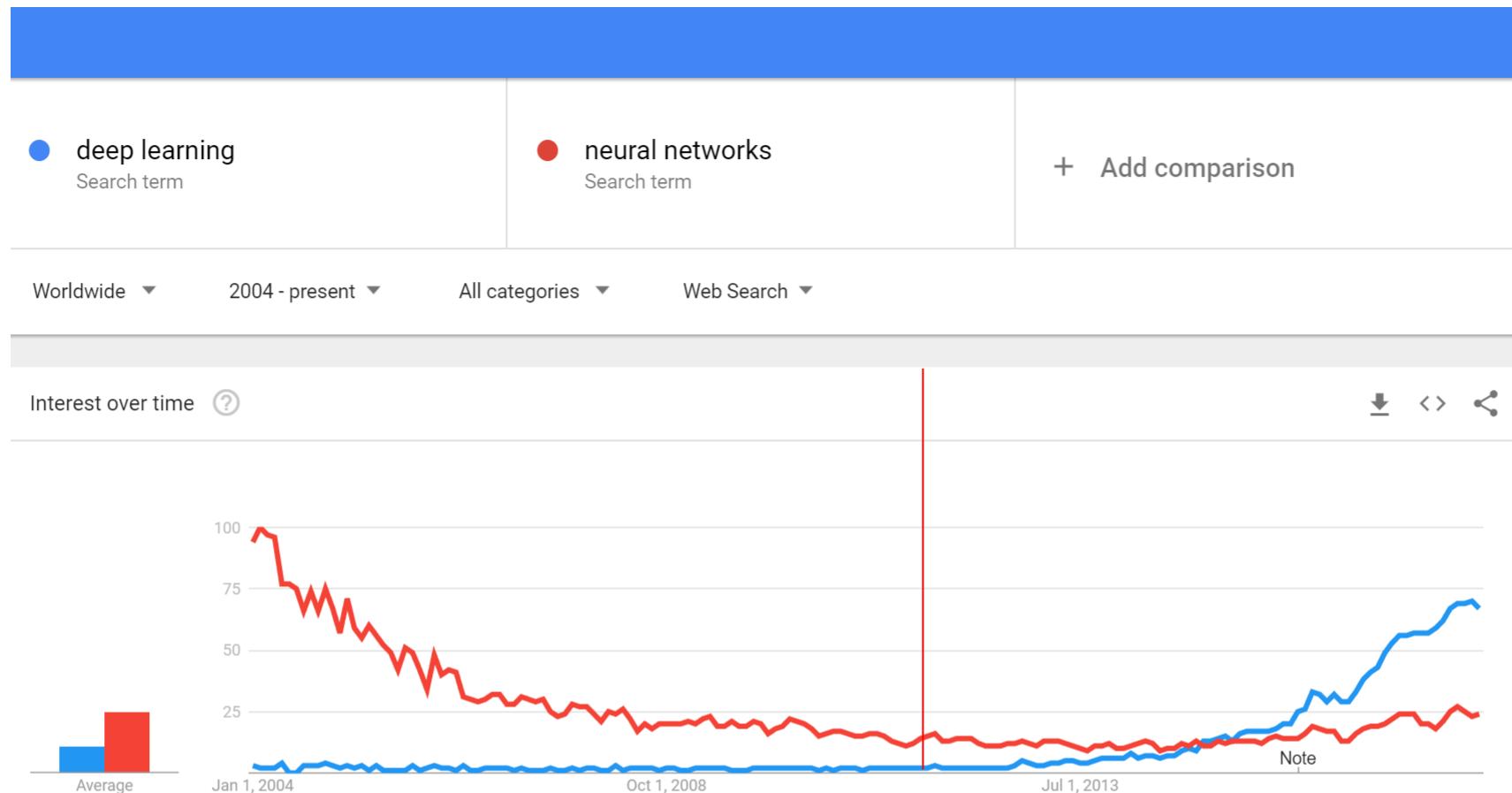


- Access to Linux computer with **3 GPUs**

# What happened in 2012?



# What was happening before 2012?



# Machine Learning

- Our focus will be on:
  - **Supervised** learning
    - “*the truth is given*”
  - **Classification** problems
    - “*we aim at labeling data*”
  - **Discriminative** models
    - “*we aim at finding a decision boundary*”

# Data in supervised learning

- We are typically given  $N$  training samples:

$$(\mathbf{x}_i, y_i), i = 1, \dots, N$$

$\mathbf{x}_i$ : feature vector for  $i^{th}$  sample

$y_i$ : label for  $i^{th}$  sample

# Data in supervised learning

- We are typically given  $N$  training samples:

$$(\mathbf{x}_i, y_i), i = 1, \dots, N$$

$\mathbf{x}_i$ : feature vector for  $i^{th}$  sample

$y_i$ : label for  $i^{th}$  sample

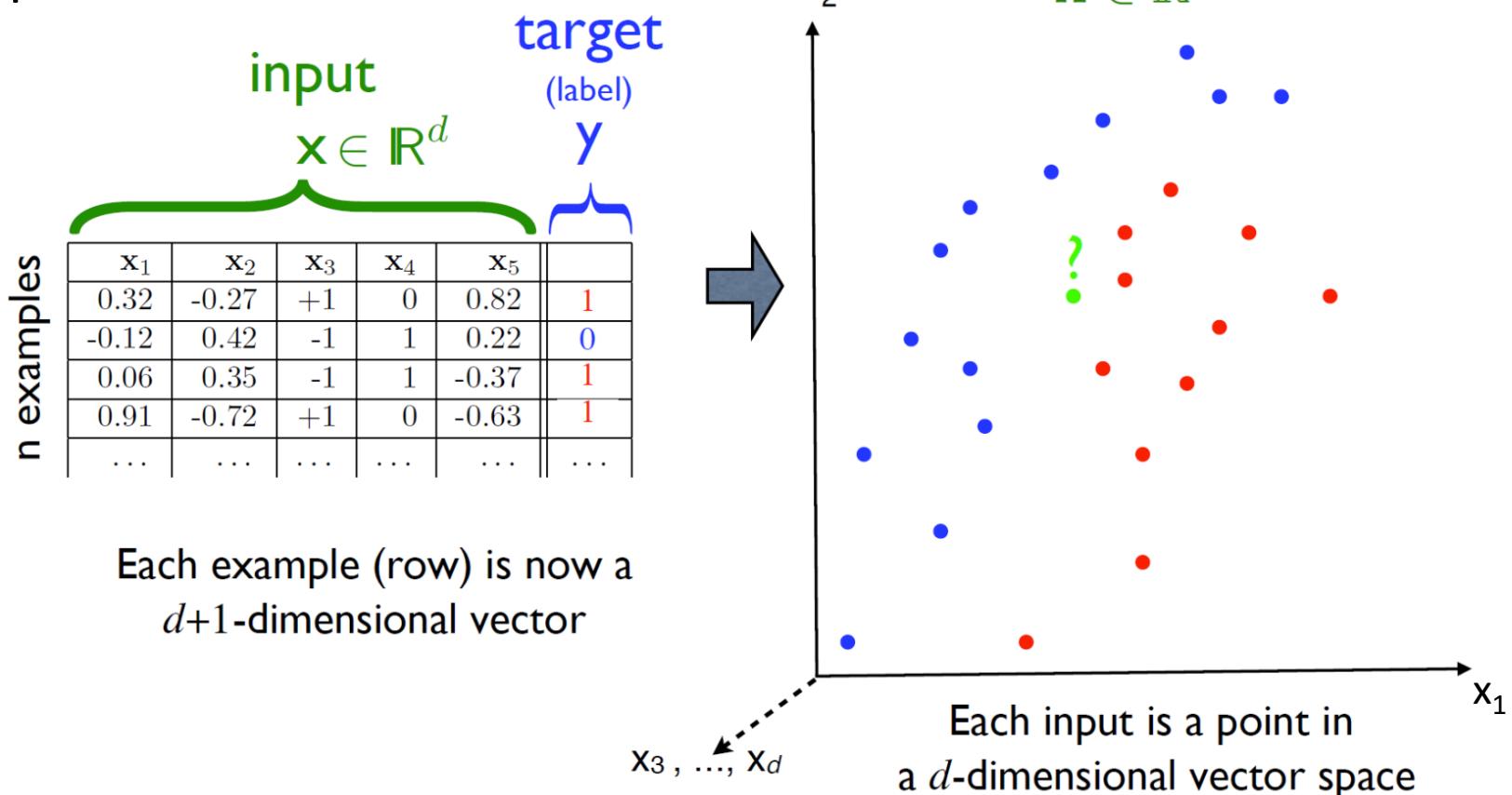
- We want to learn a model (function)

$$f(\mathbf{x}) = y$$

that does a good job at predicting the right label  $y$  for unseen data samples  $\mathbf{x}$

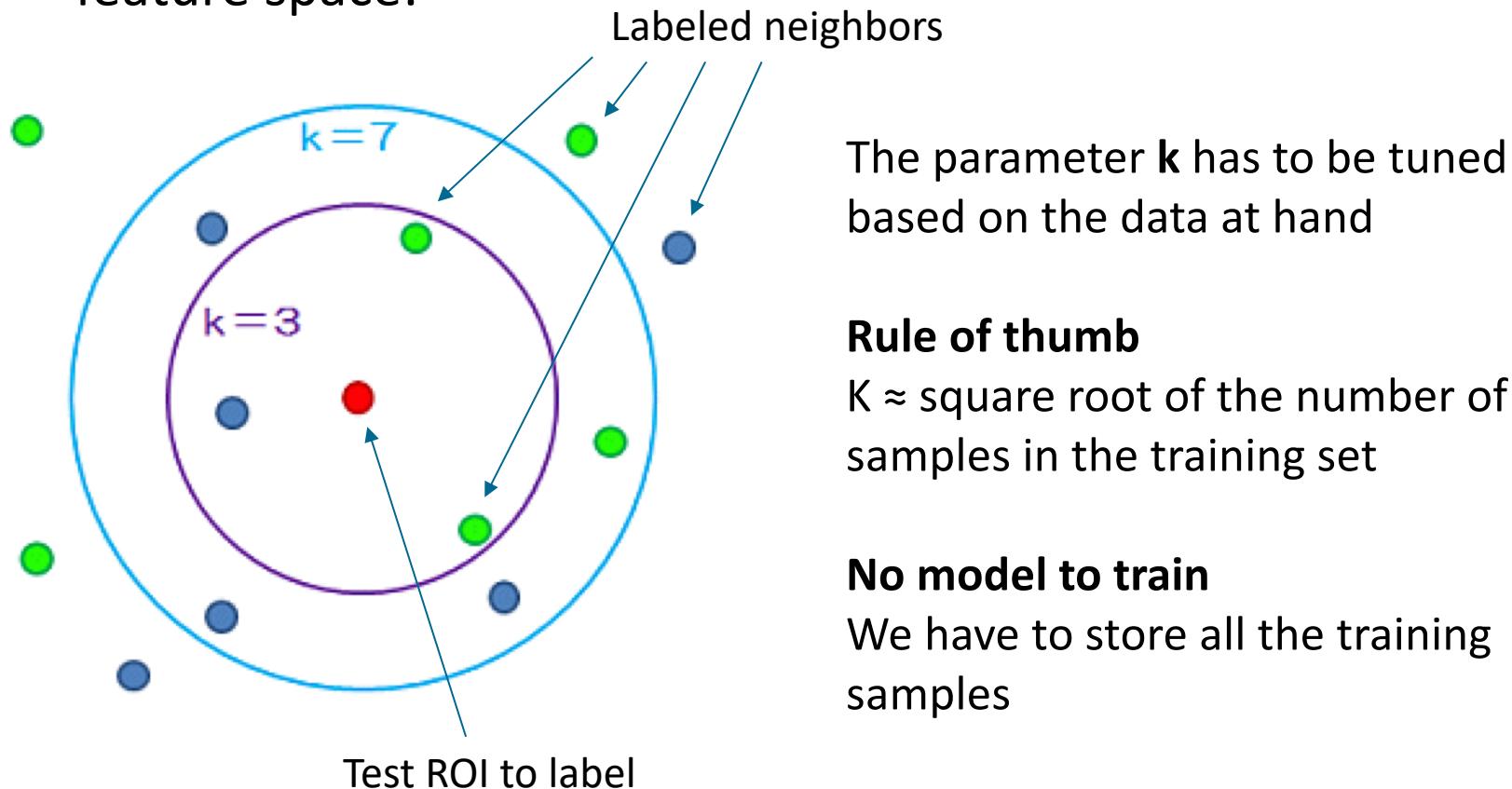
# Data in supervised learning

- Dataset imagined as a point cloud in a high-dimensional vector space



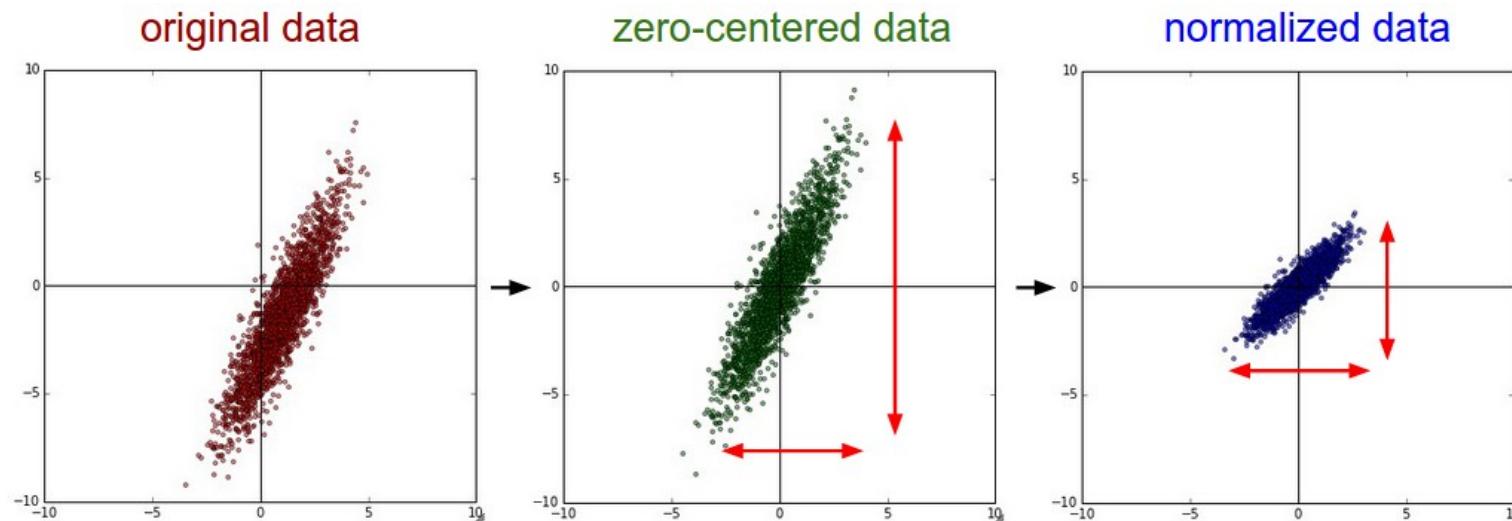
# Simple classification approach

- The label is assigned based on the labels of neighbors in the feature space:



# Feature normalization

- Features can have completely **different ranges**
- If you use **distance-based** algorithms, this can be a problem
- **Normalization** is a good idea (not only when kNN is used)
  - Zero mean, unit standard deviation
- Normalization is done **per feature**



# Supervised learning: how to split data

Labeled data ~ 10,000 samples



Training



grand-challenge.org

Test



80%

~ 8,000 samples

10%

10%

~ 1,000 samples

Training

Test



Validation

# In the Deep Learning era...

Labeled data ~ 1,000,000 samples



Training

Test



98%

~ 980,000 samples

1%

1%

~ 10,000 samples

Training

Test

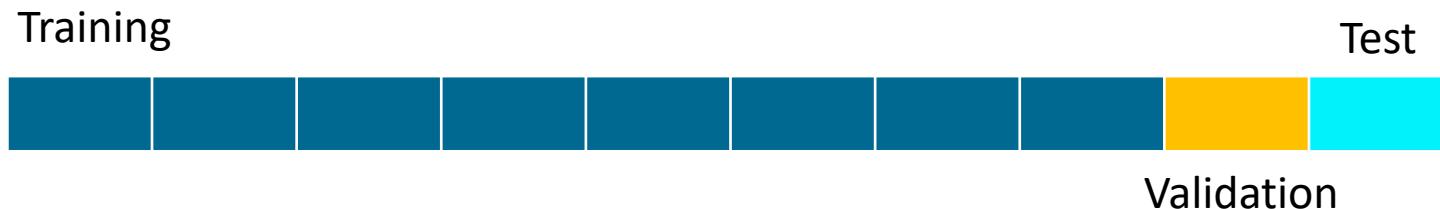


Validation

# Cross validation

*k*-fold cross validation:

- Divide the **labeled dataset** into *k* subsets (e.g., *k*=10)
- Use one subset as the **validation** set
- Use one subset as the **test** set
- Use the remaining *k* – 2 as the **training** set
- Repeat *k* times for all possible combinations



Data tested so far:

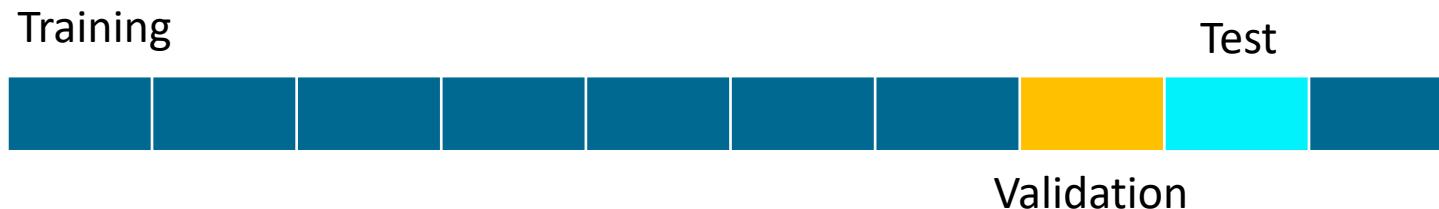


# iteration: 1

# Cross validation

*k*-fold cross validation:

- Divide the **labeled dataset** into *k* subsets (e.g., *k*=10)
- Use one subset as the **validation** set
- Use one subset as the **test** set
- Use the remaining *k* – 2 as the **training** set
- Repeat *k* times for all possible combinations



Data tested so far:



# iteration: 2

# Cross validation

*k*-fold cross validation:

- Divide the **labeled dataset** into *k* subsets (e.g., *k*=10)
- Use one subset as the **validation** set
- Use one subset as the **test** set
- Use the remaining *k* – 2 as the **training** set
- Repeat *k* times for all possible combinations



Data tested so far:



# iteration: 3

# Cross validation

*k*-fold cross validation:

- Divide the **labeled dataset** into *k* subsets (e.g., *k*=10)
- Use one subset as the **validation** set
- Use one subset as the **test** set
- Use the remaining *k* – 2 as the **training** set
- Repeat *k* times for all possible combinations



Data tested so far:

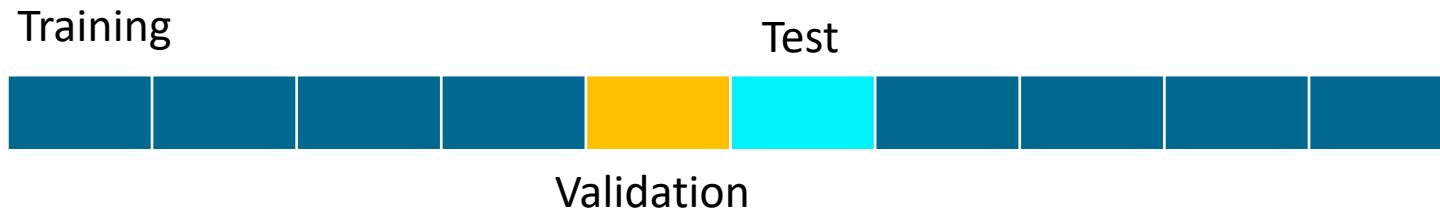


# iteration: 4

# Cross validation

*k*-fold cross validation:

- Divide the **labeled dataset** into *k* subsets (e.g., *k*=10)
- Use one subset as the **validation** set
- Use one subset as the **test** set
- Use the remaining *k* – 2 as the **training** set
- Repeat *k* times for all possible combinations



Data tested so far:

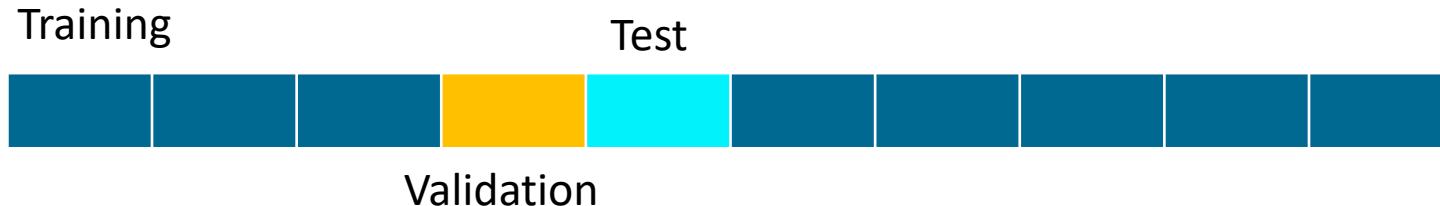


# iteration: 5

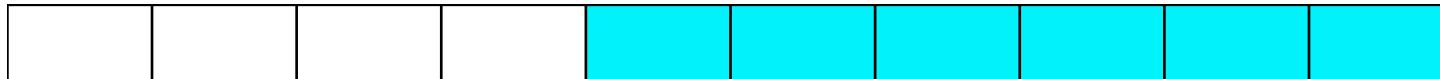
# Cross validation

*k*-fold cross validation:

- Divide the **labeled dataset** into *k* subsets (e.g., *k*=10)
- Use one subset as the **validation** set
- Use one subset as the **test** set
- Use the remaining *k* – 2 as the **training** set
- Repeat *k* times for all possible combinations



Data tested so far:

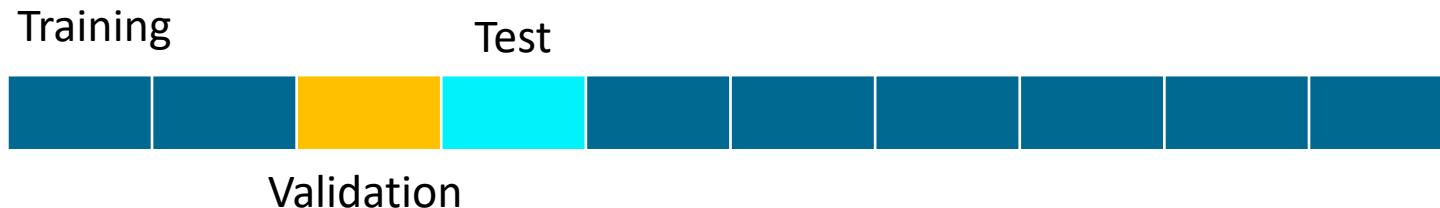


# iteration: 6

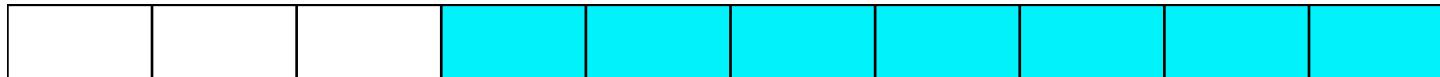
# Cross validation

*k*-fold cross validation:

- Divide the **labeled dataset** into *k* subsets (e.g., *k*=10)
- Use one subset as the **validation** set
- Use one subset as the **test** set
- Use the remaining *k* – 2 as the **training** set
- Repeat *k* times for all possible combinations



Data tested so far:

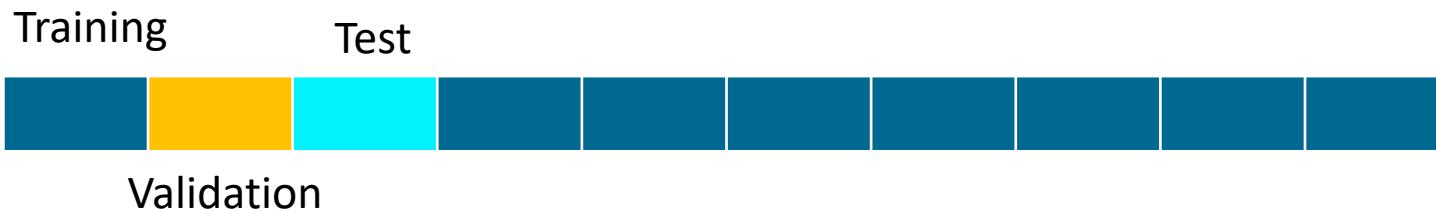


# iteration: 7

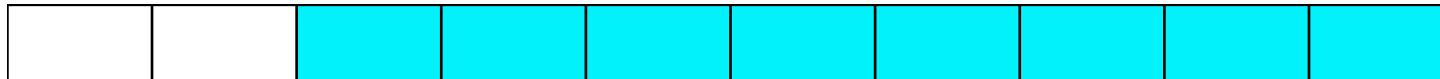
# Cross validation

*k*-fold cross validation:

- Divide the **labeled dataset** into *k* subsets (e.g., *k*=10)
- Use one subset as the **validation** set
- Use one subset as the **test** set
- Use the remaining *k* – 2 as the **training** set
- Repeat *k* times for all possible combinations



Data tested so far:



# iteration: 8

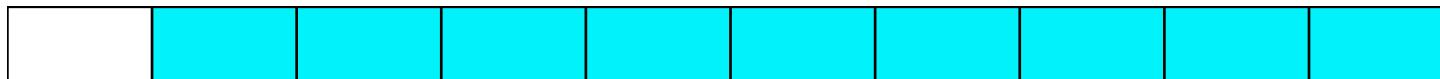
# Cross validation

*k*-fold cross validation:

- Divide the **labeled dataset** into *k* subsets (e.g., *k*=10)
- Use one subset as the **validation** set
- Use one subset as the **test** set
- Use the remaining *k* – 2 as the **training** set
- Repeat *k* times for all possible combinations



Data tested so far:



# iteration: 9

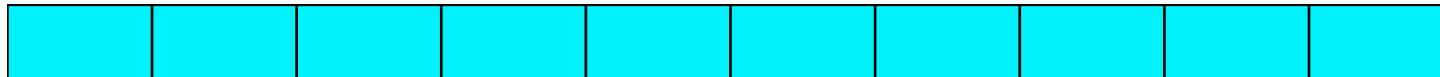
# Cross validation

*k*-fold cross validation:

- Divide the **labeled dataset** into *k* subsets (e.g., *k*=10)
- Use one subset as the **validation** set
- Use one subset as the **test** set
- Use the remaining *k* – 2 as the **training** set
- Repeat *k* times for all possible combinations



Data tested so far:

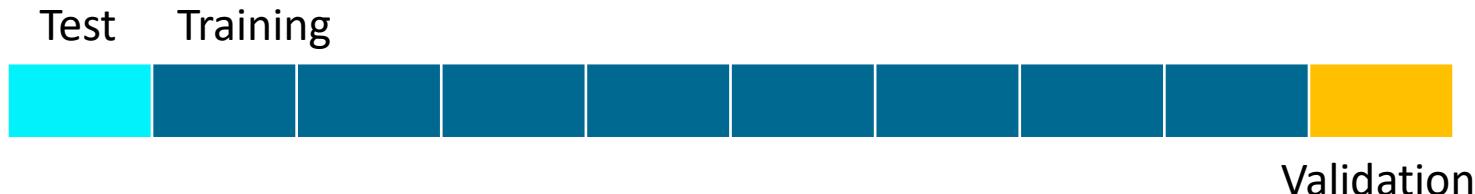


# iteration: 10

# Cross validation

*k*-fold cross validation:

- Divide the **labeled dataset** into *k* subsets (e.g., *k*=10)
- Use one subset as the **validation** set
- Use one subset as the **test** set
- Use the remaining *k* – 2 as the **training** set
- Repeat *k* times for all possible combinations



- Called *leave one-out* if  $k=|\text{trainingData}|$

In medical imaging, do not mix **patients** in training/ validation / test

**NO INTERSECTION OF PATIENT DATA ACROSS FOLDS !**

---

# Parametric models

# Parametric model

- The result of our **supervised learning** procedure is a **function**  $f(x)$  that **predicts** the label  $y$  for a given input  $x$
- Typically, the model we are building has some parameters  $W$

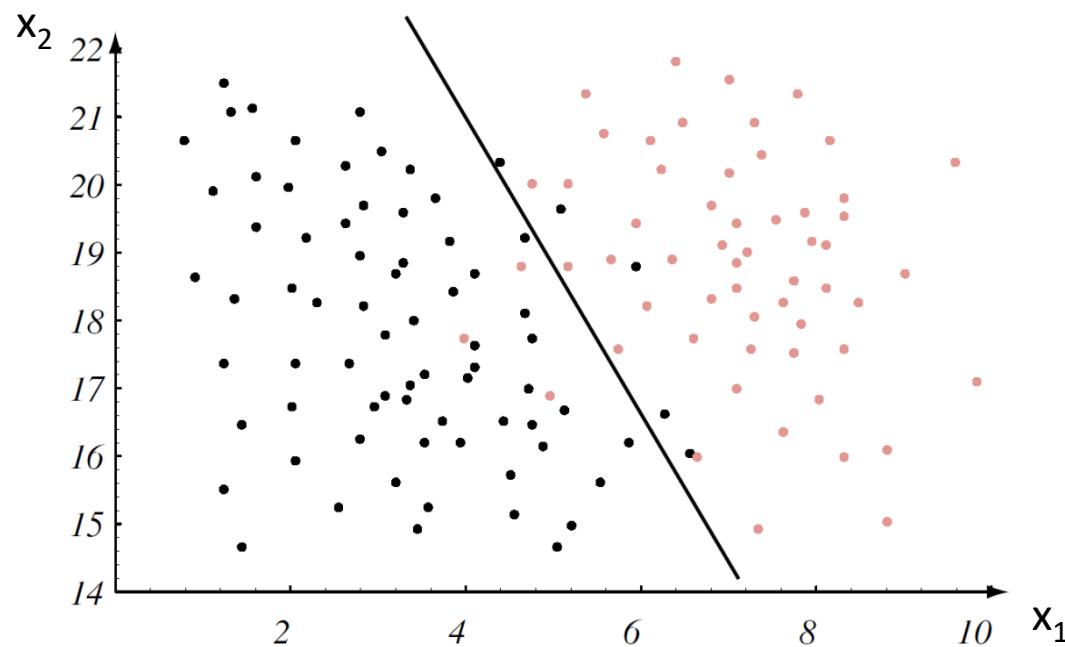
$$y = f(\textcolor{teal}{x}, \textcolor{blue}{W})$$

Data                      Parameters  
(image, feature vector)               $\textcolor{teal}{x}$                $\textcolor{blue}{W}$

- The amount of parameters of the model defines its **capacity**

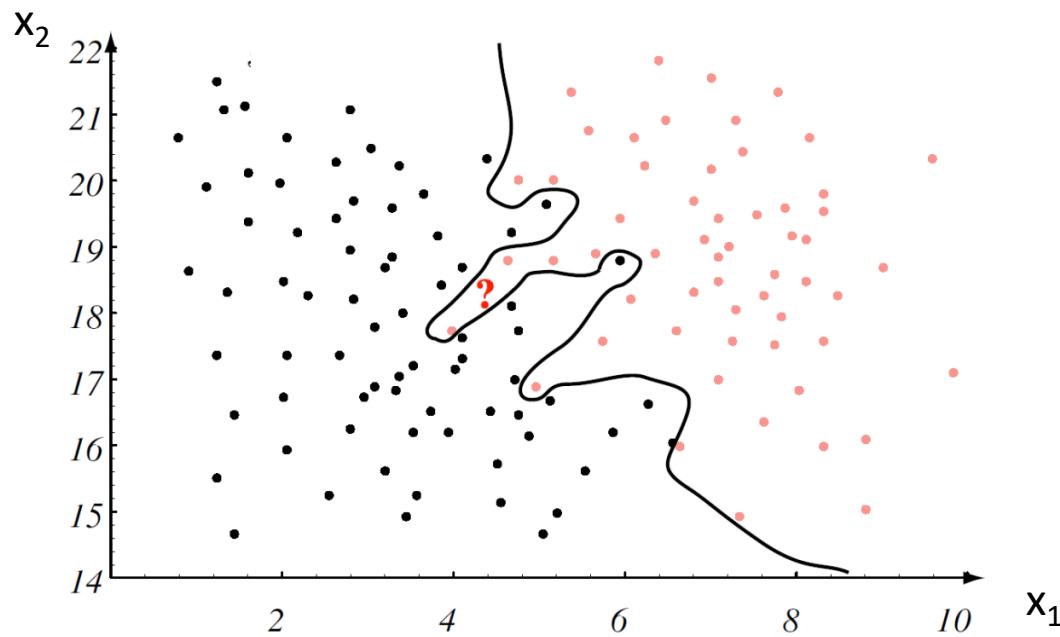
# Capacity of a model

- Function family too poor (too inflexible)
- Capacity too low
- Under-fitting



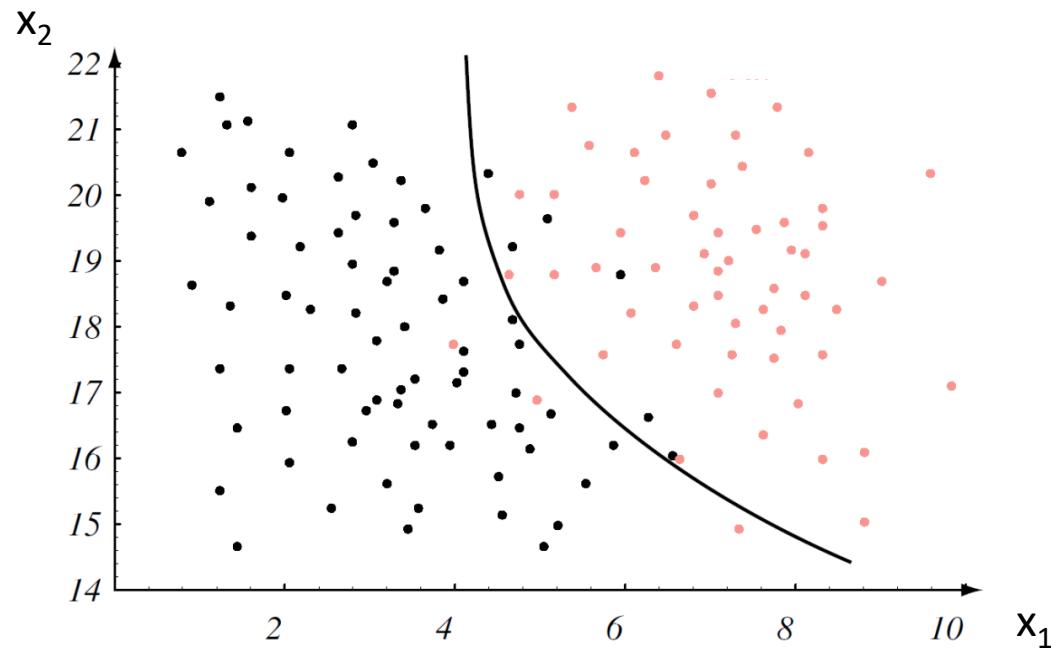
# Capacity of a model

- Function family too rich (too flexible)
- Capacity too high (for this problem, relative to the number of examples)
- Over-fitting

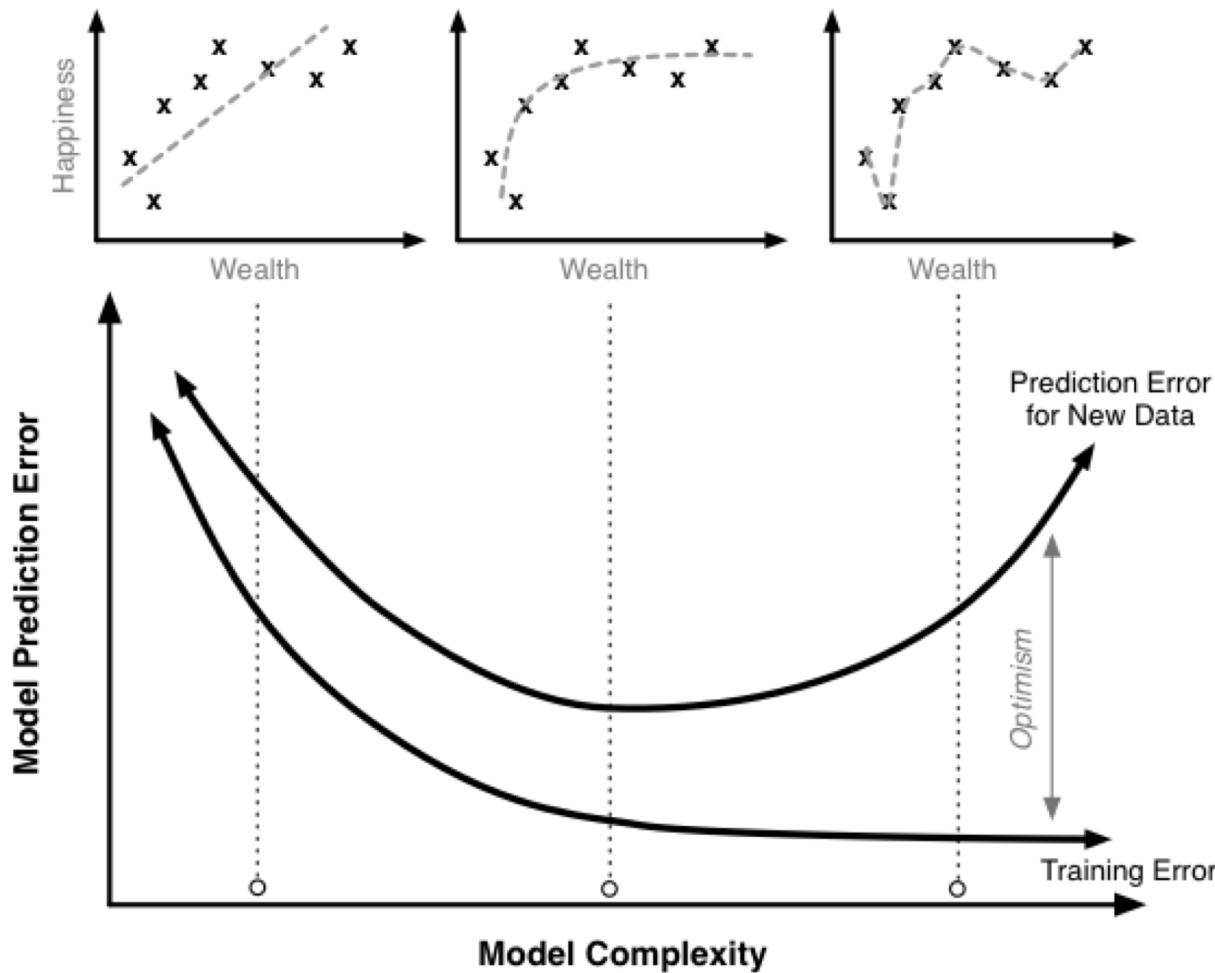


# Capacity of a model

- Optimal capacity (for this problem)
- **Best generalization** (on future test points)



# Capacity of a model



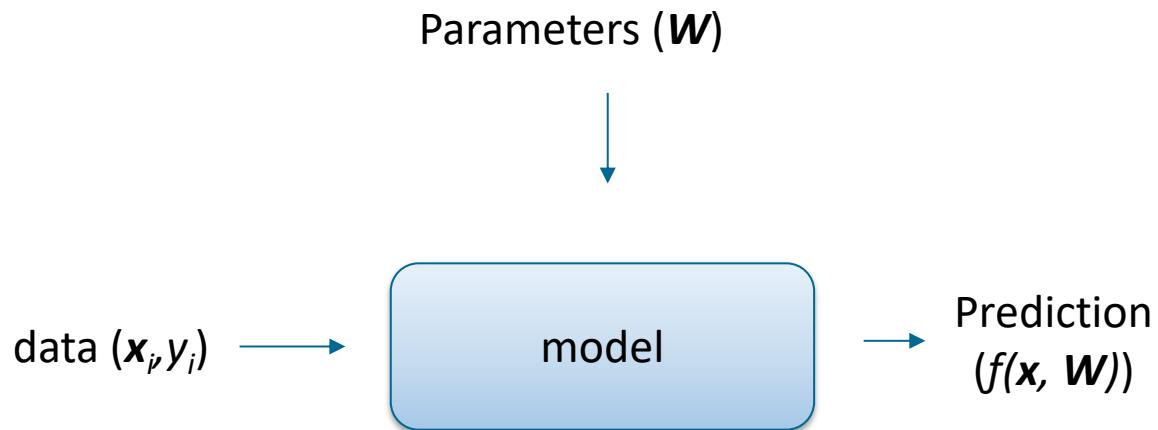
---

# Training a parametric model

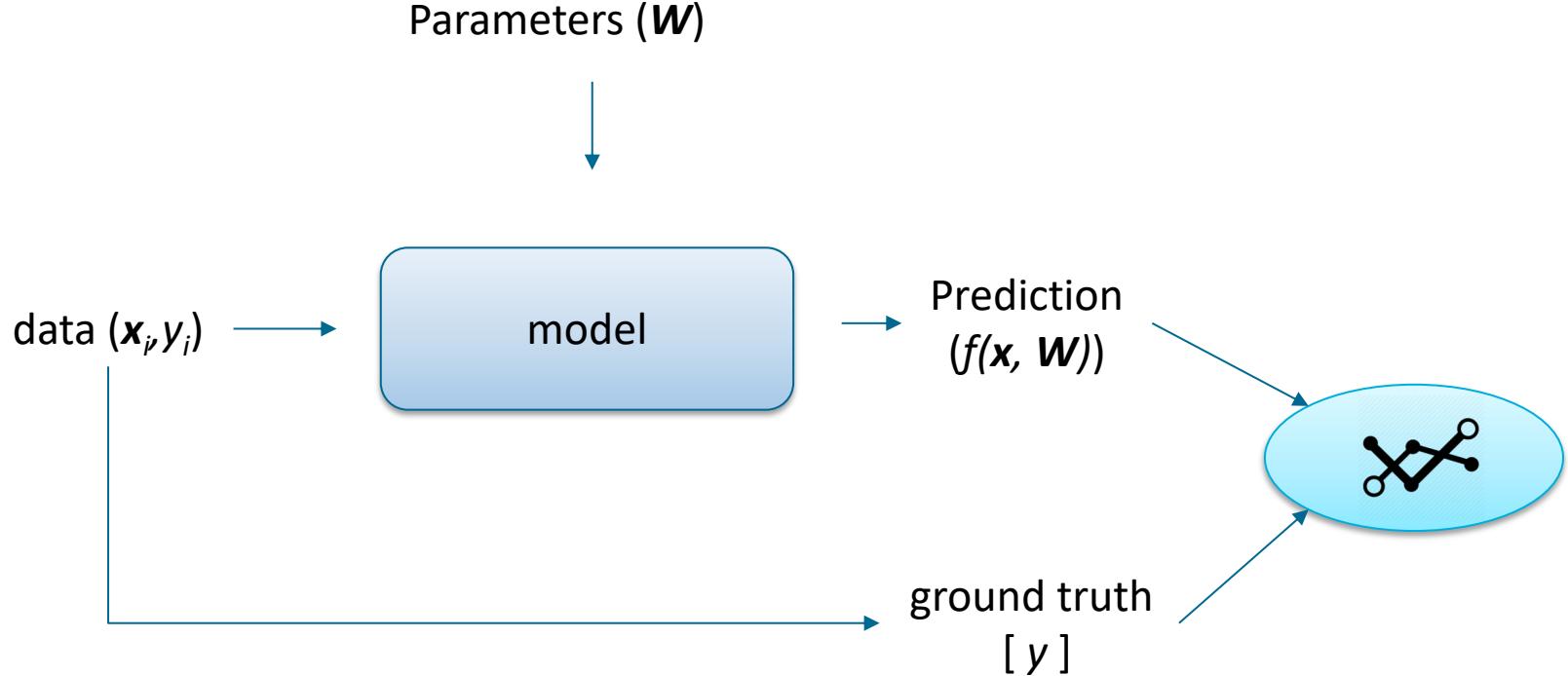
# Supervised learning

- The goal of a supervised learning algorithm (**and ours**):
  1. Learn from training data (achieve low training error)
  2. Generalize to unseen data (achieve low validation error)
  3. Avoid overfitting (when both condition (1) and (2) are achieved)
- Let's start from the first question:
  - *“How can we learn from data?”*

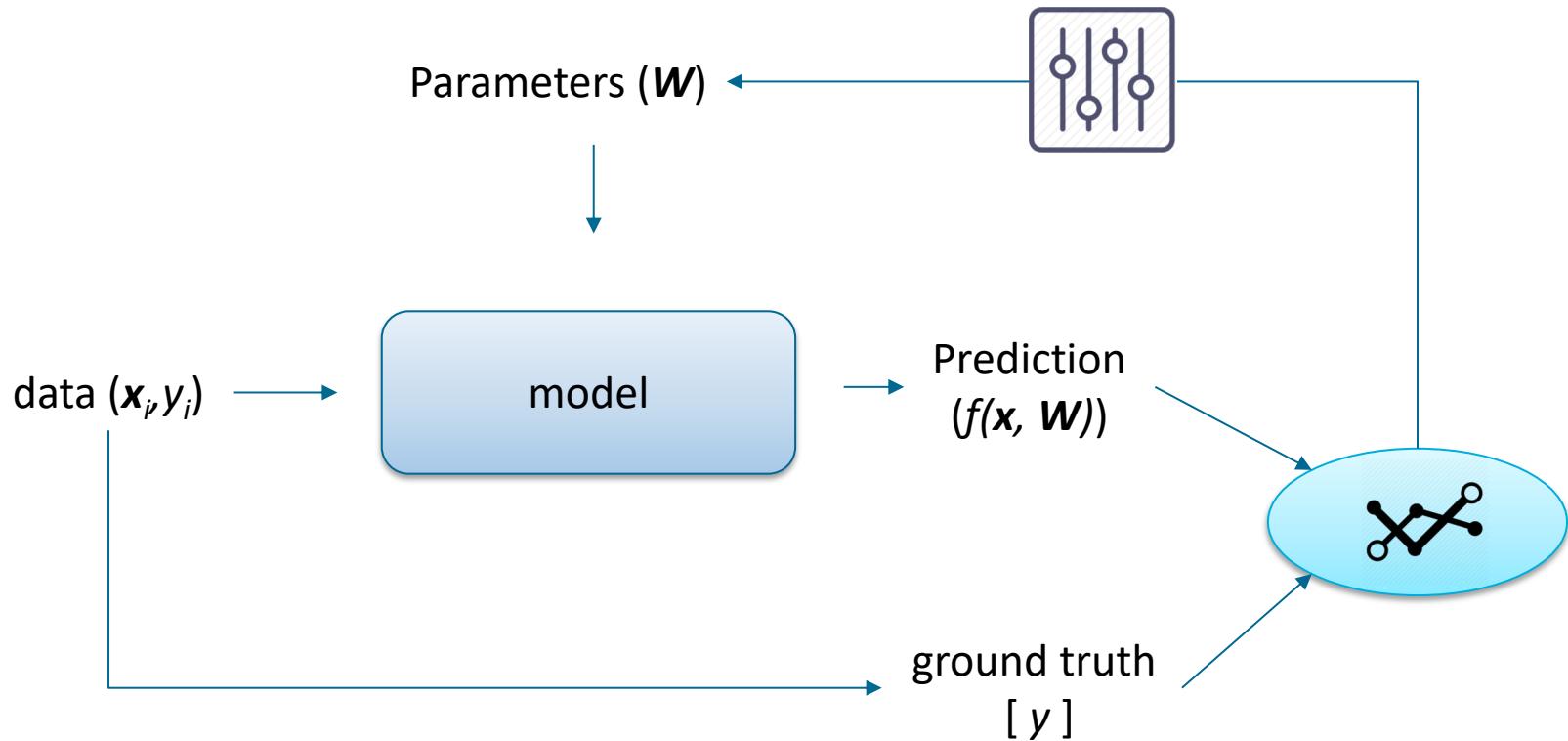
# Supervised learning procedure



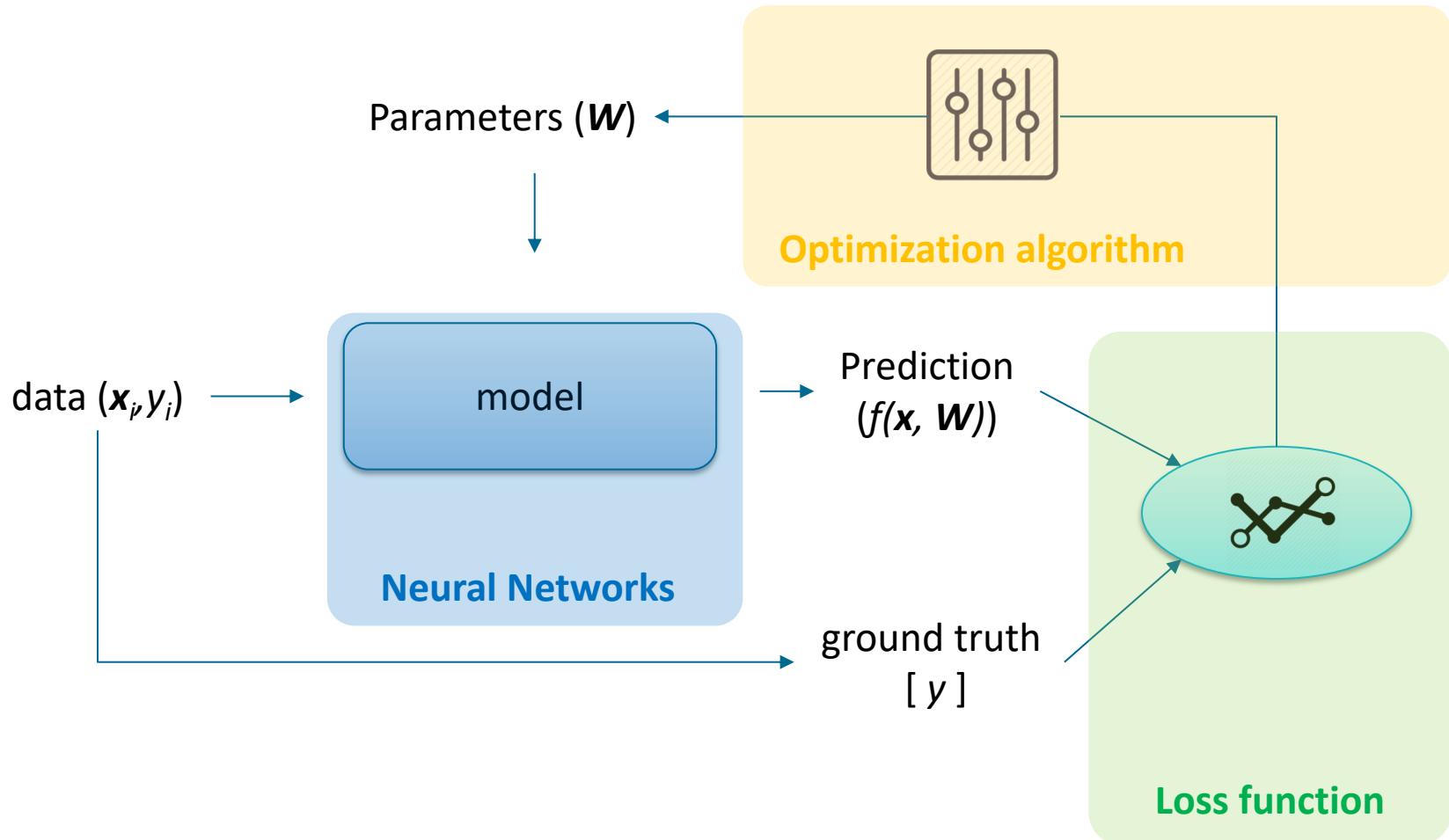
# Supervised learning procedure



# Supervised learning procedure

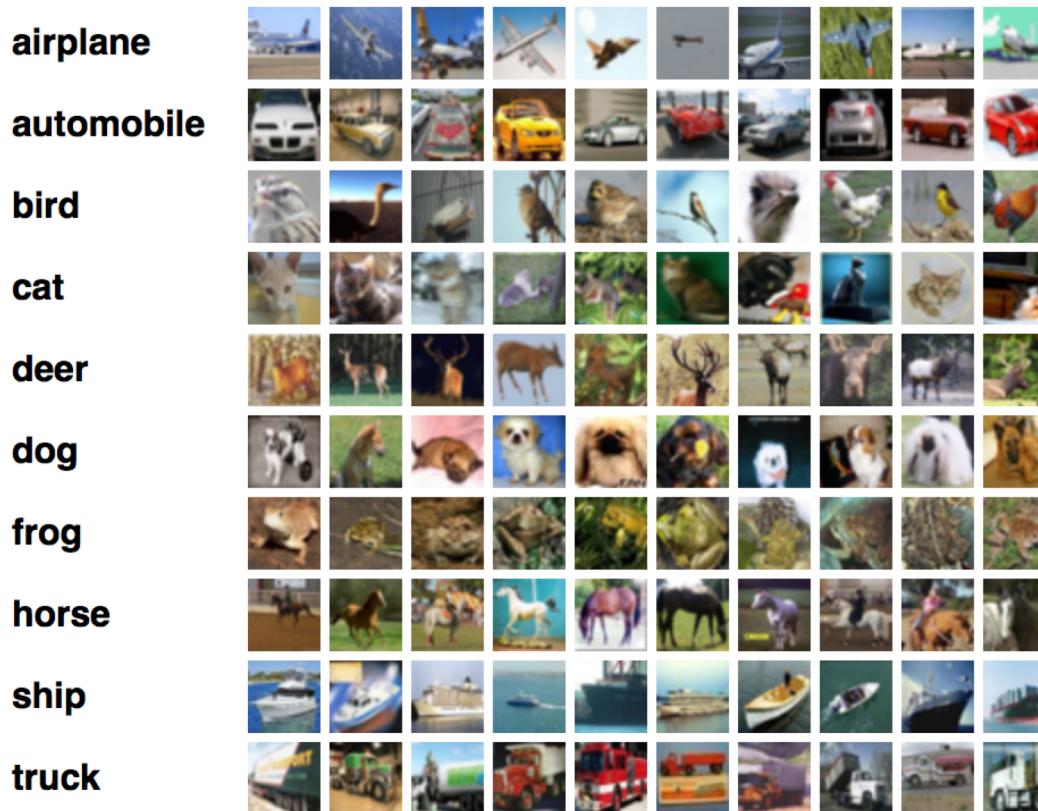


# Supervised learning procedure



# A simple model: linear classifier

- Suppose we want to classify images in the CIFAR10 dataset



# A simple model: linear classifier

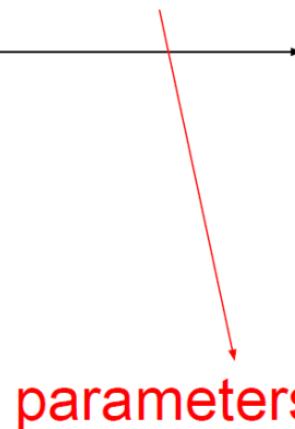


[32x32x3]  
array of numbers 0...1

$$f(x, W) = \boxed{W} \boxed{x} \quad 3072 \times 1$$

**10x1**      **10x3072**

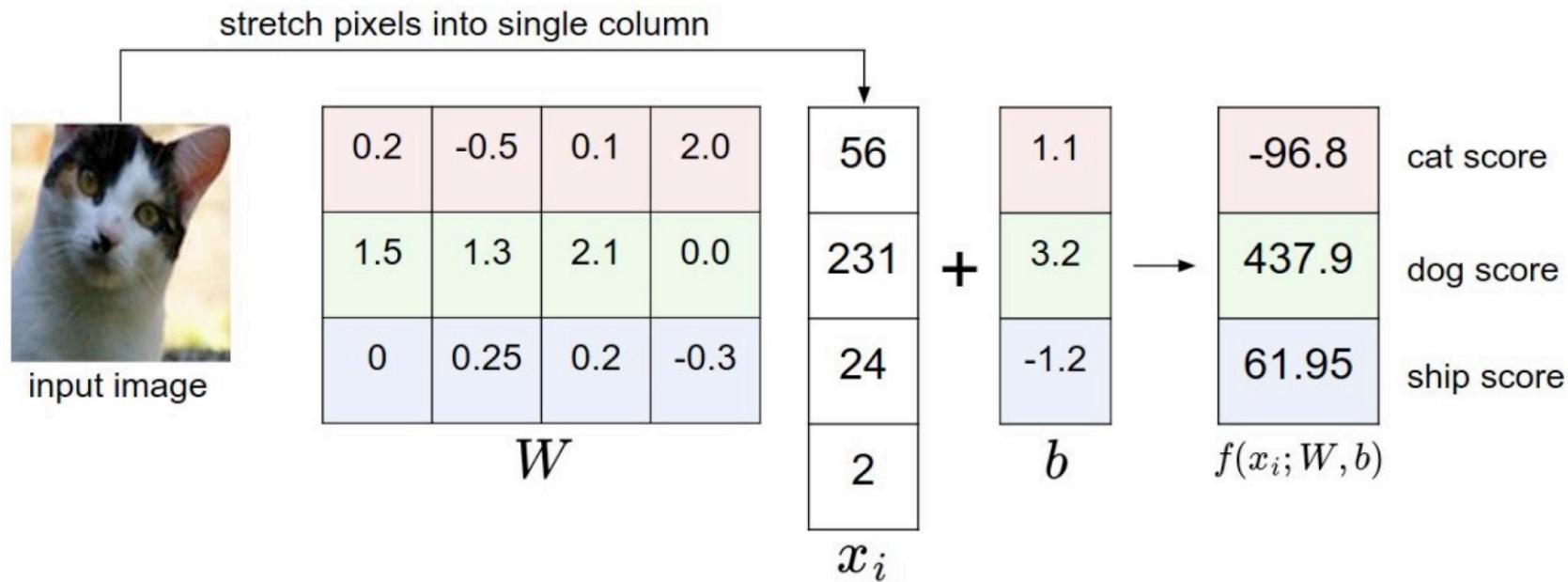
A diagram showing the multiplication of the input vector  $x$  (10x1) and the weight matrix  $W$  (10x3072). A horizontal arrow points from the input vector to the result, which is a 10x1 vector of class scores.



10 numbers,  
indicating class  
scores

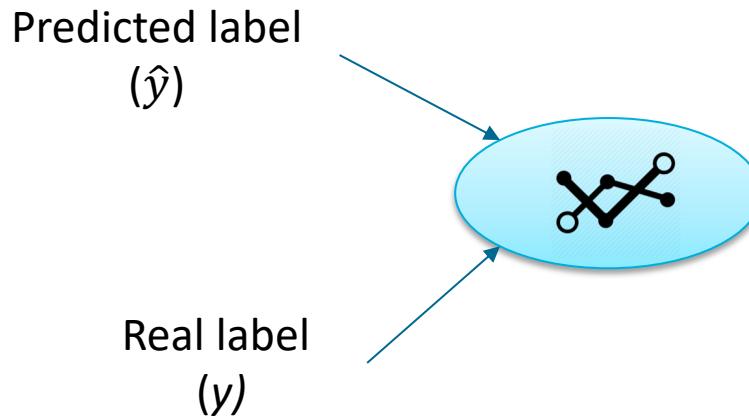
# A simple model: linear classifier

- We simplify the example: image with 4 pixels, 3 classes (**cat/dog/ship**)



# Loss function

- The loss function measures the difference between the **predicted label** and the **real label**



- During training, we want to **minimize the loss**
- It can be an **arbitrary** function, but **differentiable!!!**

# Loss function

The screenshot shows the Keras Documentation website. The left sidebar contains a navigation menu with categories like Home, Why use Keras, Getting started, Guide to the Sequential model, Guide to the Functional API, FAQ, Models, Layers, and various sub-sections for each. The main content area is titled "Loss functions". It includes sections for "Arguments" and "Returns". Below these, there is a list of loss functions, each with its name in bold and a code snippet in a code block:

- categorical\_crossentropy**  
`categorical_crossentropy(y_true, y_pred)`
- sparse\_categorical\_crossentropy**  
`sparse_categorical_crossentropy(y_true, y_pred)`
- binary\_crossentropy**  
`binary_crossentropy(y_true, y_pred)`
- kullback\_leibler\_divergence**  
`kullback_leibler_divergence(y_true, y_pred)`
- poisson**  
`poisson(y_true, y_pred)`
- cosine\_proximity**  
`cosine_proximity(y_true, y_pred)`

A green oval highlights the first three items in the list: categorical\_crossentropy, sparse\_categorical\_crossentropy, and binary\_crossentropy.

# Loss function

- **Binary cross-entropy**

- Real label:  $y_i \in [0,1]$
- Predicted label:  $\hat{y}_i \in [0,1]$

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

# Loss function

- **Binary cross-entropy**

- Real label:  $y_i \in [0,1]$
- Predicted label:  $\hat{y}_i \in [0,1]$

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

## Intuition

- If  $y = 1$ : your loss function will force  $\hat{y}$  to be **large**
- If  $y = 0$ : your loss function will force  $\hat{y}$  to be **small**

# Loss function

- **Categorical cross-entropy loss**

- Real label:  $y \in [0,1]$
- Predicted label:  $\hat{y} \in [0,1]$

$$L_i = - \sum_j y_{i,j} \log(\hat{y}_{i,j})$$

# Loss function

- **Categorical cross-entropy loss**

- Real label:  $y \in [0,1]$
- Predicted label:  $\hat{y} \in [0,1]$

$$L_i = - \sum_j y_{i,j} \log(\hat{y}_{i,j})$$

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

- Extension of the binary case to multi-class

# One-hot representation

- **Categorical cross-entropy for multi-class problems**
- How can we have multiple labels with only  $y \in [0,1]$ ?

# One-hot representation

- Categorical cross-entropy for multi-class problems
- How can we have multiple labels with only  $y \in [0,1]$ ?
- Suppose we have 4 classes, for each sample, the format of the target will be actually a vector:
  - If  $y = 1$ 

1	0	0	0
---	---	---	---
  - If  $y = 2$ 

0	1	0	0
---	---	---	---
  - If  $y = 3$ 

0	0	1	0
---	---	---	---
  - If  $y = 4$ 

0	0	0	1
---	---	---	---

# One-hot representation

- Categorical cross-entropy for multi-class problems
- How can we have multiple labels with only  $y \in [0,1]$ ?
- Suppose we have 4 classes, for each sample, the format of the target will be actually a vector:
  - If  $y = 1$ 

1	0	0	0
---	---	---	---
  - If  $y = 2$ 

0	1	0	0
---	---	---	---
  - If  $y = 3$ 

0	0	1	0
---	---	---	---
  - If  $y = 4$ 

0	0	0	1
---	---	---	---

One-hot  
representation



# Cost function

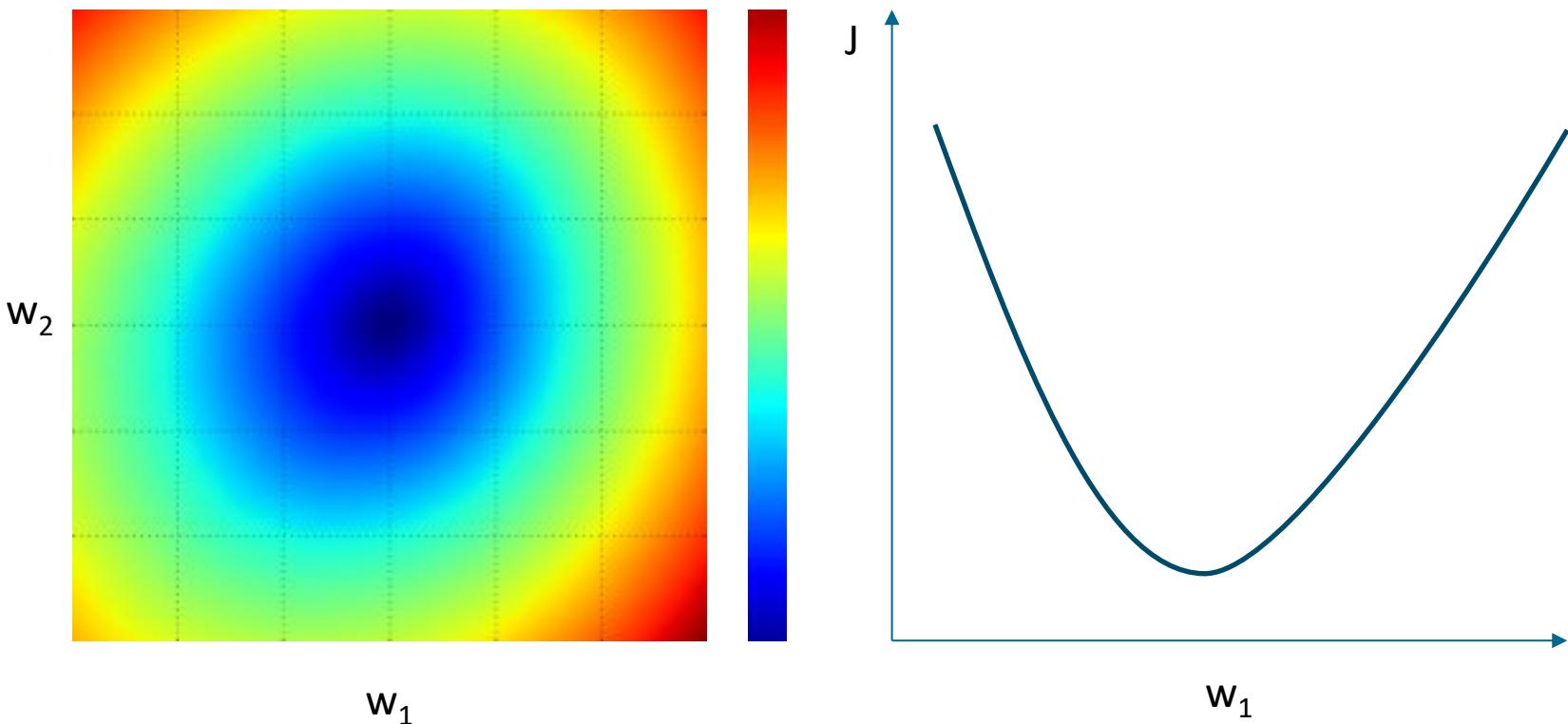
- Given a training set  $(x_i, y_i), i = 1, \dots, N$ :
  - We can define a **loss** function  $L_i$  for each training sample
  - We can define a **cost** function  $J$  for the entire training set

$$J = \frac{1}{N} \sum_{i=1}^N L_i$$

- The final cost function is the **average across the training set**

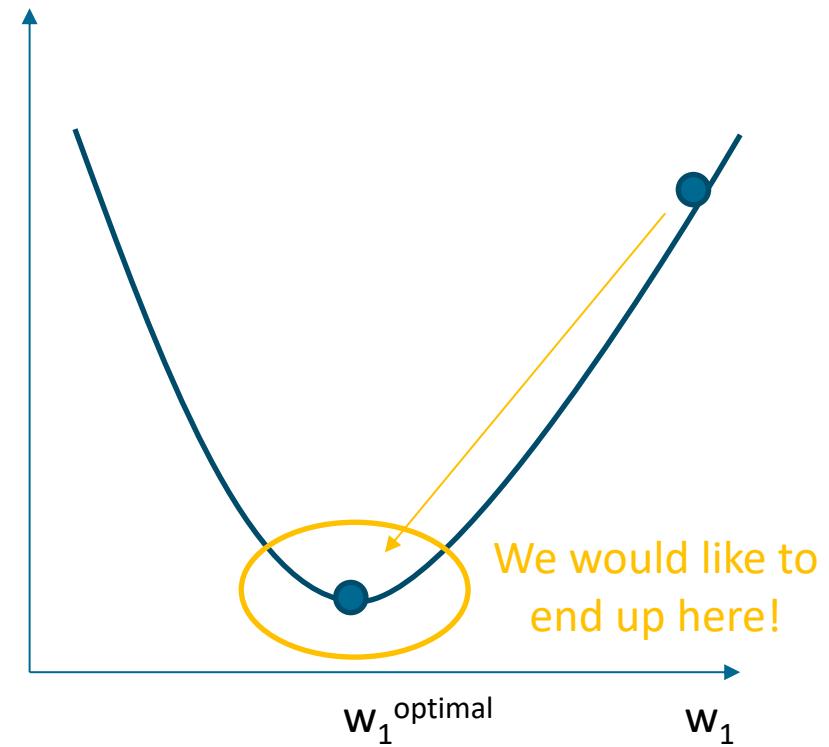
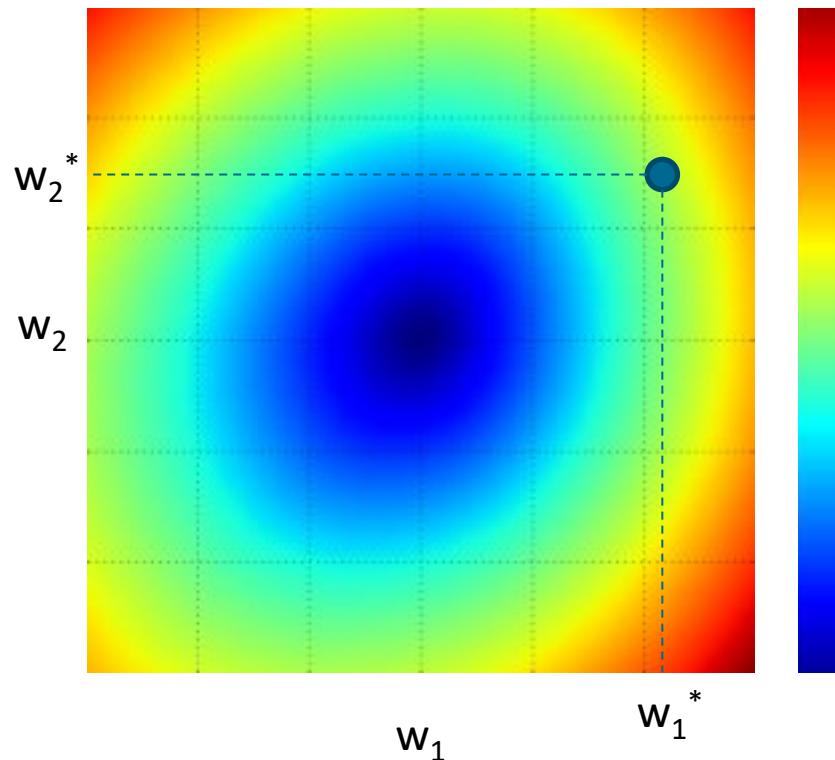
# Cost function

- The shape of the cost function could be something like this:



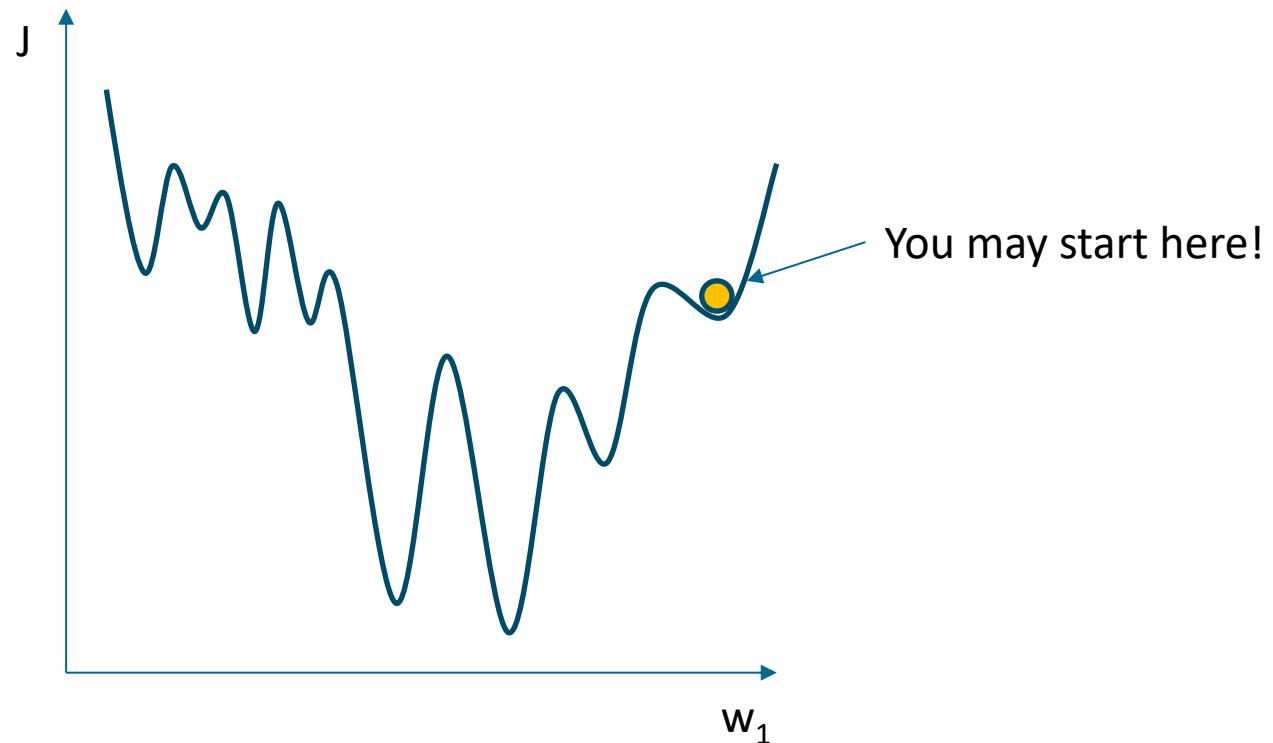
# Cost function

- When we initialize our model (typically with random parameters), we are in a random position of this bowl:



# Cost function

- How can we get to the minimum?
- Consider that the typical profile of cost function is like this:



# Gradient

- Given the initial position, we have to find the **best direction** along which we should step to change our weights
- The best direction is given by the **gradient** of the cost function
- Computing the gradient is like computing the “**slope**” of the function given a certain position
  - **Slope** → one-dimensional function, single number
  - **Gradient** → vector of numbers (parameters)

# Gradient descent

- Let us define the **gradient** of the loss function  $L$  as:

$$\nabla_{\mathbf{W}} L(\mathbf{W}) = \left[ \frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_N} \right] \quad \text{Vector of partial derivatives}$$

- We can use the gradient to update the parameters of the model in an iterative way:

$$\mathbf{W}_t = \mathbf{W}_{t-1} - \eta \nabla_{\mathbf{W}} J(\mathbf{W}) \Big|_{\mathbf{W}=\mathbf{W}_{t-1}}$$

step size  
a.k.a. learning rate

# Gradient descent

- Let us define the **gradient** of the loss function  $L$  as:

$$\nabla_{\mathbf{W}} L(\mathbf{W}) = \left[ \frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_N} \right] \quad \text{Vector of partial derivatives}$$

- We can use the gradient to update the parameters of the model in an iterative way:

$$\mathbf{W}_t = \mathbf{W}_{t-1} - \eta \nabla_{\mathbf{W}} J(\mathbf{W}) \Big|_{\mathbf{W}=\mathbf{W}_{t-1}}$$

???

# Loss/cost function

- Do we apply gradient descent to the loss function or to the cost function?

$$\mathbf{W}_t = \mathbf{W}_{t-1} - \eta \nabla_{\mathbf{W}} J(\mathbf{W}) \Big|_{\mathbf{W}=\mathbf{W}_{t-1}}$$

- Remember that:  $J = \frac{1}{N} \sum_{i=1}^N L_i$
- In practice, we compute  $\nabla_{\mathbf{W}}$  for each sample and cumulate it
- In the end, we update by  $\frac{\nabla_{\mathbf{W}}}{N}$  (linear operations)

# Gradient descent

- **Batch gradient descent**

$$\mathbf{W}_t = \mathbf{W}_{t-1} - \eta \nabla_{\mathbf{W}} J(\mathbf{W}) \Big|_{\mathbf{W}=\mathbf{W}_{t-1}}$$

- **Stochastic gradient descent** (online, 1 random sample)

$$\mathbf{W}_t = \mathbf{W}_{t-1} - \eta \nabla_{\mathbf{W}} J(\mathbf{W}, \boxed{\mathbf{x}_i, y_i}) \Big|_{\mathbf{W}=\mathbf{W}_{t-1}}$$

- **Mini-batch gradient descent** (mini-batches of  $n$  random samples)

$$\mathbf{W}_t = \mathbf{W}_{t-1} - \eta \nabla_{\mathbf{W}} J(\mathbf{W}, \boxed{\mathbf{x}_{i:i+n}, y_{i:i+n}}) \Big|_{\mathbf{W}=\mathbf{W}_{t-1}}$$

What people now call “stochastic gradient descent” (SGD) is actually a “mini-batch gradient descent”

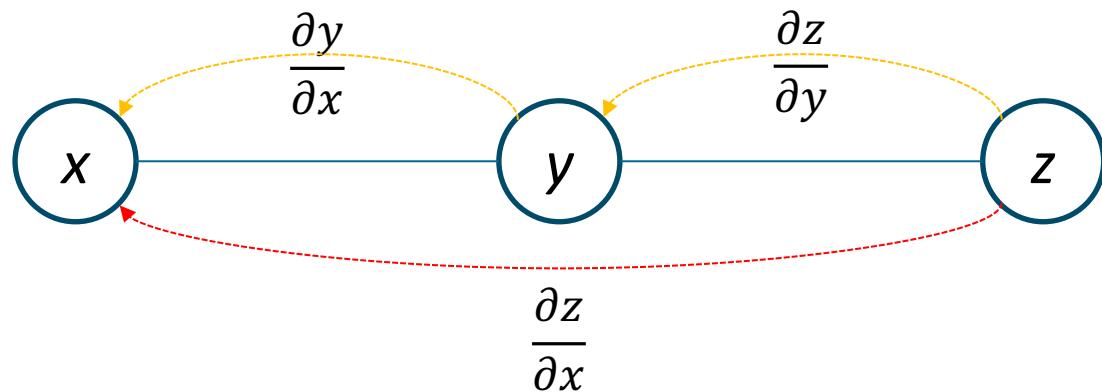
# Backpropagation

- Algorithm to **compute gradients** of expression through recursive application of the chain rule

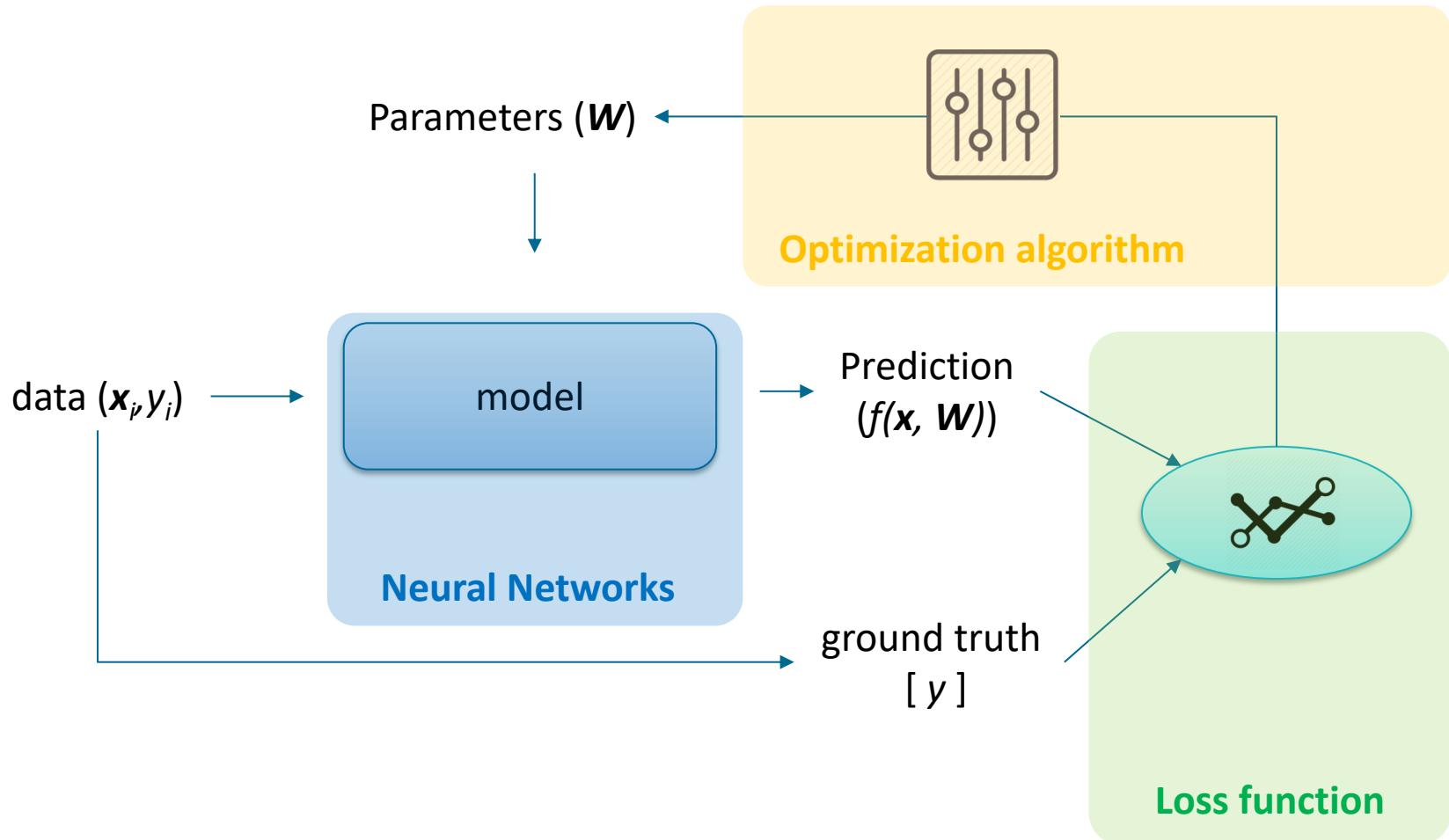
## CHAIN RULE

- Given three variables,  $(x, y, z)$ , such that  $z = f(y)$  and  $y = f(x)$ , it is possible to compute the derivative of  $z$  with respect to  $x$  as:

$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y} \frac{\partial y}{\partial x}$$



# Supervised learning procedure



# Training using mini-batches

- We select a subset of training samples
  - Mini-batch
- Forward pass
  - Compute predictions
  - Compute loss
  - Compute cost as a combination of the losses (e.g., average)
- Backpropagation
  - Compute derivatives

---

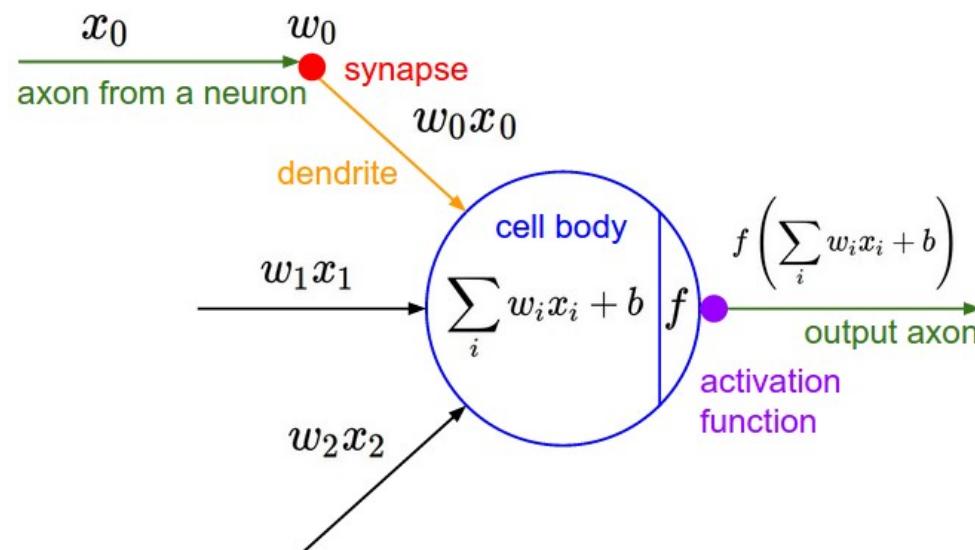
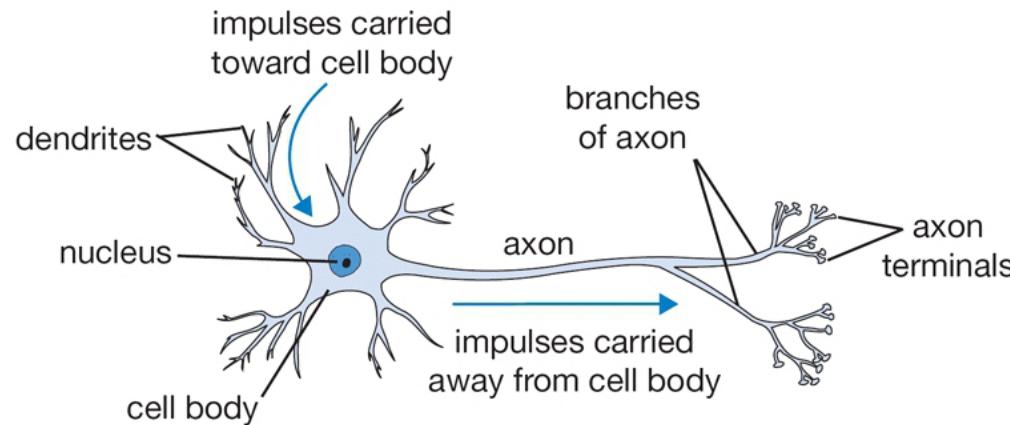
# Break?



---

# Neural Networks

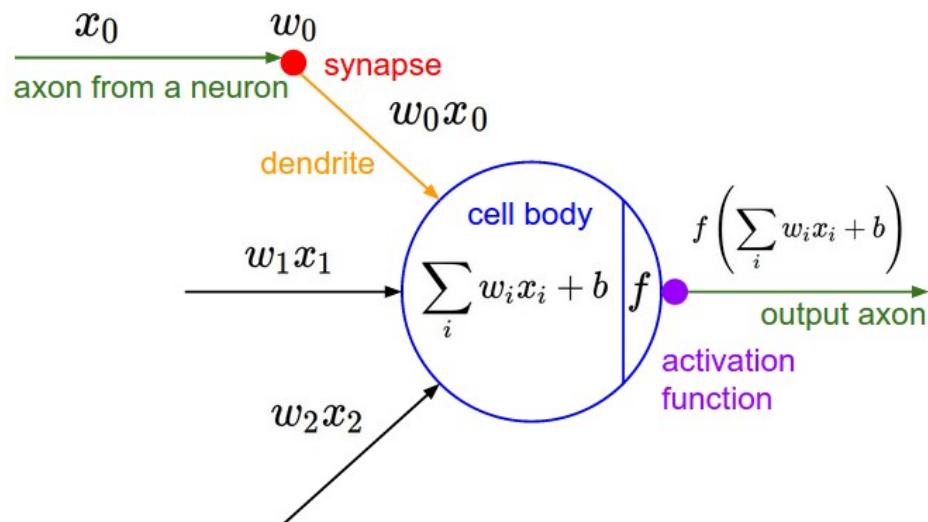
# Neural Networks - Inspiration



# Neural Networks

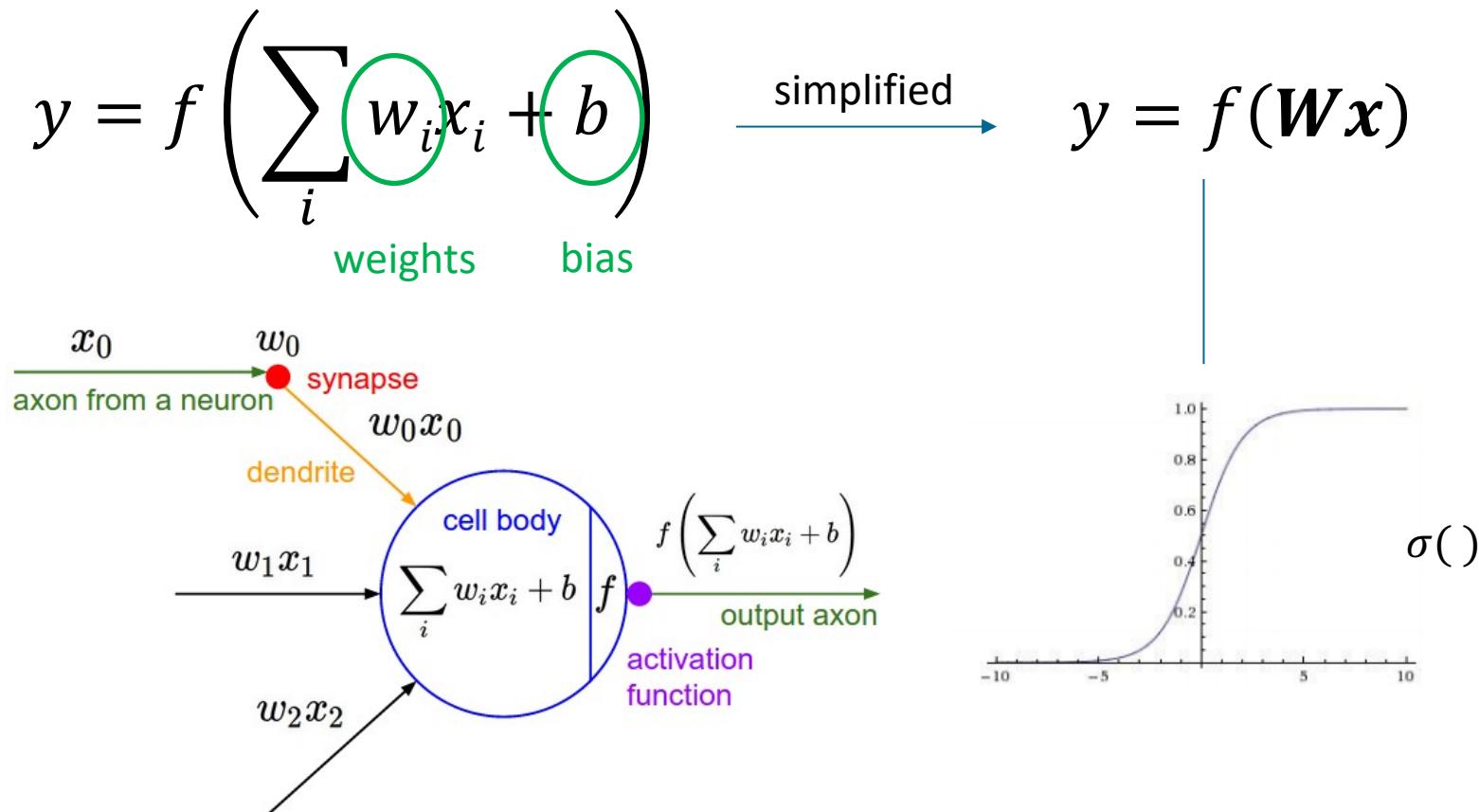
- Artificial Neuron

$$y = f \left( \sum_i w_i x_i + b \right)$$



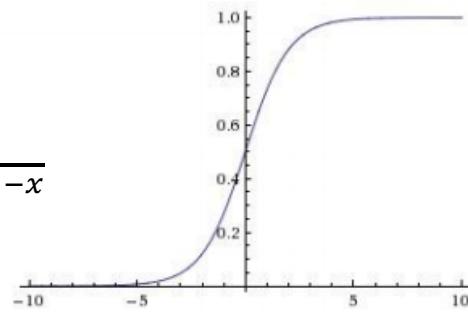
# Neural Networks

- Artificial Neuron

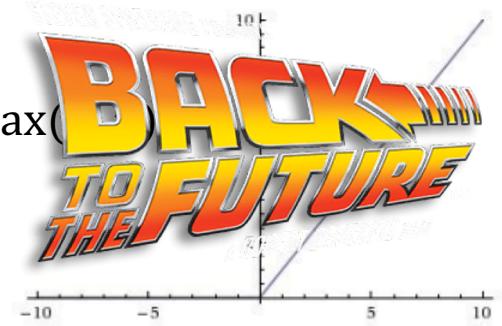


# Non linearity

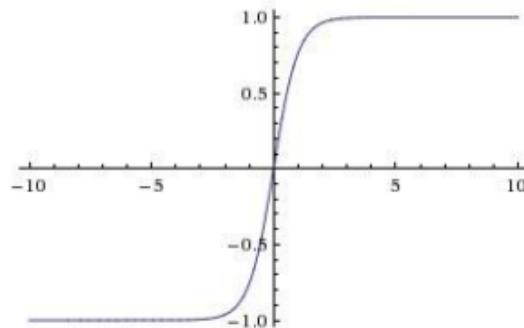
$$\text{Sigmoid } \sigma(x) = \frac{1}{1+e^{-x}}$$



ReLU max(

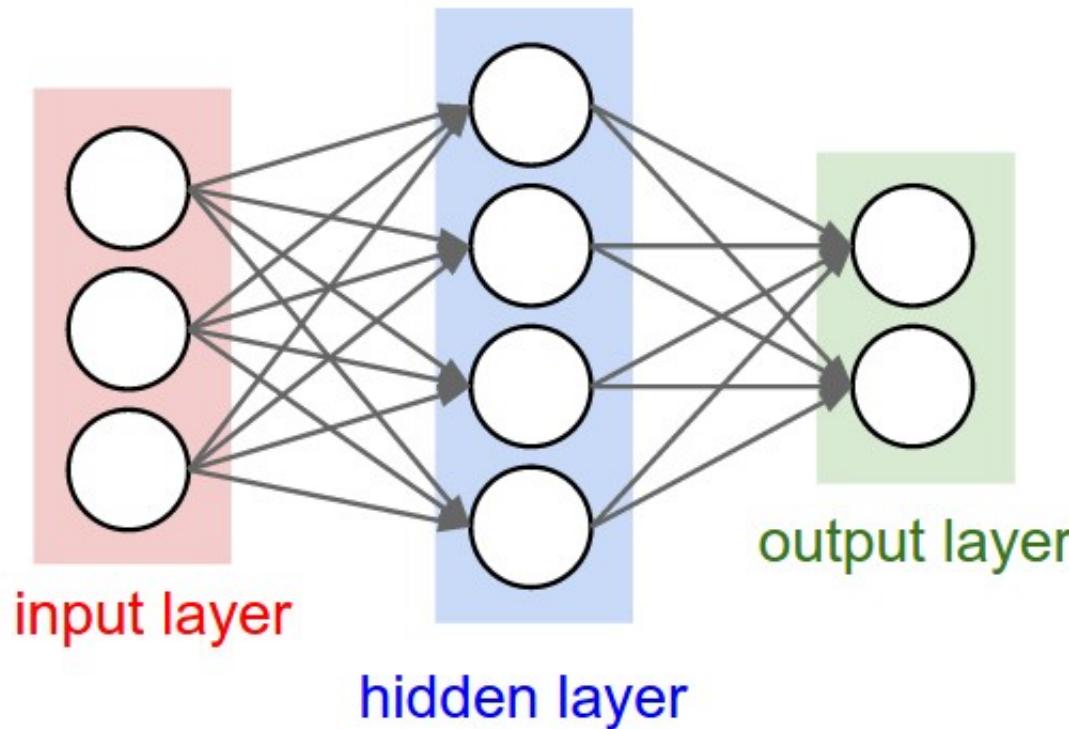


Tanh       $\text{Tanh}(x)$



# Neural Networks

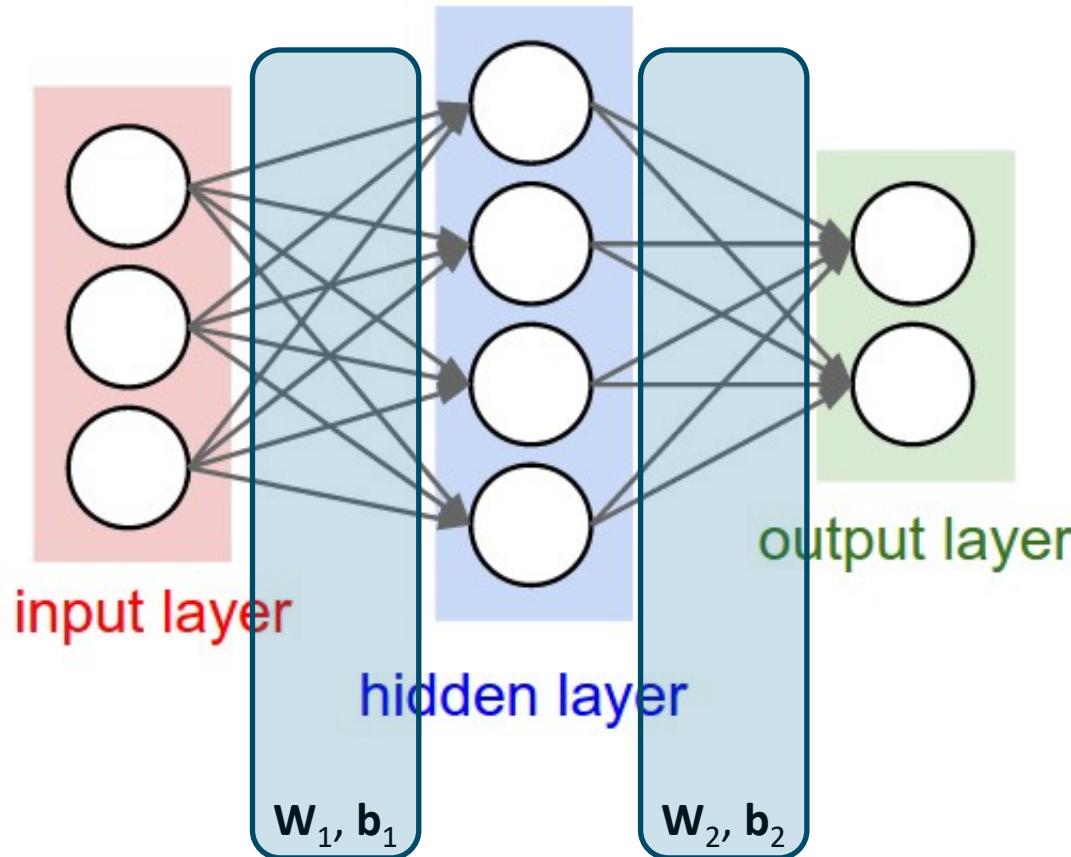
- More than one neuron and more than one layer



Two-layer  
neural network

# Neural Networks

- More than one neuron and more than one layer

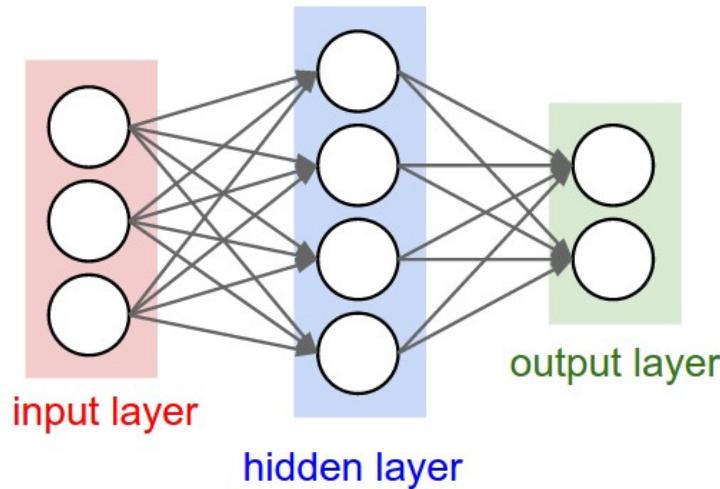


Two-layer  
neural network

$$f = \sigma(W_2\sigma(W_1x))$$

# Neural Networks

- How many parameters?



## Hidden layer

- 4 neurons connected to 3 inputs
- $4 \times 3$  weights + 4 biases = 16

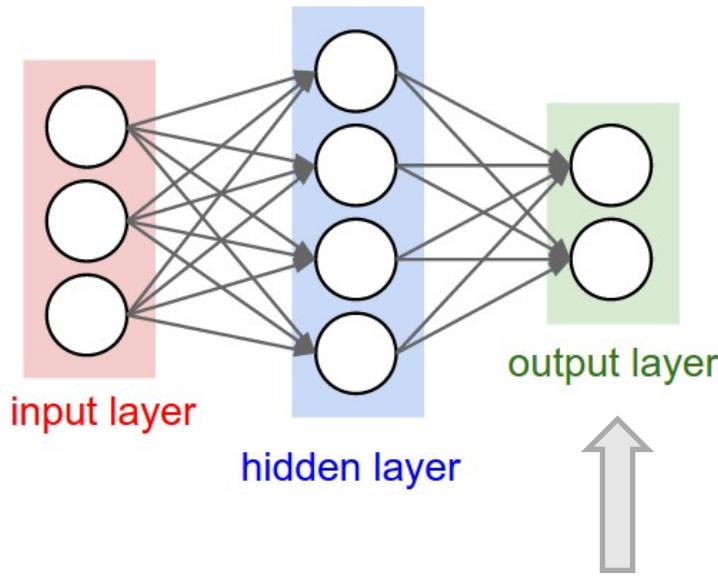
## Output layer

- 2 neurons connected to 4 neurons
- $2 \times 4$  weights + 2 biases = 10

TOTAL = 26 parameters

# Neural Networks

- How many parameters?



The number of neurons in the output layer depends on the classification problem (# classes)

## Hidden layer

- 4 neurons connected to 3 inputs
- $4 \times 3$  weights + 4 biases = 16

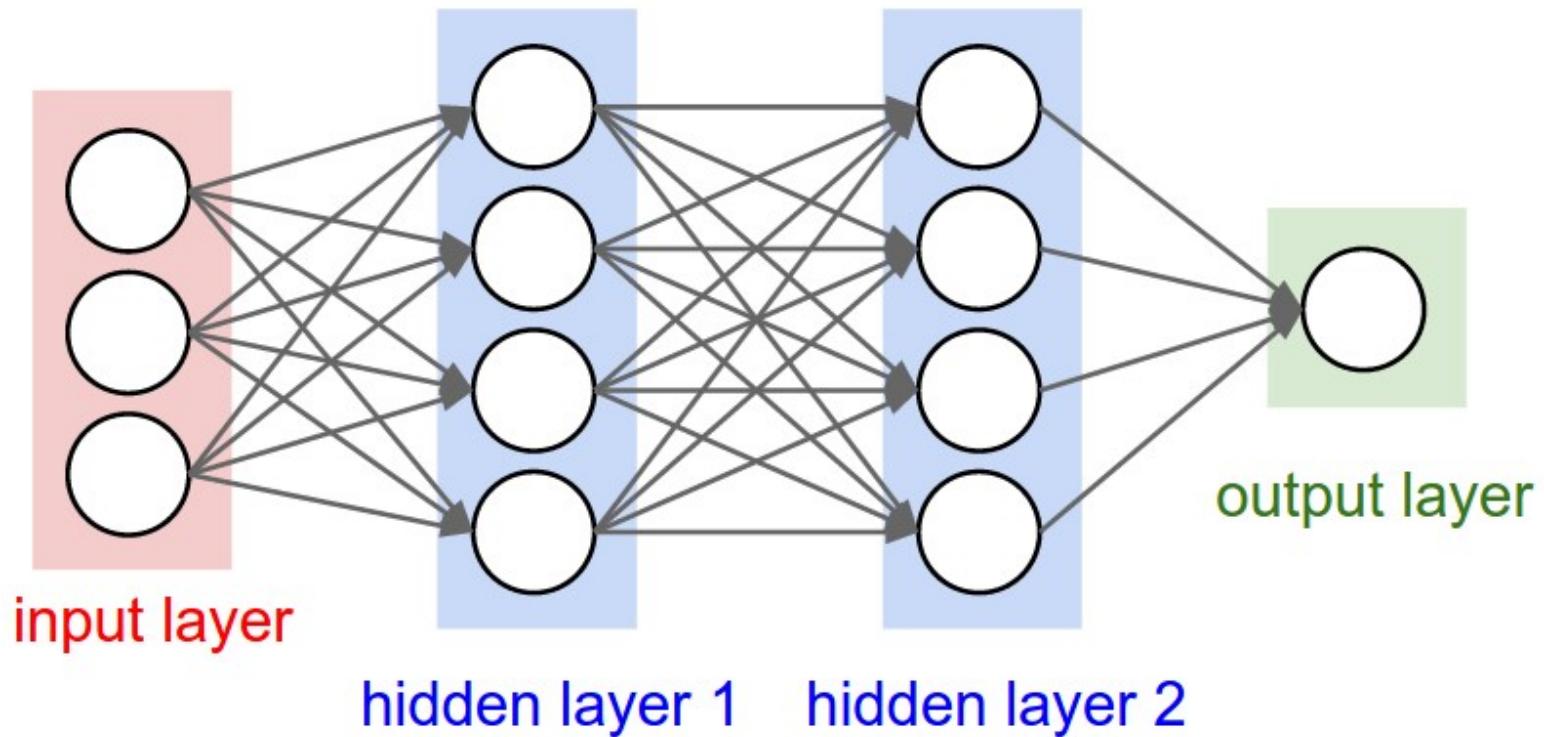
## Output layer

- 2 neurons connected to 4 neurons
- $2 \times 4$  weights + 2 biases = 10

TOTAL = 26 parameters

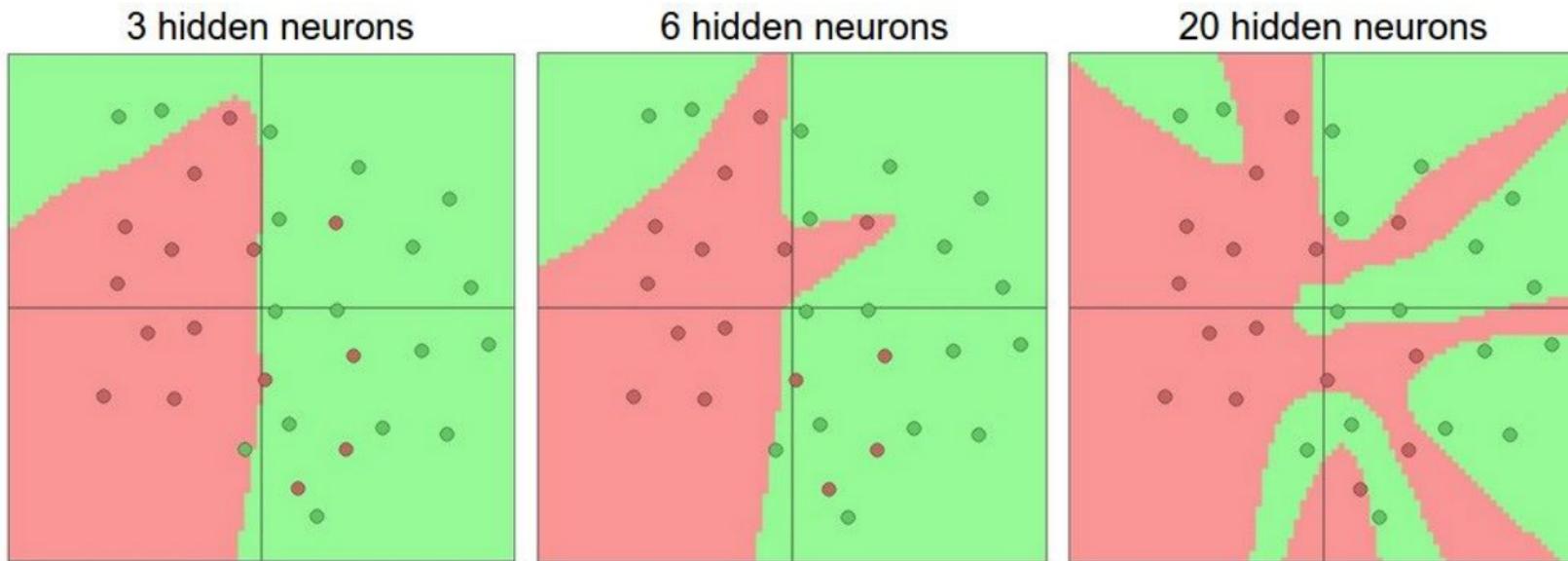
# Neural Networks

- $f = \sigma(W3\sigma(W2\sigma(W1x)))$



# Neurons in hidden layer

Demo = <http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>



more neurons = more capacity

# Non-linearity

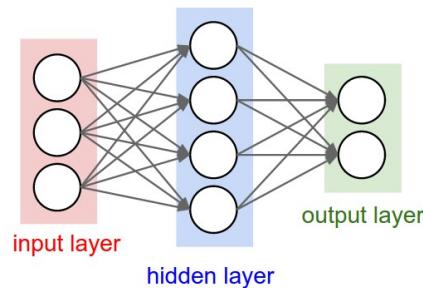
- What if we do not incorporate non-linearity?

$$\bullet \quad f = \underbrace{W_2(W_1x)}_{\text{2nd Layer}} = (W_2W_1)x = W'x \quad \text{Still a linear transformation!!}$$

1<sup>st</sup> layer

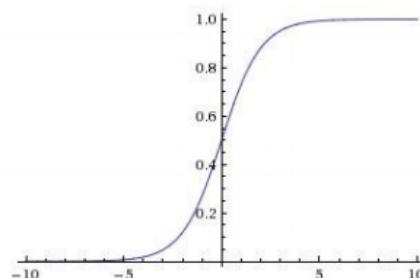
# Output layer

- The desired output for a multi-class classification problem is:
  - One classification score per class



Achieved by using one neuron per class

- Each score indicates the probability for its class ( $P \in [0,1]$ )



A sigmoid function guarantees output in the range  $[0,1]$

# Softmax

- What if we have more than one class? We use Softmax!
- Softmax is a generalization of logistic regression (sigmoid function) to multiple classes
- If we have K neurons in the output layer, the probability for class with index  $i$  is, based on softmax:

$$P_i = \frac{e^{z_i}}{\sum_k e^{z_k}}$$

Always positive

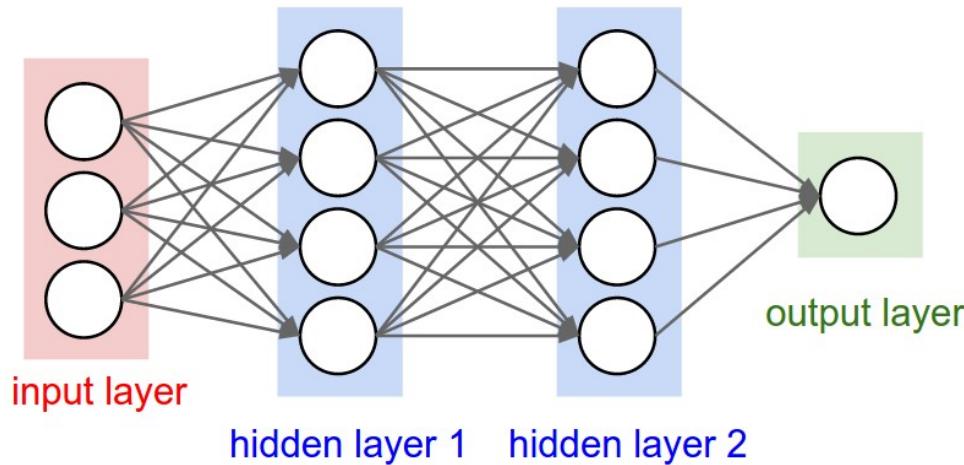
Does the normalization

# Training a Neural Networks

## How to

# 1. Define an architecture

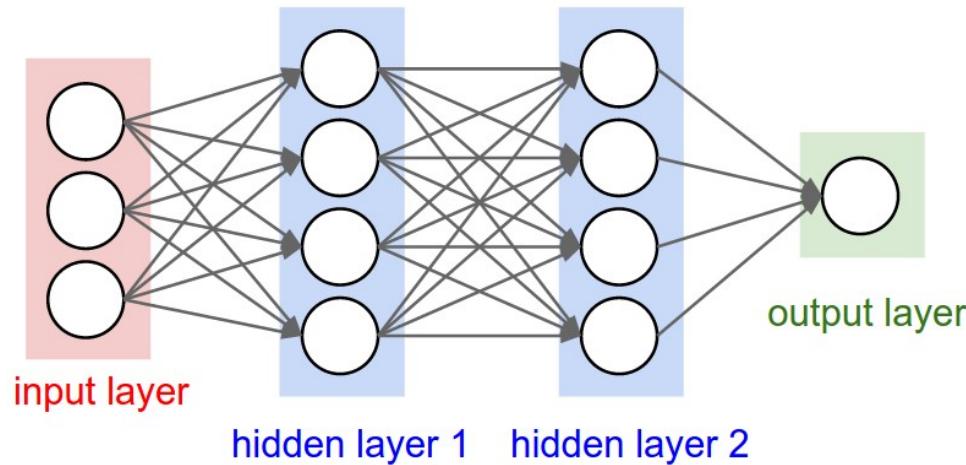
- Let's pick this architecture and count its parameters



Layer	# input connections	# output connections	# parameters
Input layer	N/A	3	0
Hidden layer 1	$4 * 3$	3	$12 + 4 = 16$
Hidden layer 2	$4 * 4$	4	$16 + 4 = 20$
Output layer	4	1	$4 + 1 = 5$

## 2. Initialization

- How to initialize the parameters that have to be trained?



Layer	# input connections	# output connections	# parameters	Initial values
Input layer	N/A	3	0	N/A
Hidden layer 1	$4 * 3$	3	$12 + 4 = 16$	???
Hidden layer 2	$4 * 4$	4	$16 + 4 = 20$	???
Output layer	4	1	$4 + 1 = 5$	???

# Initialization

- How do we initialize **weights** and **biases**?

## 1. All parameters set to zero

- When we compute the gradient, all parameters will get the same update!

## 2. Random initialization with normal distribution

- Zero mean and unit standard deviation

The variance of the output of a neuron increases with the number of inputs → we can normalize the variance accordingly:

$$w = np.random.randn(n) / \text{sqrt}(n)$$

number of inputs  
for that neuron

# Initialization

- This idea has been further developed, even recently:
  - Glorot10
  - He15

## glorot\_normal

```
glorot_normal(seed=None)
```

Glorot normal initializer, also called Xavier normal initializer.

It draws samples from a truncated normal distribution centered on 0 with

`stddev = sqrt(2 / (fan_in + fan_out))` where `fan_in` is the number of input units in the weight tensor and `fan_out` is the number of output units in the weight tensor.

### Arguments

- `seed`: A Python integer. Used to seed the random generator.

### Returns

An initializer.

### References

Glorot & Bengio, AISTATS 2010 - <http://jmlr.org/proceedings/papers/v9/glorot10a/glorot10a.pdf>

## he\_normal

```
he_normal(seed=None)
```

He normal initializer.

It draws samples from a truncated normal distribution centered on 0 with

`stddev = sqrt(2 / fan_in)` where `fan_in` is the number of input units in the weight tensor.

### Arguments

- `seed`: A Python integer. Used to seed the random generator.

### Returns

An initializer.

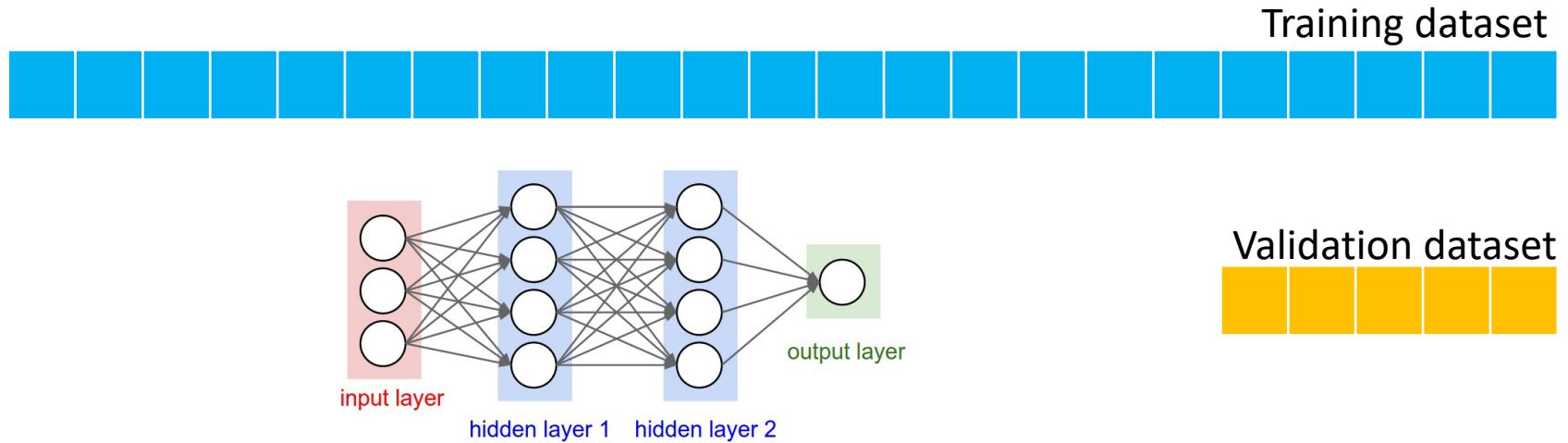
### References

He et al., <http://arxiv.org/abs/1502.01852>



### 3. Training and validation set

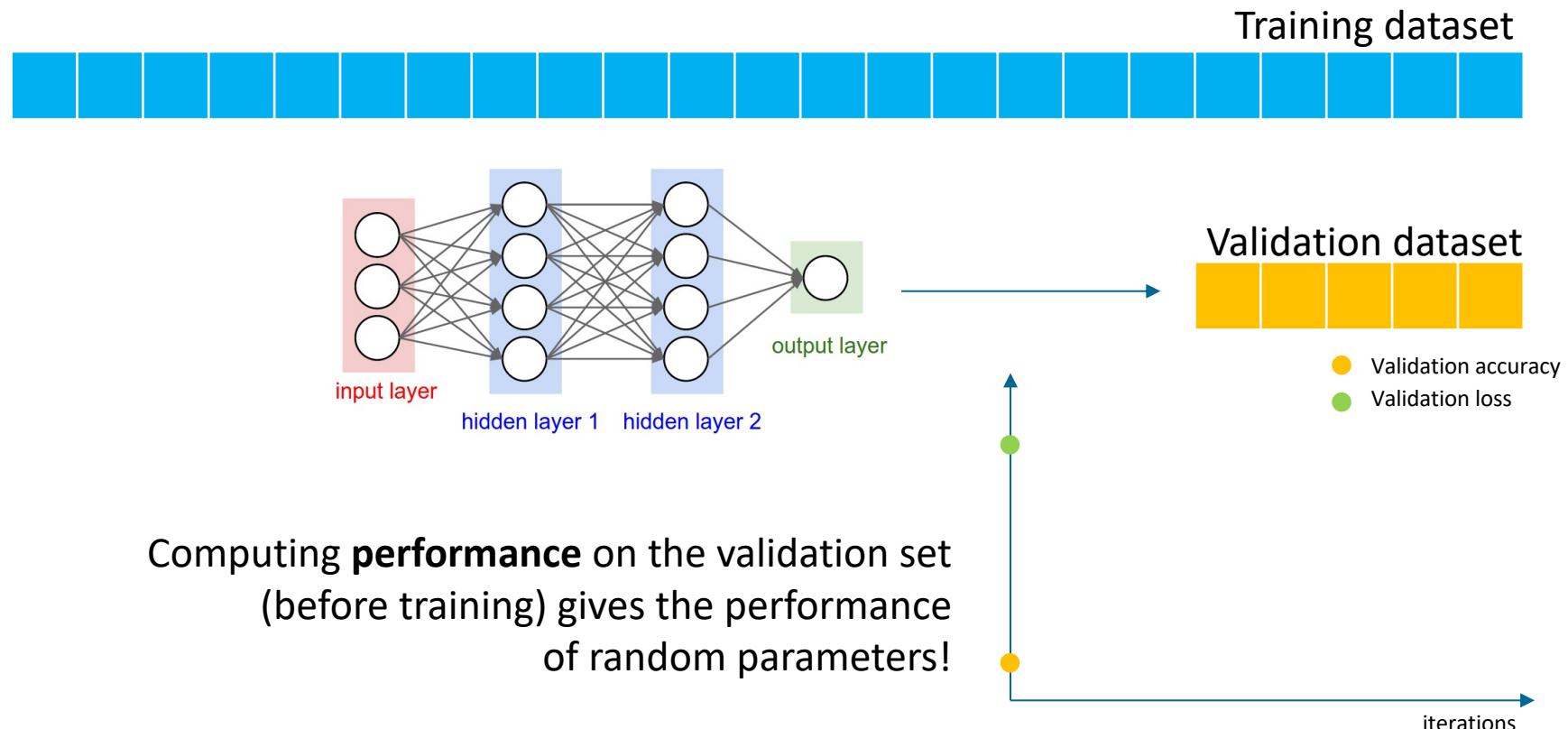
- We need data to train our network -> training set
- We need data to assess its performance -> validation set



**SHUFFLE YOUR TRAINING SET BEFORE TRAINING!**

## 3b. Initial validation performance

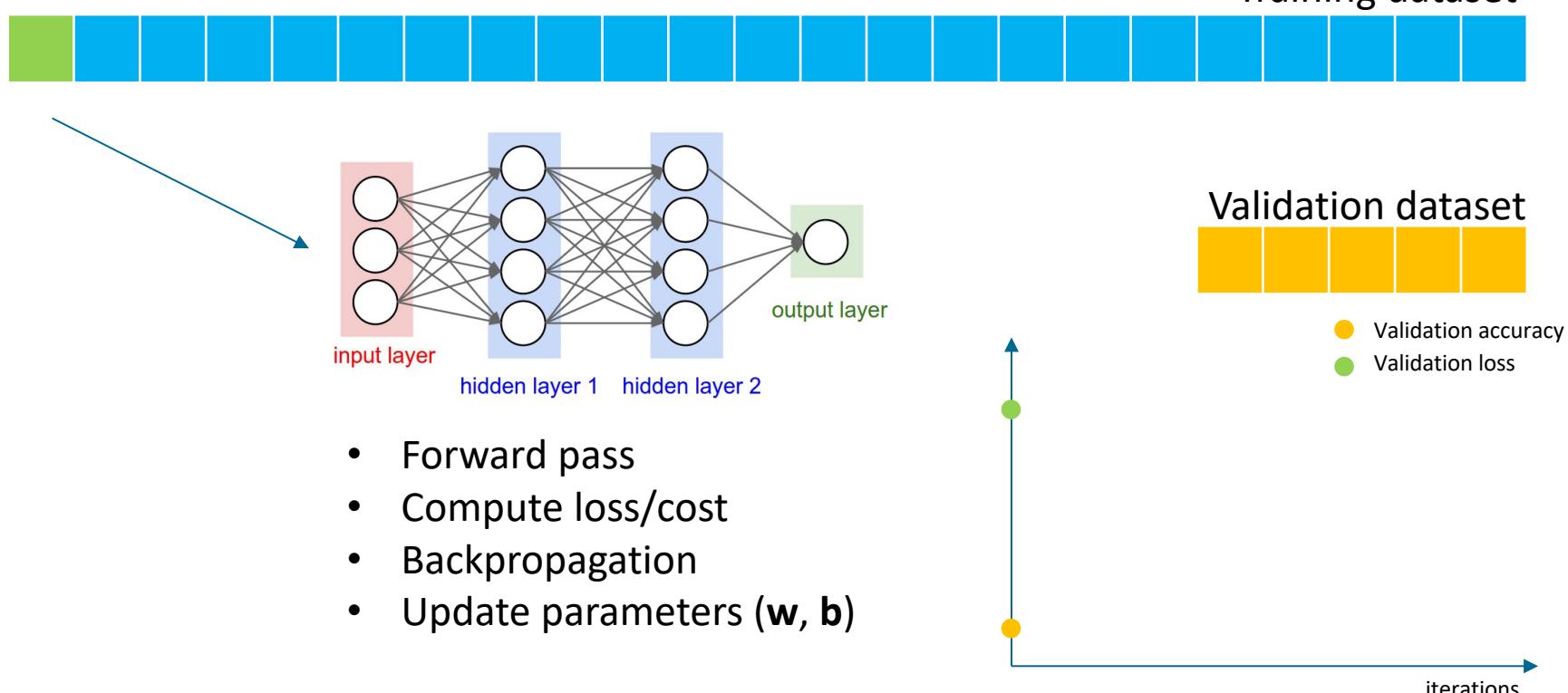
- We need data to train our network -> training set
- We need data to assess its performance -> validation set



# 4. Training

- We train neural networks with **mini-batch** gradient descent

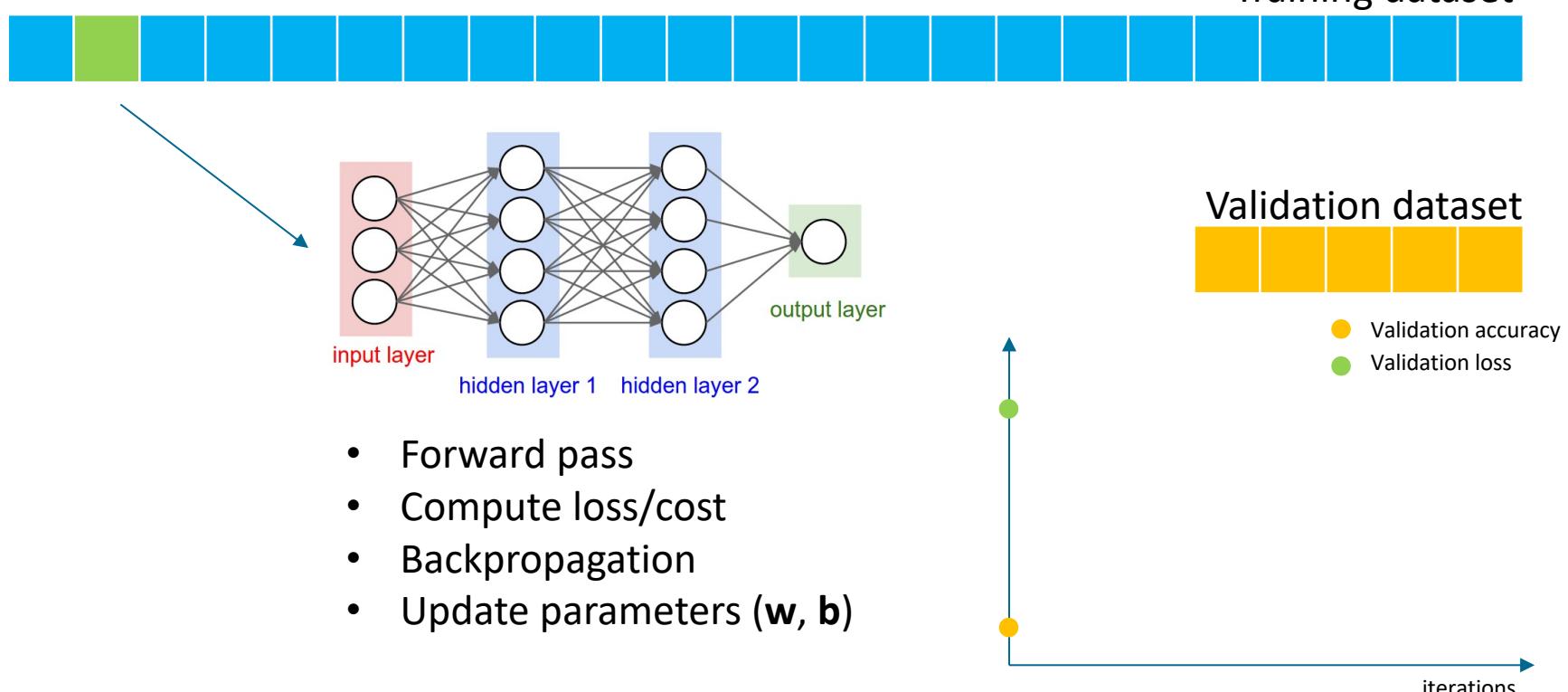
mini-batch = set of (few) training samples



# 4. Training

- We train neural networks with **mini-batch** gradient descent

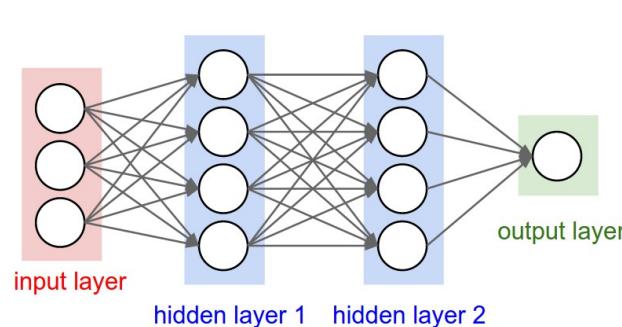
mini-batch = set of (few) training samples



# 4. Training

- We train neural networks with **mini-batch** gradient descent

mini-batch = set of (few) training samples

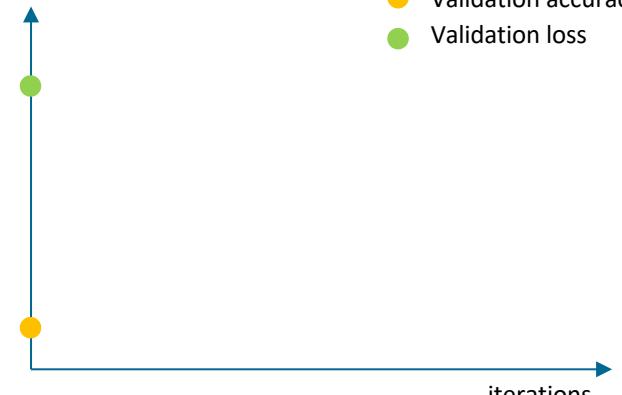


- Forward pass
- Compute loss/cost
- Backpropagation
- Update parameters ( $w, b$ )

Training dataset



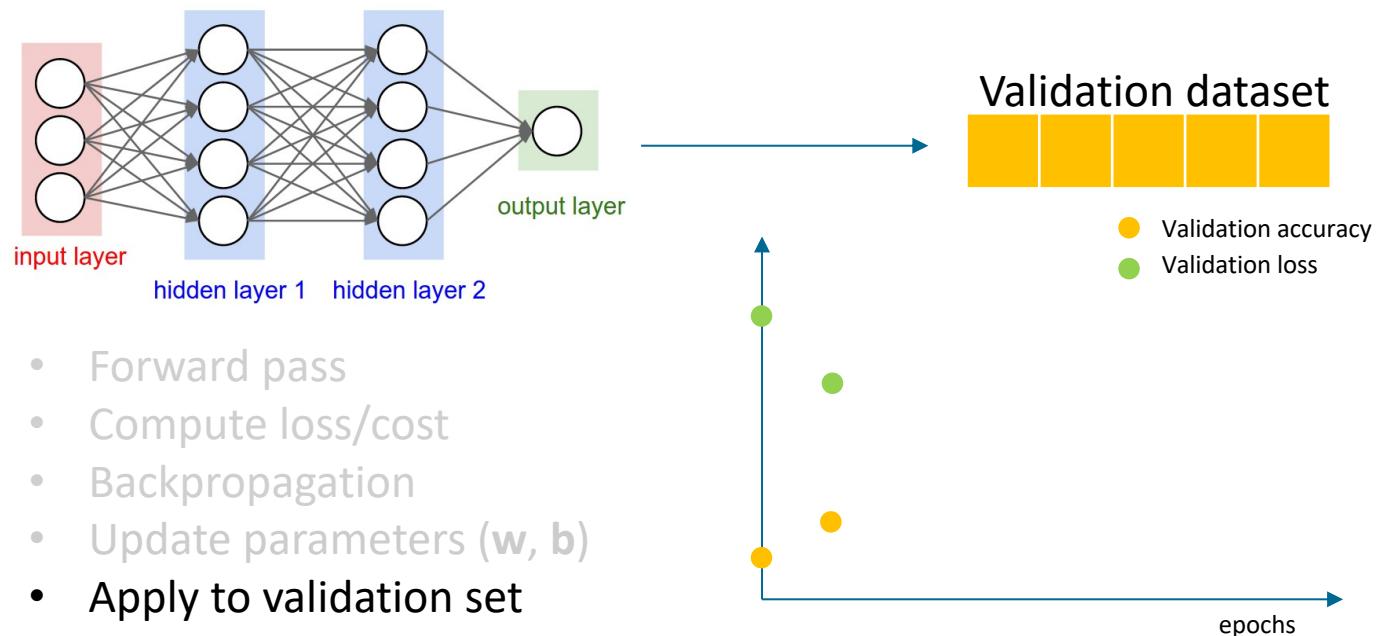
Validation accuracy  
Validation loss



# 4. Training

- When the network has seen all the mini-batches, it has done an **epoch**

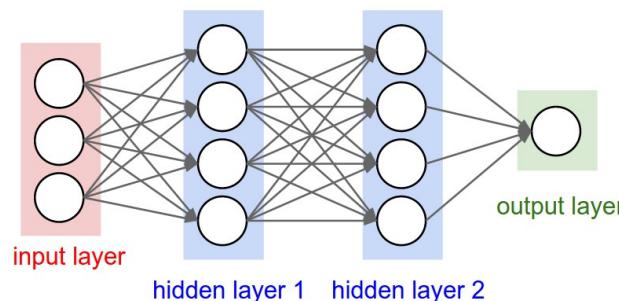
mini-batch = set of (few) training samples



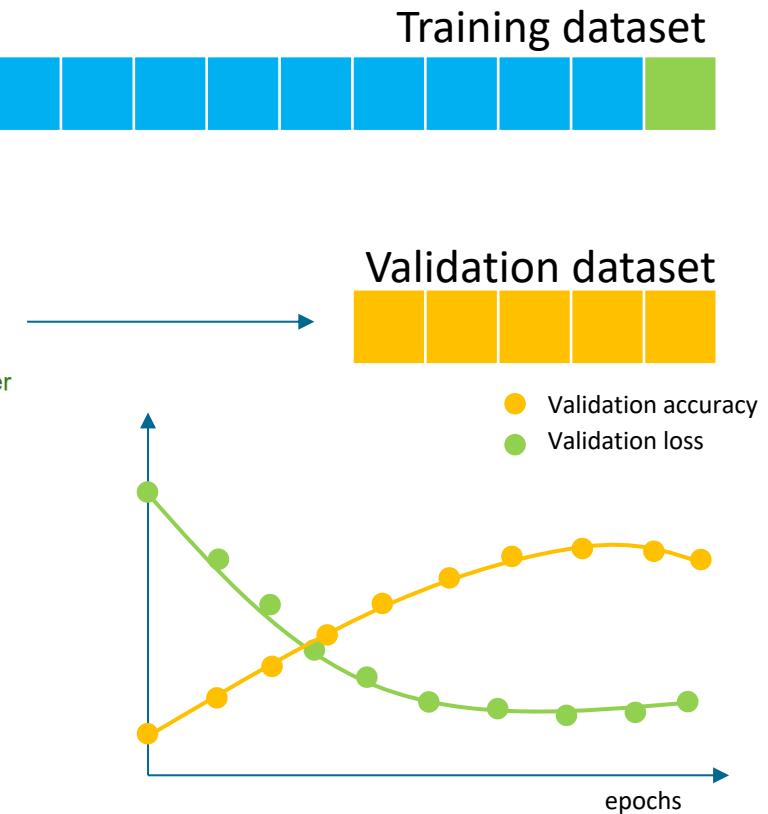
# 4. Training

- Repeat this several epochs

mini-batch = set of (few) training samples



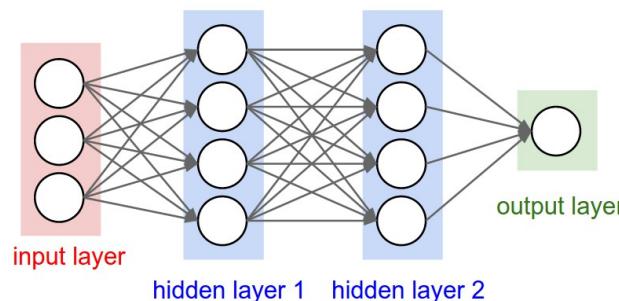
- Forward pass
- Compute loss/cost
- Backpropagation
- Update parameters ( $w, b$ )
- Apply to validation set



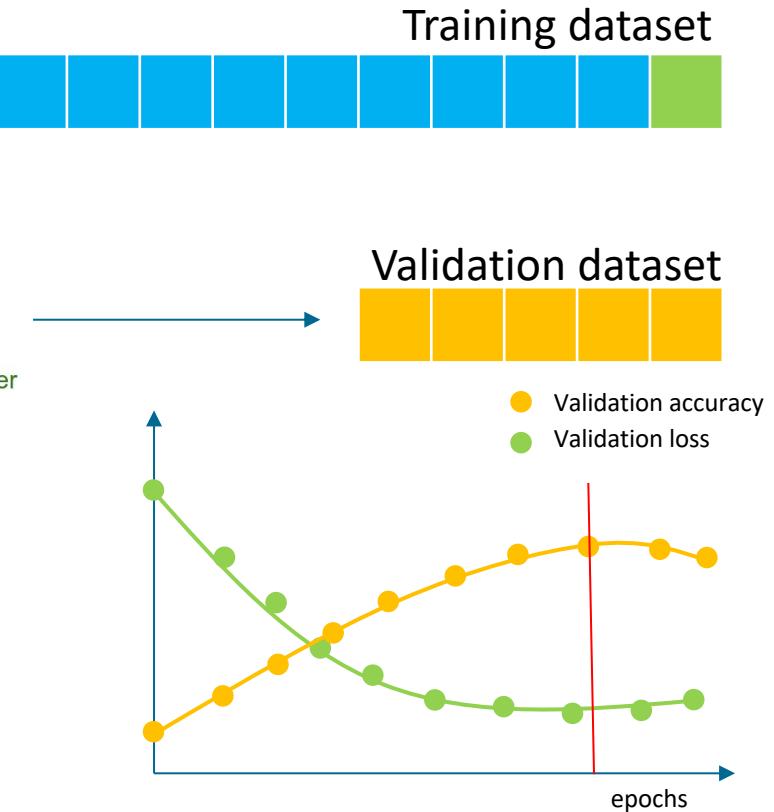
# 4. Training

- Save the parameters of the network at each iteration
- Pick the ones that maximize the accuracy (for example)

mini-batch = set of (few) training samples



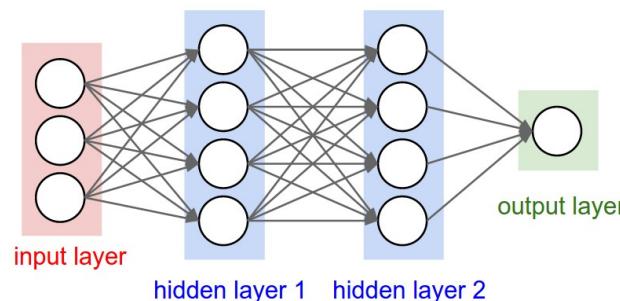
- Forward pass
- Compute loss/cost
- Backpropagation
- Update parameters ( $w, b$ )
- Apply to validation set



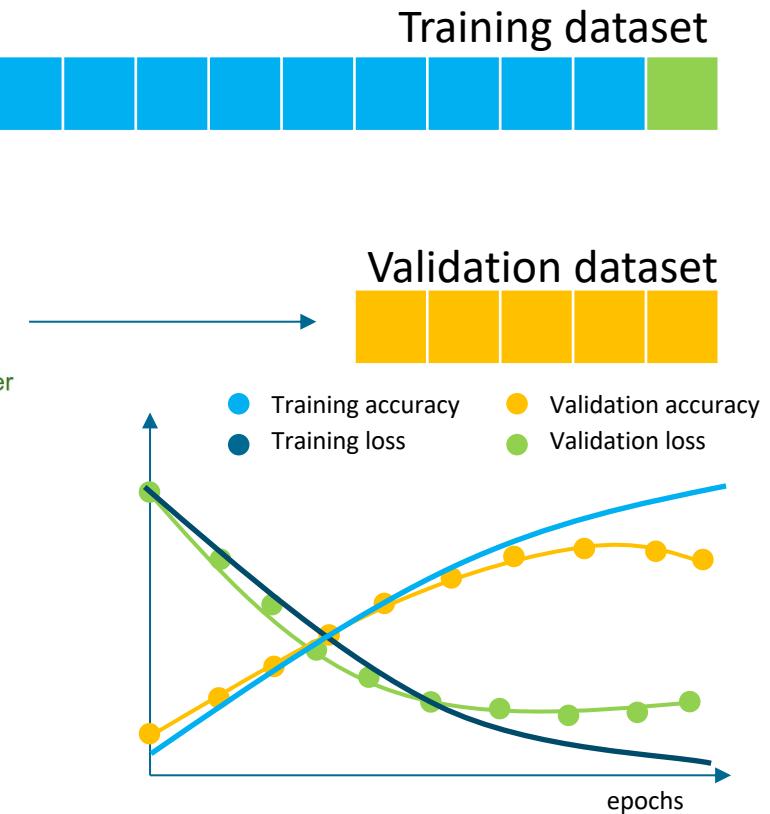
# 4. Training

- During training you also monitor **training loss** and **accuracy**

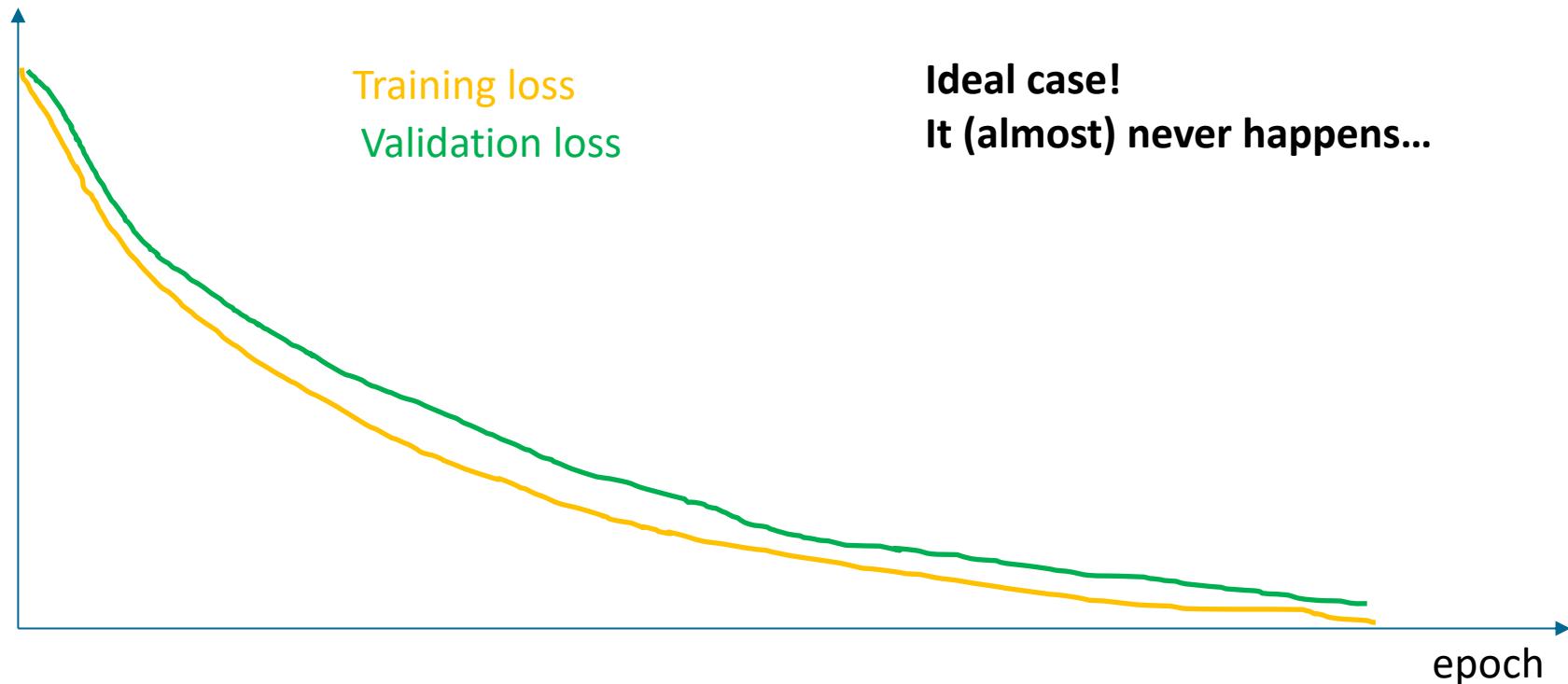
mini-batch = set of (few) training samples



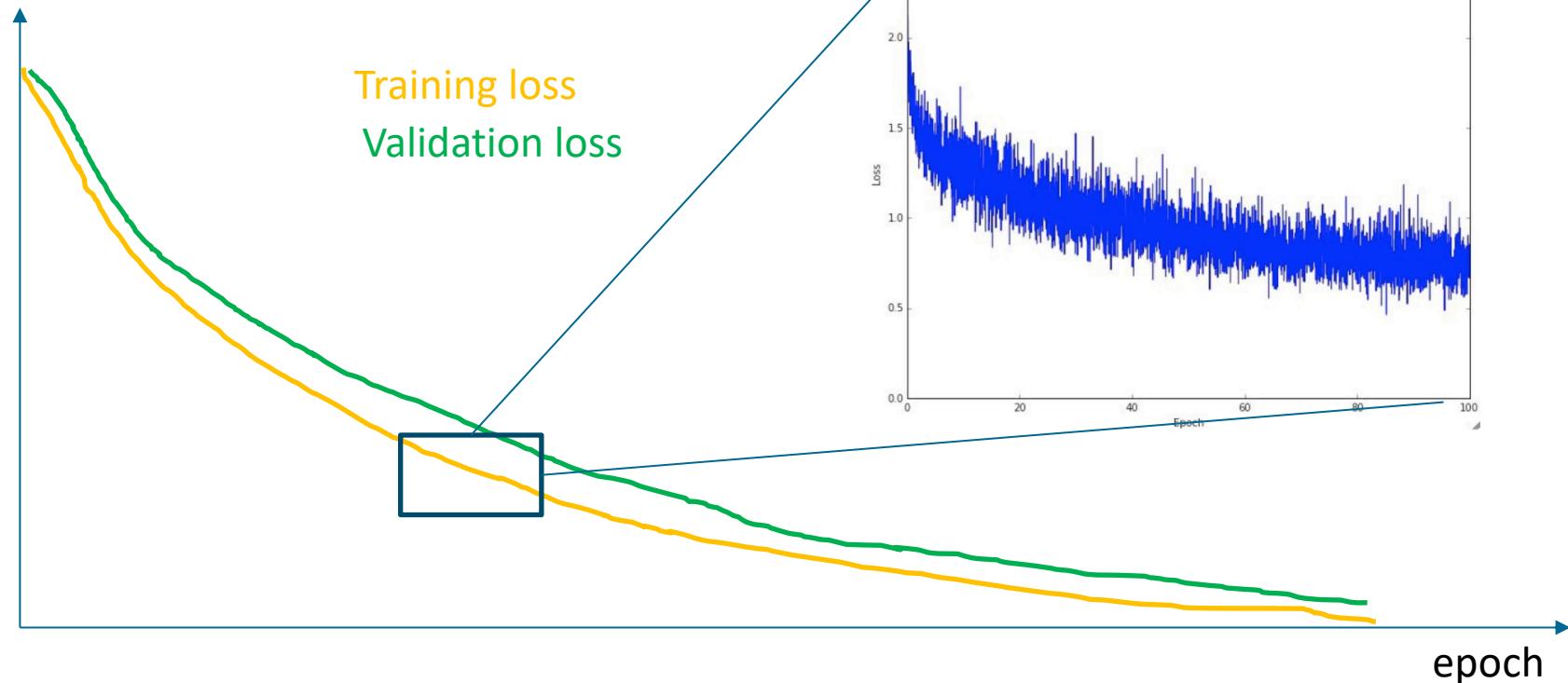
- Forward pass
- Compute loss/cost
- Backpropagation
- Update parameters ( $w, b$ )
- Apply to validation set



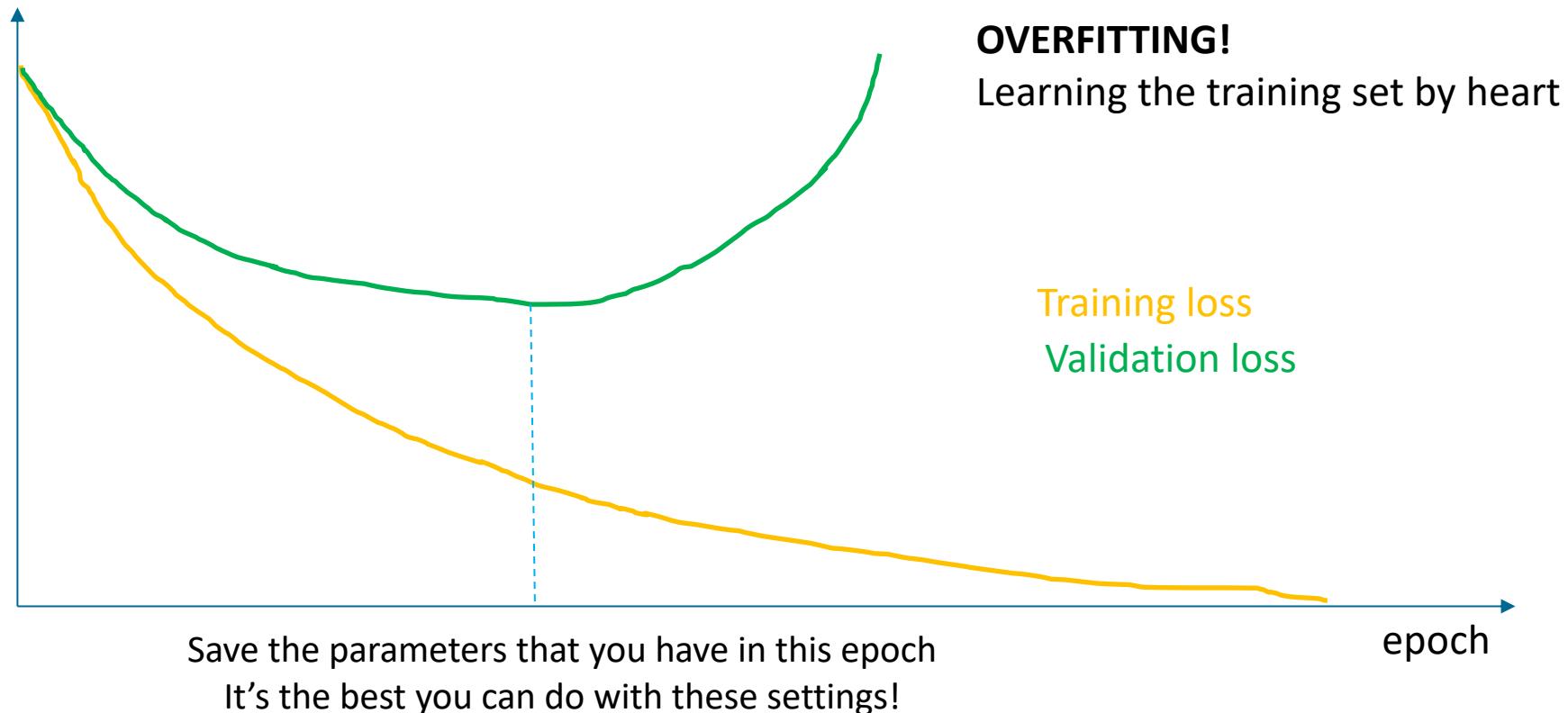
# Learning curves



# Learning curves

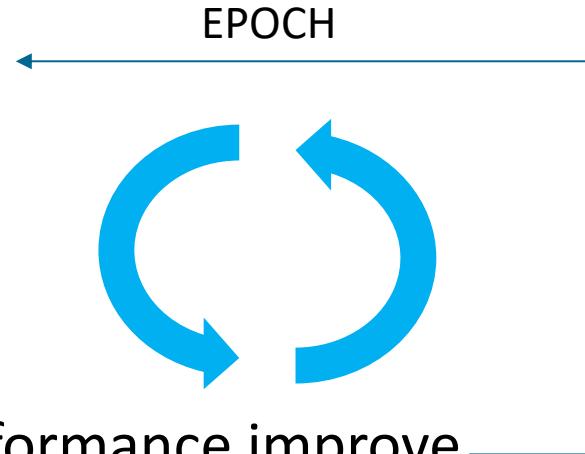


# Learning curves



# Training overview

1. Define network architecture
2. Initialize parameters
3. Define training and validation sets
4. Do training
  - a. Sample training mini-batches
  - b. Forward pass
  - c. Backpropagation
  - d. Update parameters
  - e. Apply on validation set
  - f. Save network parameters if performance improve

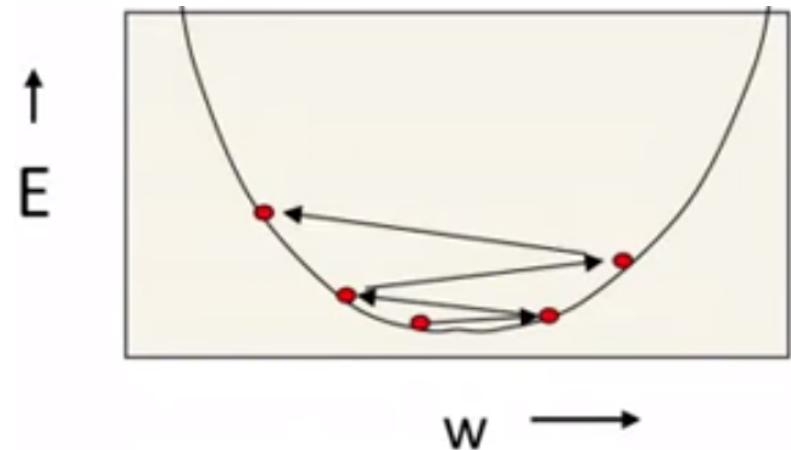
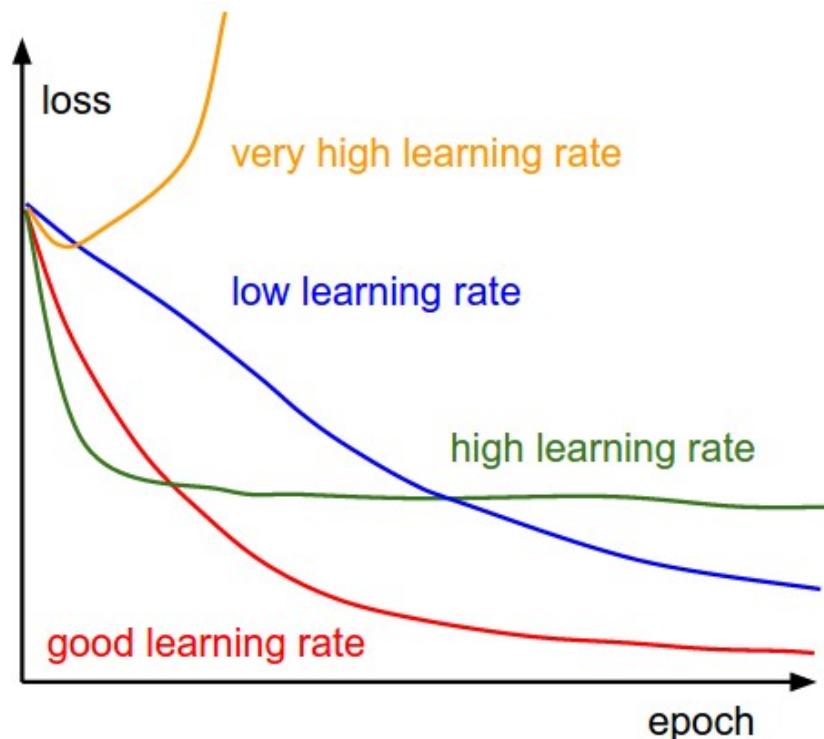


# Design choices

1. Define network architecture ← # layers, # neurons
2. Initialize parameters ← algorithm
3. Define training and validation sets ← data split
4. Do training
  - a. Sample training mini-batches ← mini-batch size
  - b. Forward pass
  - c. Backpropagation ← { loss function  
regularization }
  - d. Update parameters ← { Learning rate  
update rule (SGD) }
  - e. Apply on validation set
  - f. Save network parameters if performance improve

# Learning rate

- Learning rate is an important parameter
  - Too low: very slow convergence
  - Too high: No convergence!



# Learning rate

- How do we pick a good learning rate?
- In practice, we only pick the **initial one**, then we implement a strategy to decrease it over the epochs
  - Linear decrease
  - Exponential decrease
  - ...
- You can decrease it every M epochs, or when the performance of the network doesn't improve for N epochs

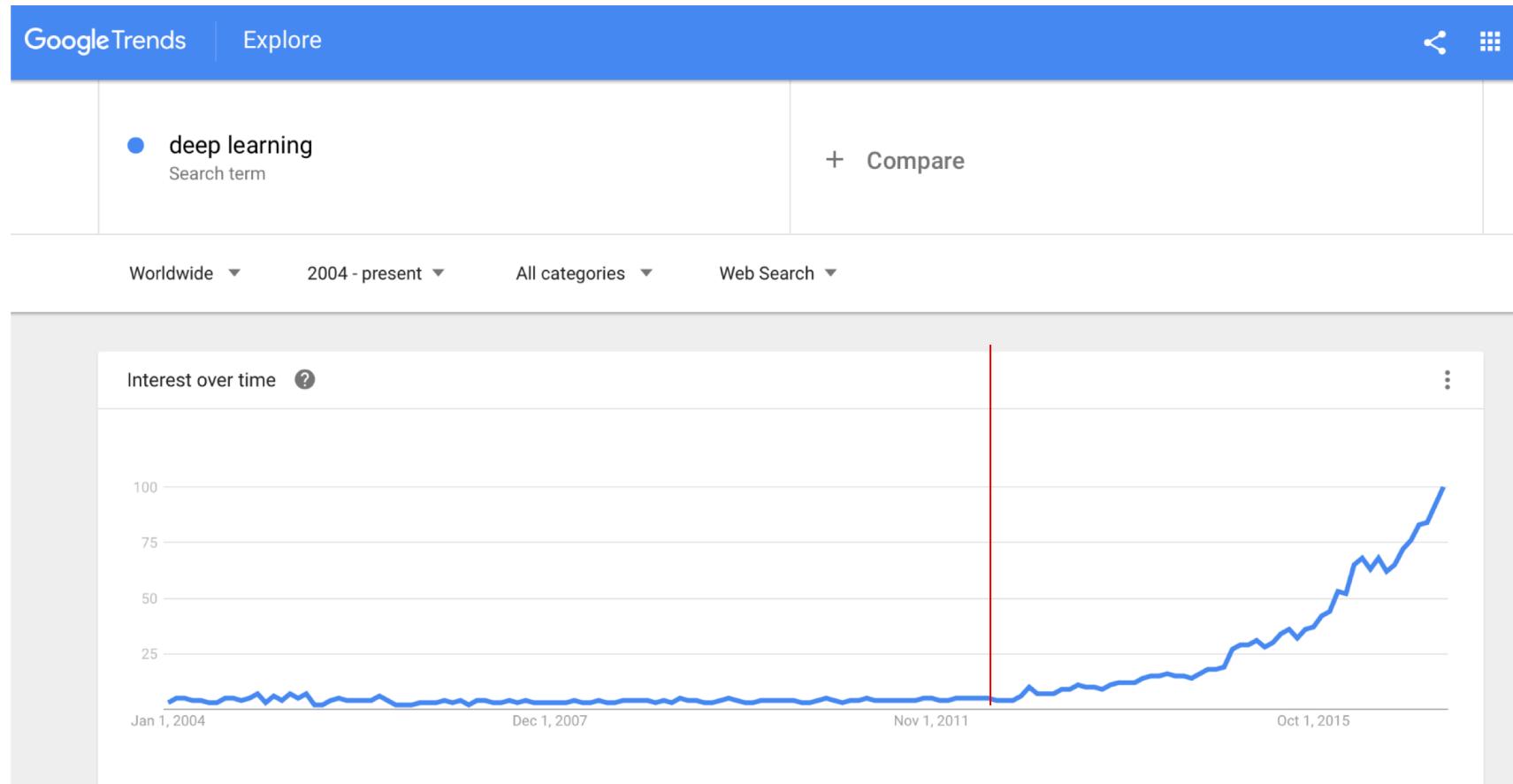
# Useful tricks and sanity checks

- **Sanity checks**
  - Train and validate on your (small) training set -> you should get loss  $\approx 0$  and accuracy  $\approx 1$  (over-fit training set)
  - Run your network just initialized on the validation set. If balanced, you should know the accuracy that you get
- **Tricks**
  - Make sure your batches are balanced
  - Make a balanced validation set, if accuracy is the target
  - **Shuffle your dataset!**
  - Penalize samples in the loss function in a different way

---

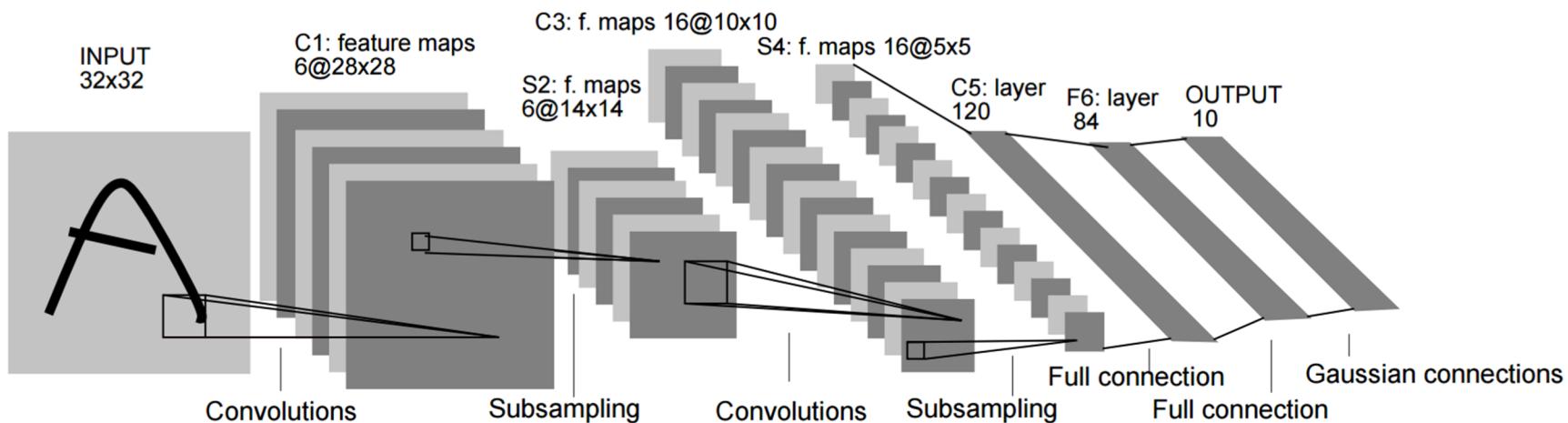
# Convolutional Neural Networks

# What happened in 2012?



# Yann LeCun, 1998

- **LeNet5**
  - Used backpropagation
  - Used stochastic gradient descent



# Lots of data

14,197,122 images, 21841 synsets indexed

Explore Download Challenges Publications CoolStuff About

Not logged in. Login | Signup

ImageNet is an image database organized according to the [WordNet](#) hierarchy (currently only the nouns), in which each node of the hierarchy is depicted by hundreds and thousands of images. Currently we have an average of over five hundred images per node. We hope ImageNet will become a useful resource for researchers, educators, students and all of you who share our passion for pictures.

[Click here](#) to learn more about ImageNet, [Click here](#) to join the ImageNet mailing list.

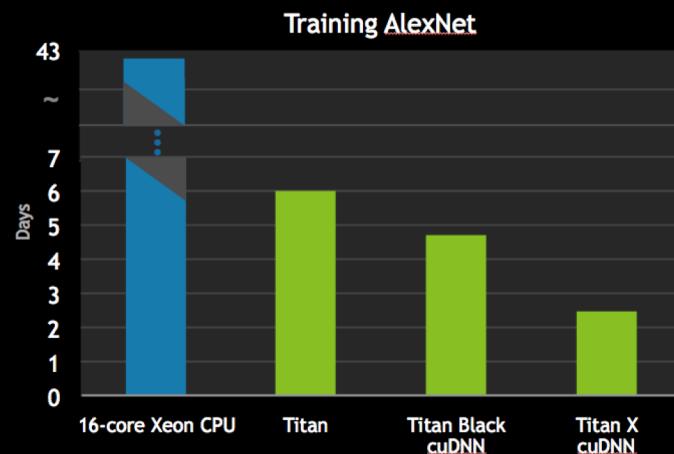


What do these images have in common? [Find out!](#)

# GPUs

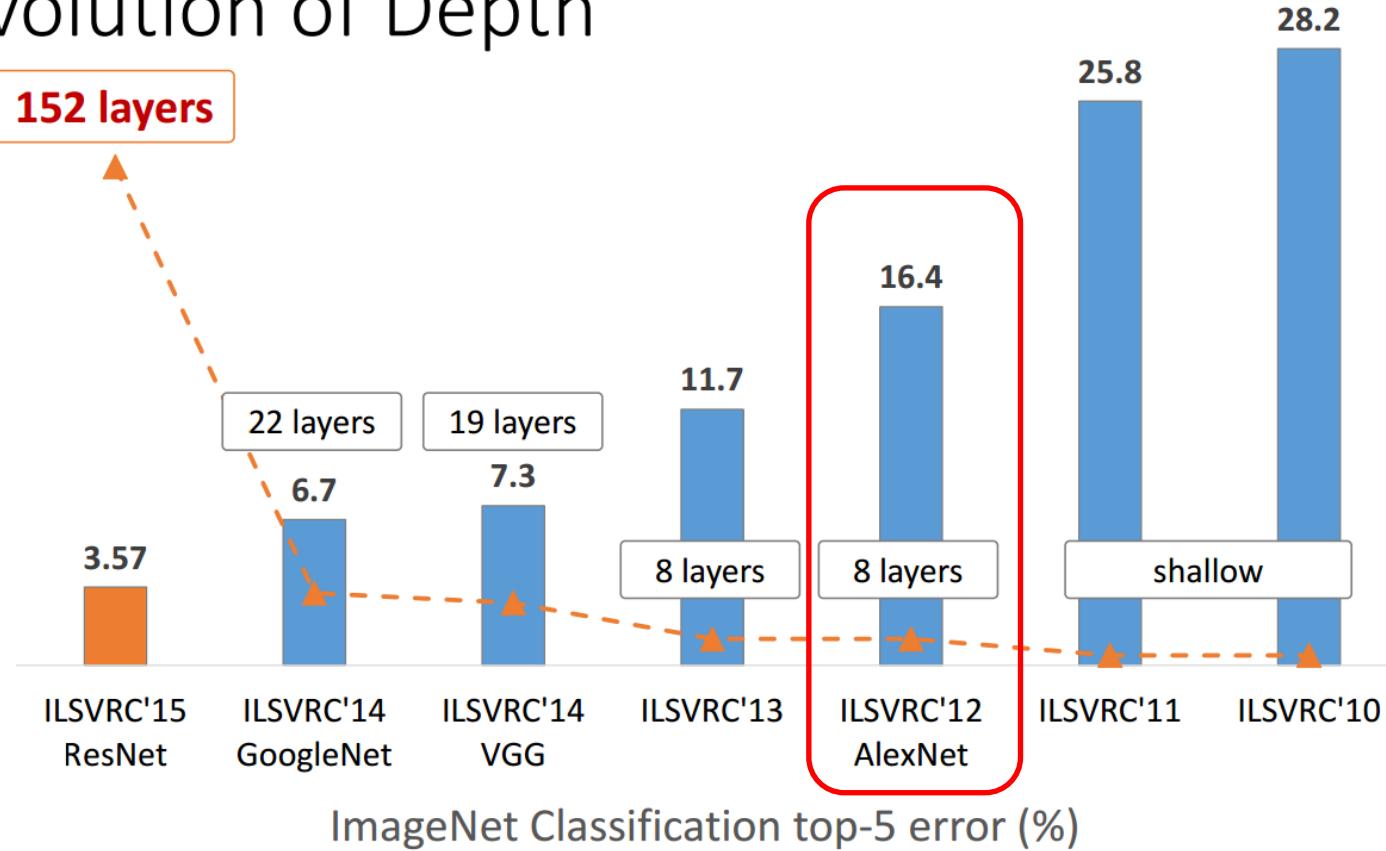


## TITAN X FOR DEEP LEARNING



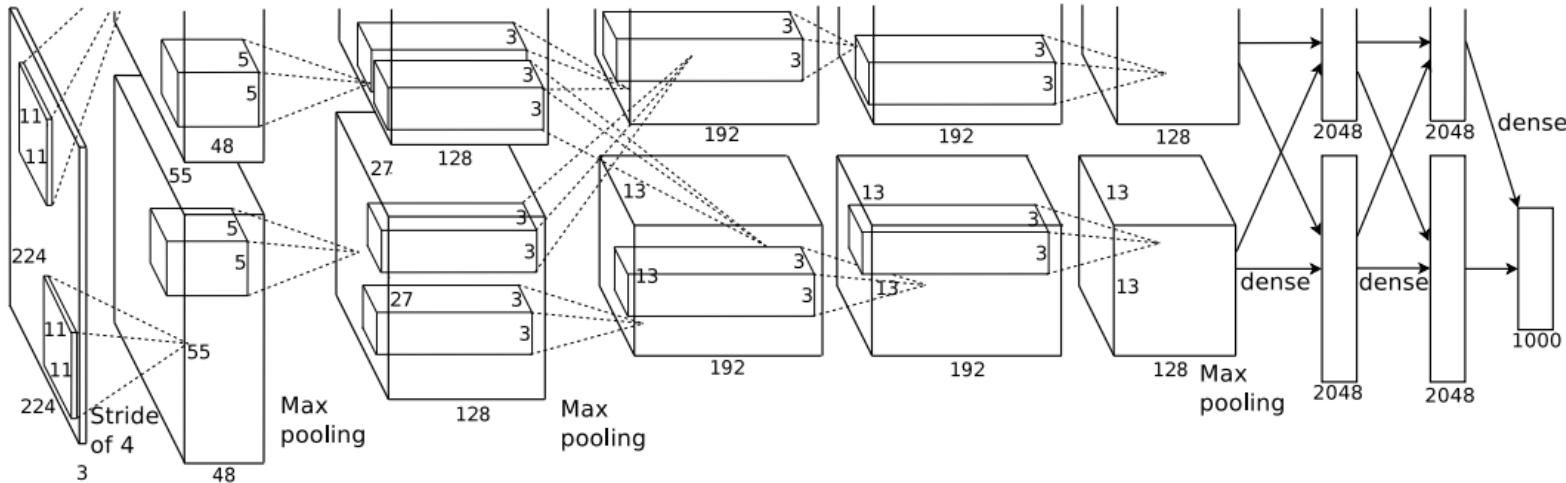
# In 2012 (and later...)

## Revolution of Depth



# Alex Krizhevsky, 2012

- AlexNet
  - Used backpropagation, SGD, ReLU
  - Used Dropout
  - Used (2x)GPU
  - Used 7 layers
  - Used ImageNet data (1.5M images, 1k classes)



---

# Convolutional

---

# Convolutional Neural

---

# Convolutional Neural Networks

# ~~Convolutional Neural Networks~~

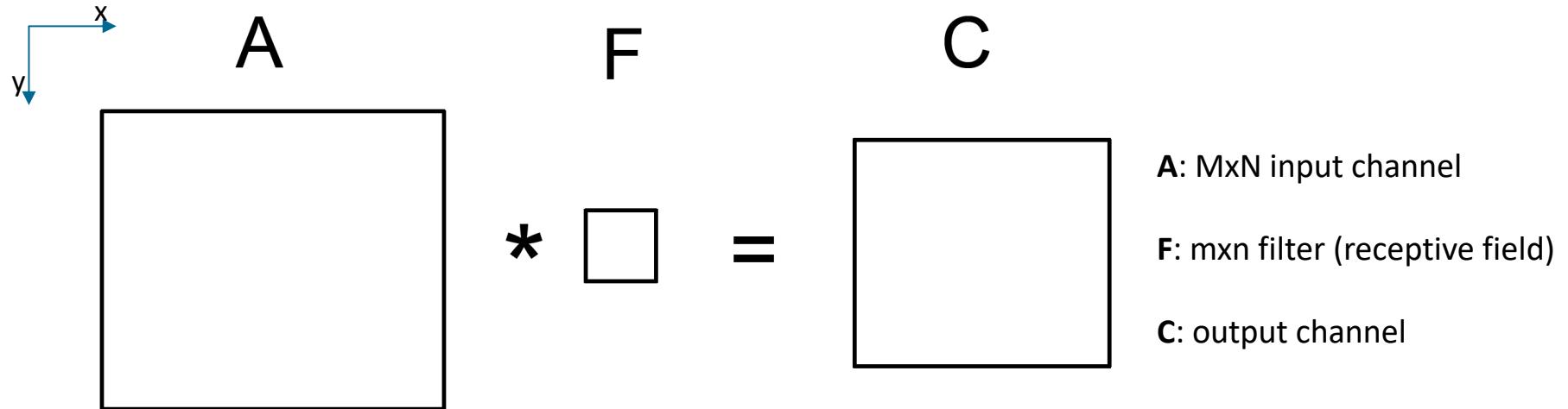
CNN  $\neq$  

ConvNets

---

# Convolutional Neural Networks

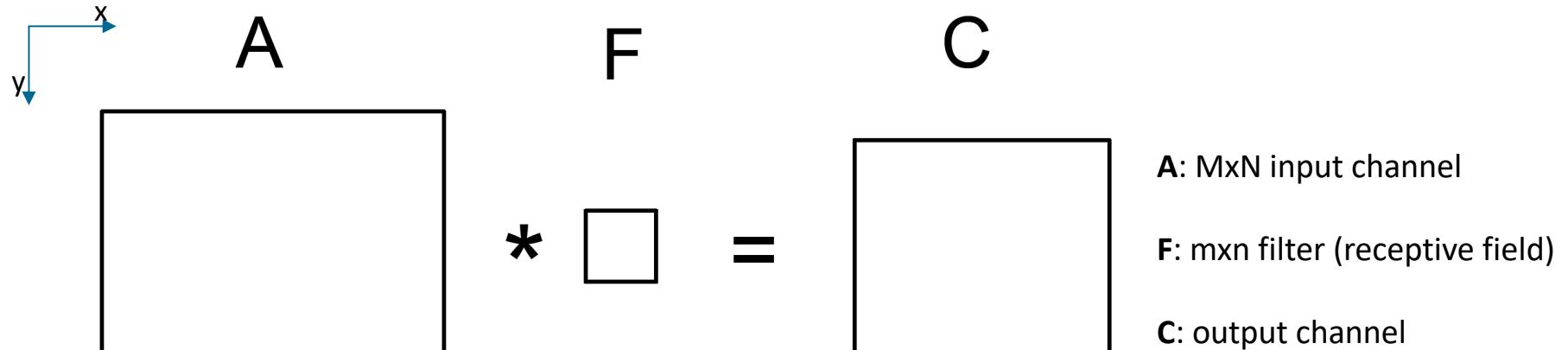
# Convolutional Neural Network



**Convolution:  $C = A * F$**

$$C(i, j) = \sum_{x=-m/2}^{m/2} \sum_{y=-n/2}^{n/2} A(i - x, j - y)F(x, y)$$

# Convolutional Neural Network



**Convolution:  $C = A * F$**

$$C(i, j) = \sum_{x=-m/2}^{m/2} \sum_{y=-n/2}^{n/2} A(i - x, j - y)F(x, y)$$

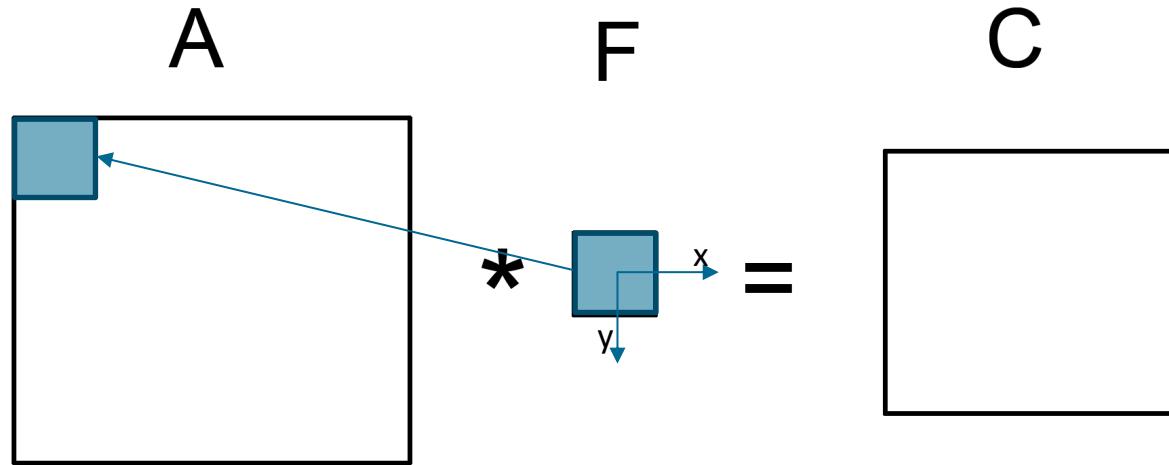
associative:  $(A * (B * F)) = ((A * B) * F)$

**Correlation**

$$C(i, j) = \sum_{x=-m/2}^{m/2} \sum_{y=-n/2}^{n/2} A(i + x, j + y)F(x, y)$$

not associative

# Convolutional Neural Network



**A:** MxN input channel

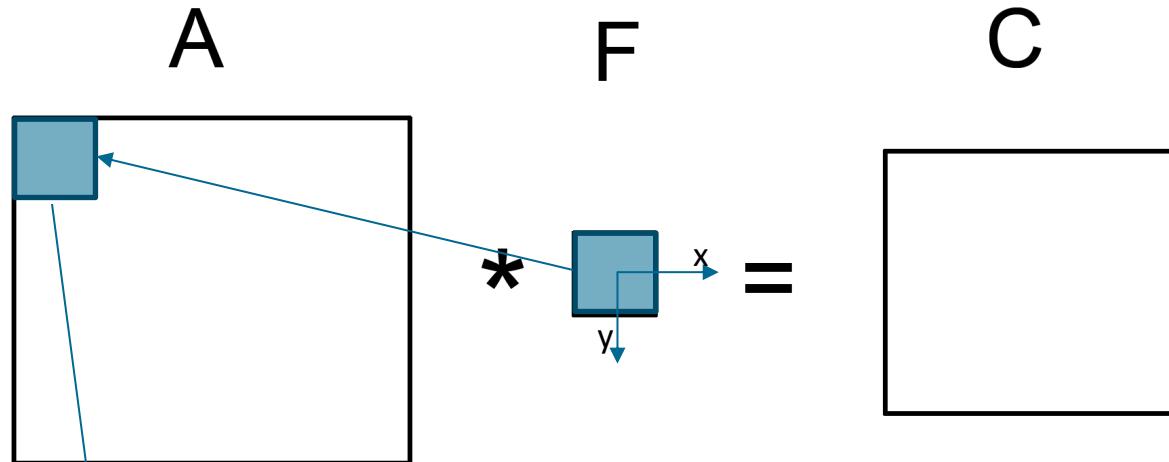
**F:** mxn filter (receptive field)

**C:** output channel

## Convolution

$$C(i, j) = \sum_{x=-m/2}^{m/2} \sum_{y=-n/2}^{n/2} A(i - x, j - y)F(x, y)$$

# Convolutional Neural Network



**A:** MxN input channel

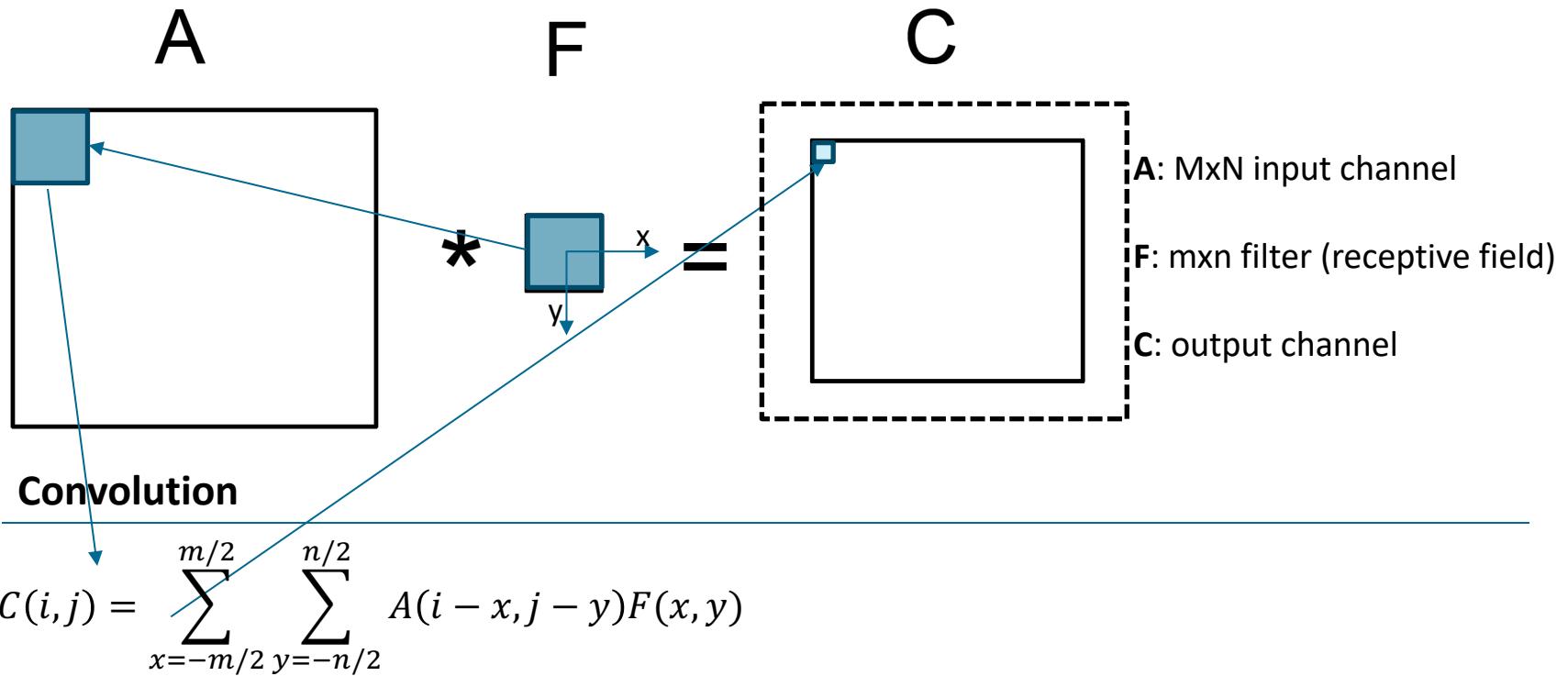
**F:** mxn filter (receptive field)

**C:** output channel

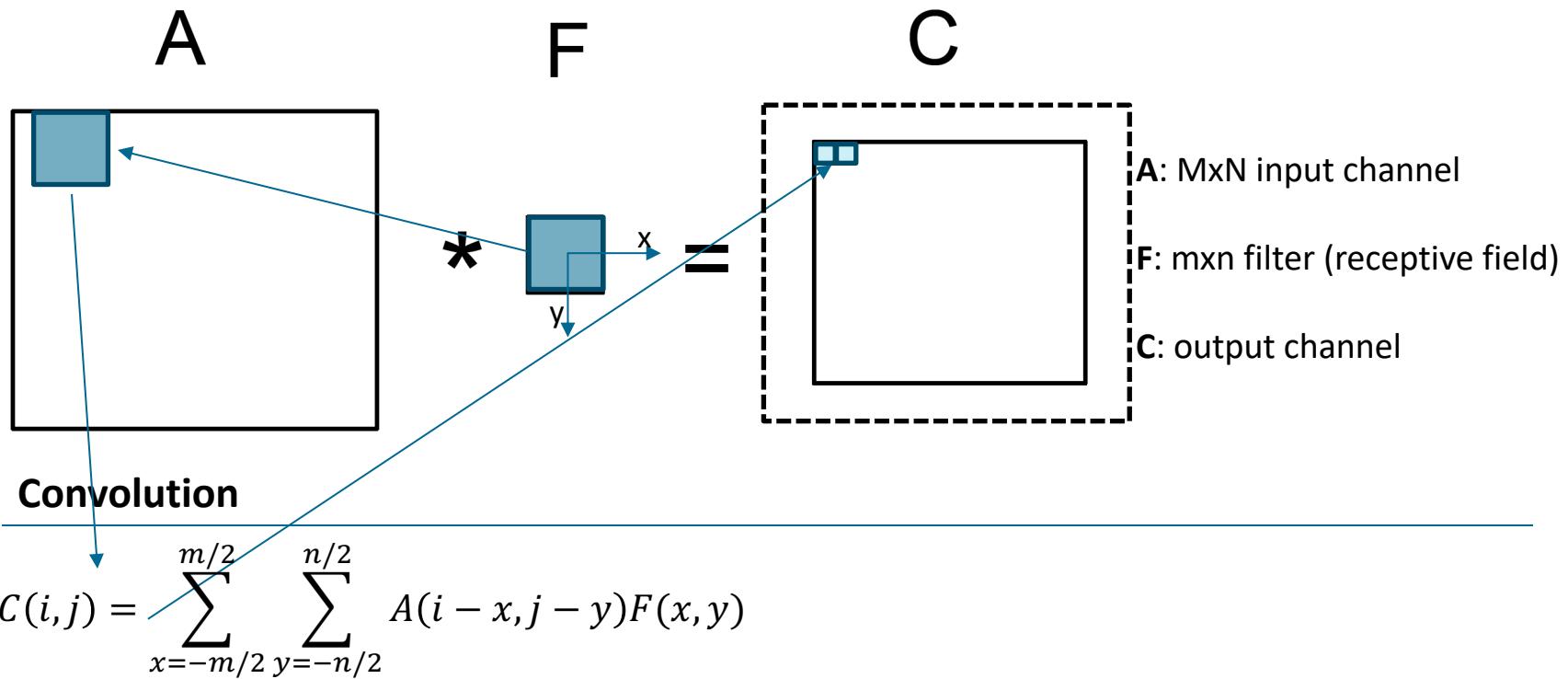
**Convolution**

$$C(i, j) = \sum_{x=-m/2}^{m/2} \sum_{y=-n/2}^{n/2} A(i - x, j - y)F(x, y)$$

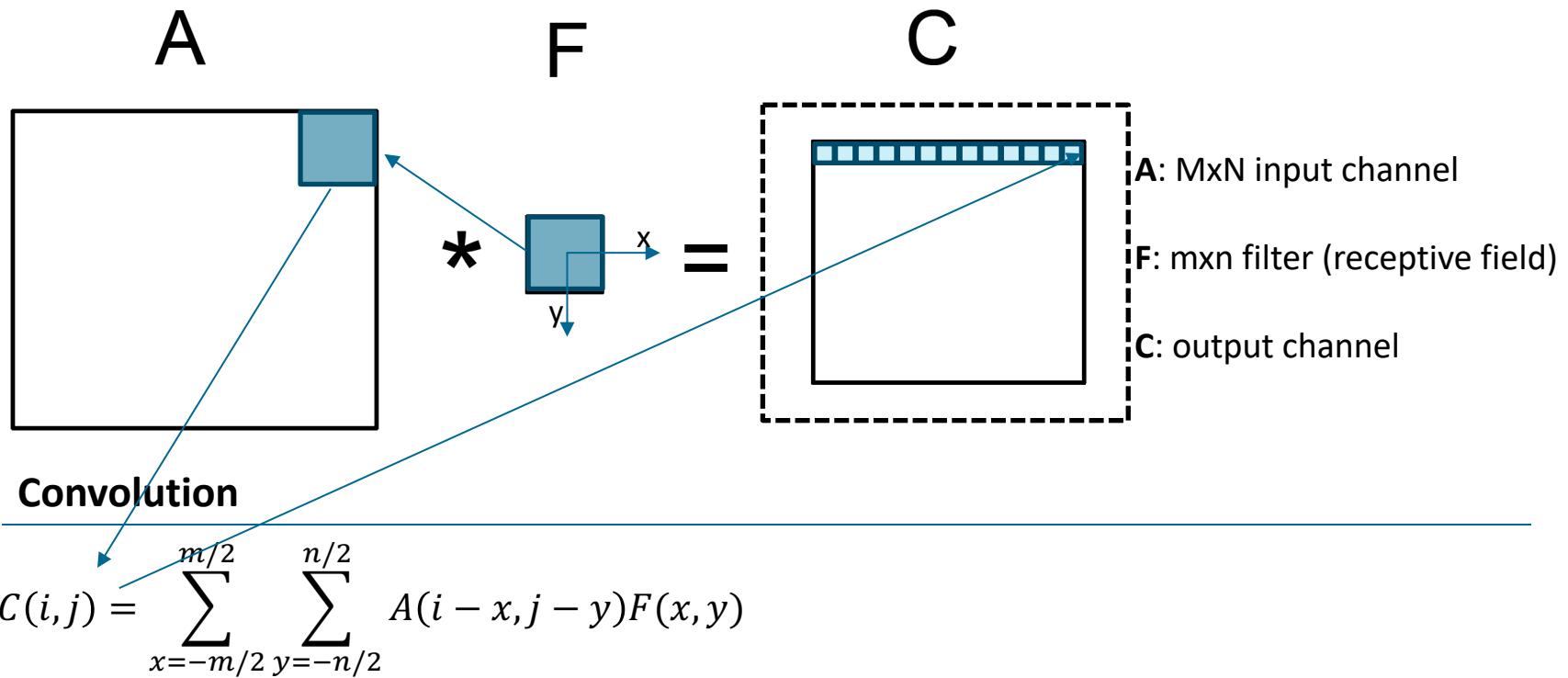
# Convolutional Neural Network



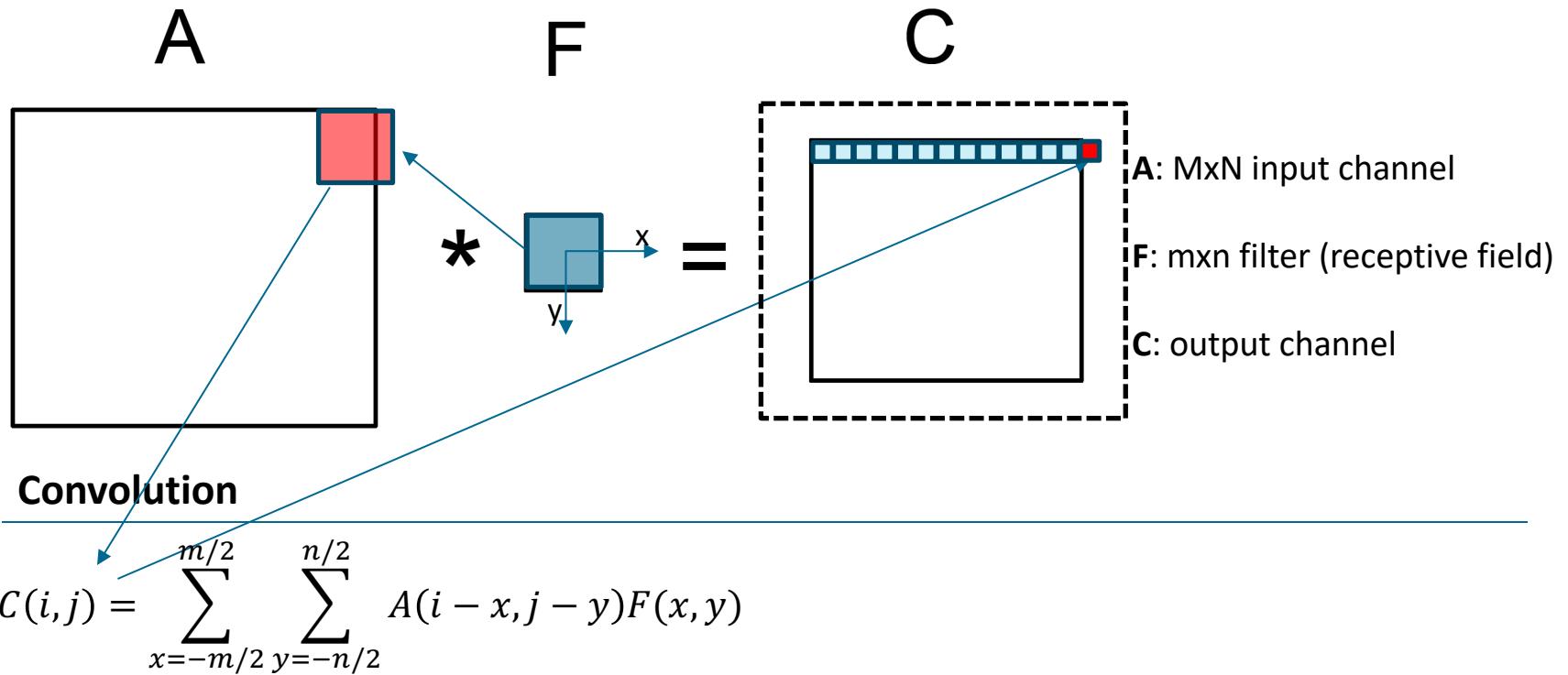
# Convolutional Neural Network



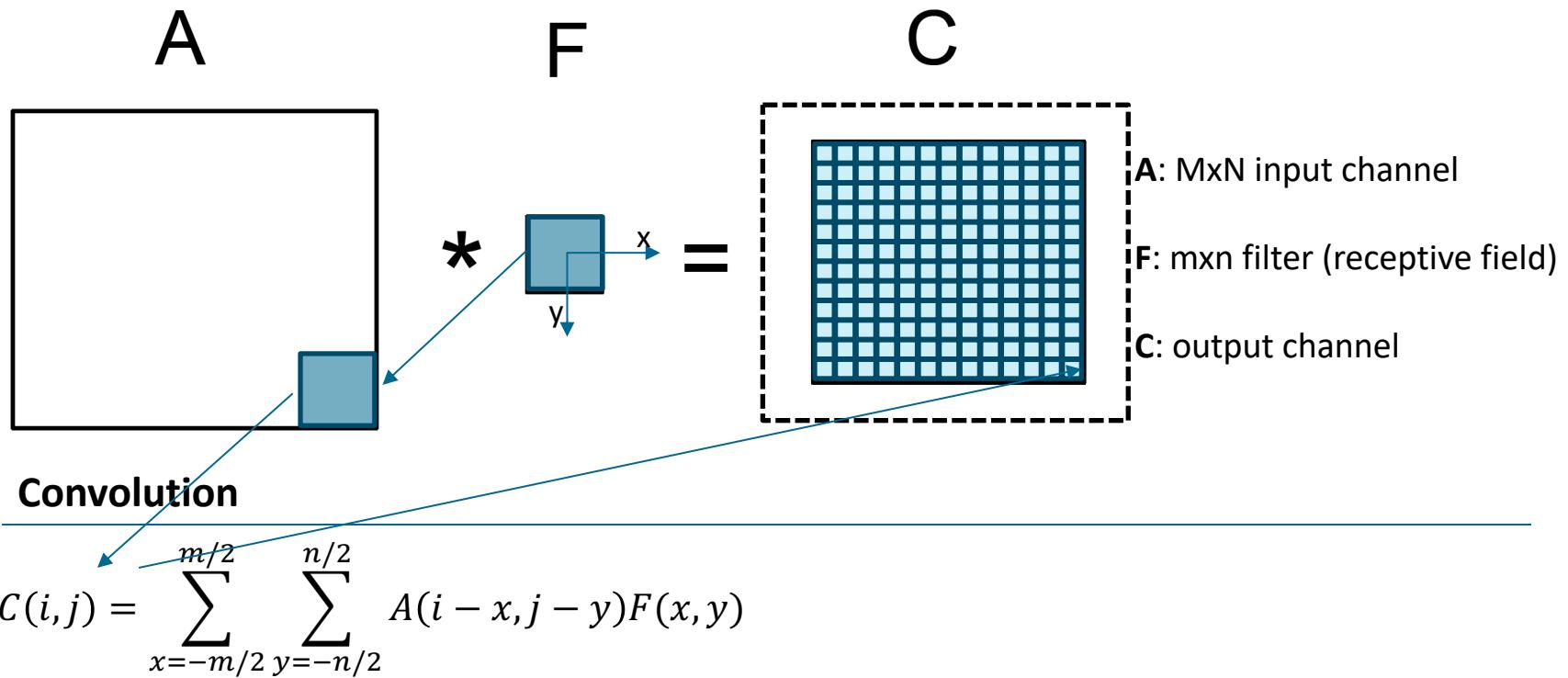
# Convolutional Neural Network



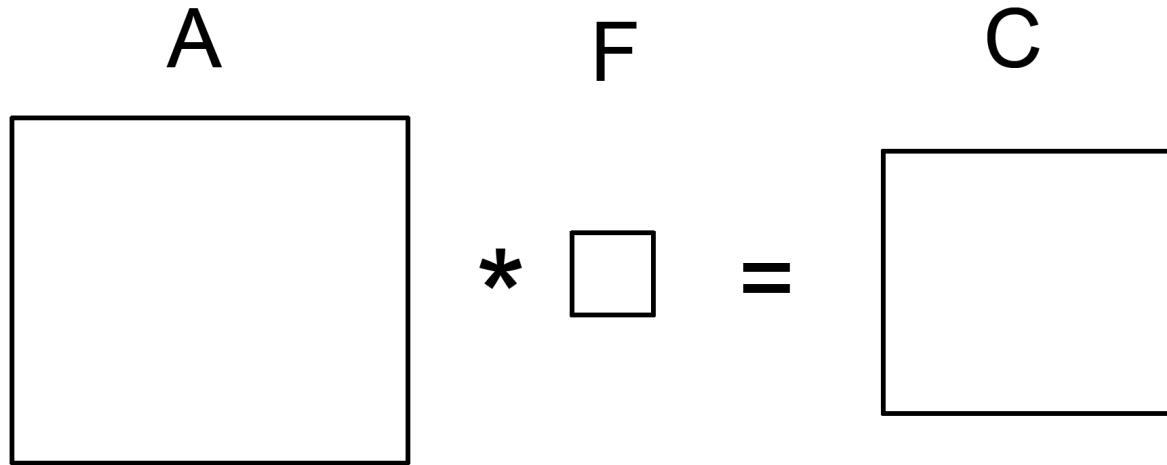
# Convolutional Neural Network



# Convolutional Neural Network



# Convolutional Neural Network



A: MxN input channel

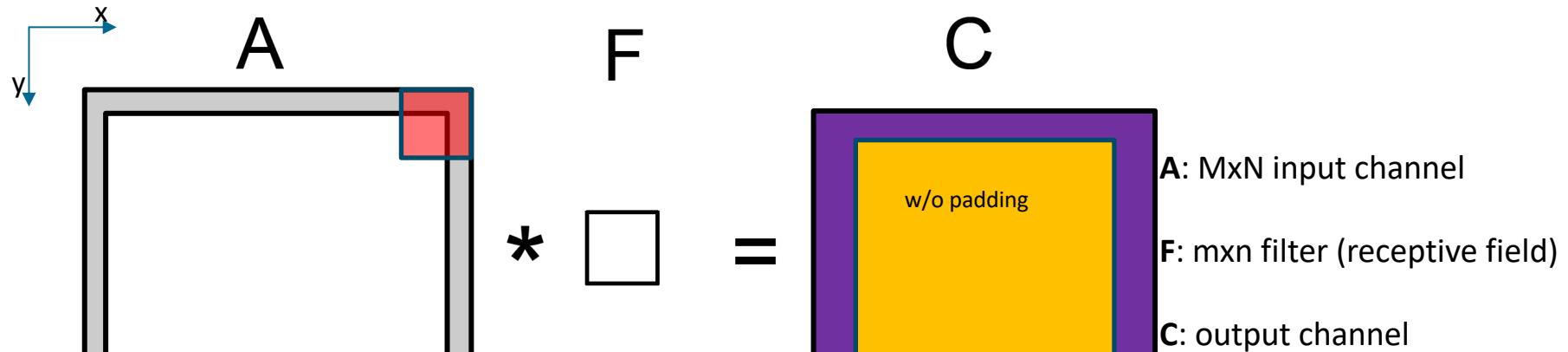
F: mxn filter (receptive field)

C: output channel

## Convolution

$$C(i, j) = \sum_{x=-m/2}^{m/2} \sum_{y=-n/2}^{n/2} A(i - x, j - y)F(x, y)$$

# Convolutional Neural Network

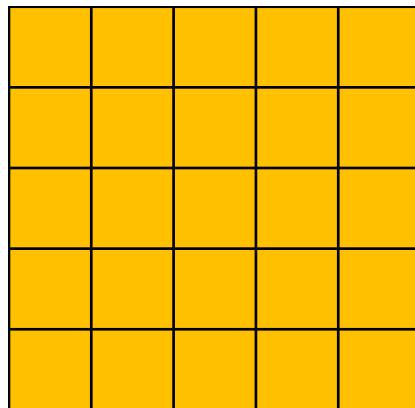


- The size of the output channel of a convolution might be **smaller** than the size of the input channel.
- In order to avoid that, we can apply **padding** to the input channel

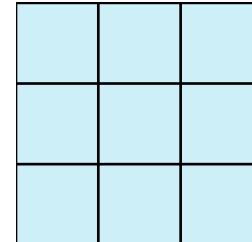
# Padding

- How much do we need to pad?

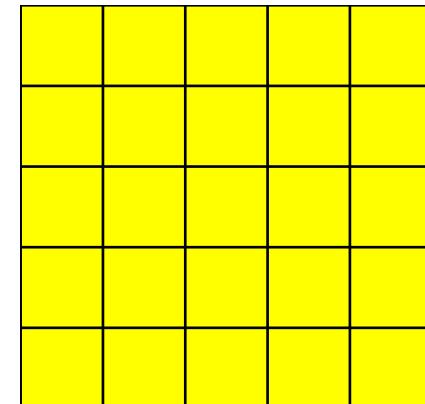
5x5 input image



3x3 filter



5x5 output



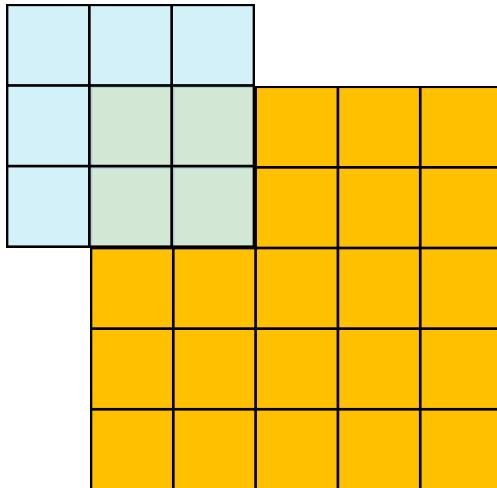
$$\text{Padding} = (3-1)/2 = 1 \text{ px}$$

“same” convolution

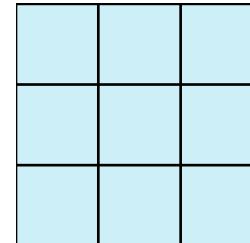
# Padding

- How much do we need to pad?

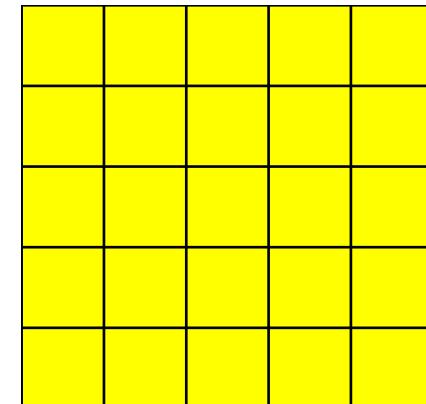
5x5 input image



3x3 filter



5x5 output

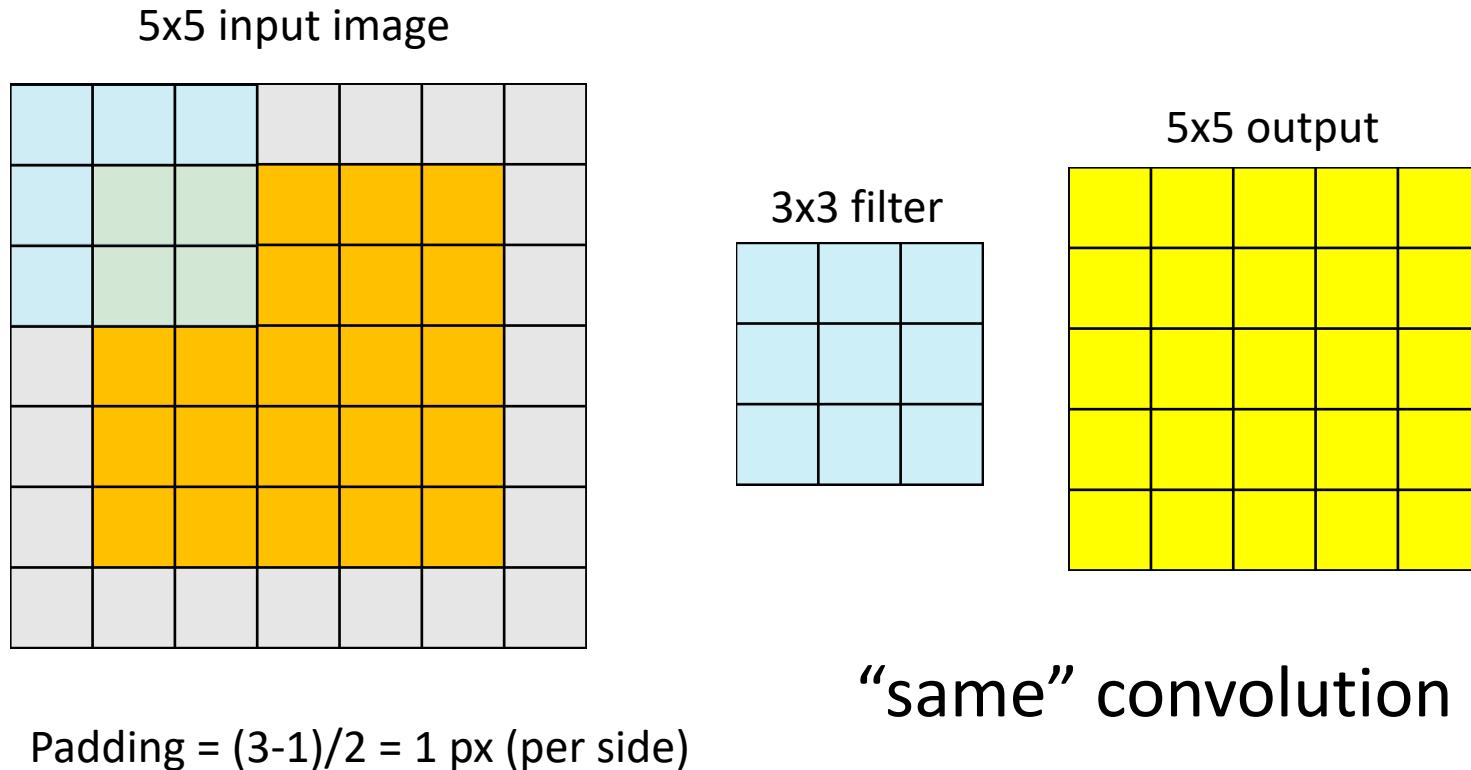


“same” convolution

$$\text{Padding} = (3-1)/2 = 1 \text{ px (per side)}$$

# Padding

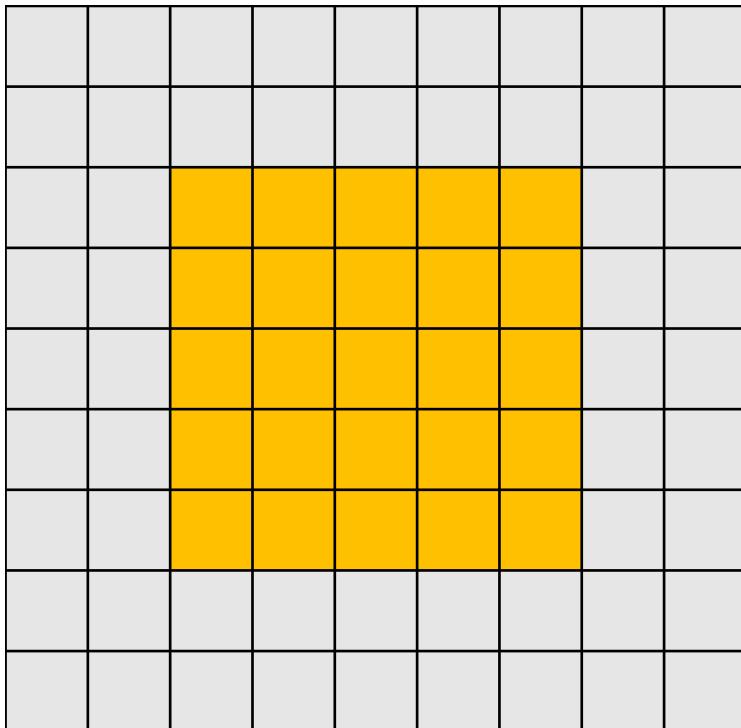
- How much do we need to pad?



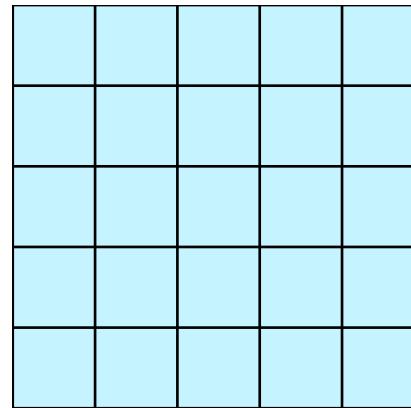
# Padding

- How much do we need to pad?

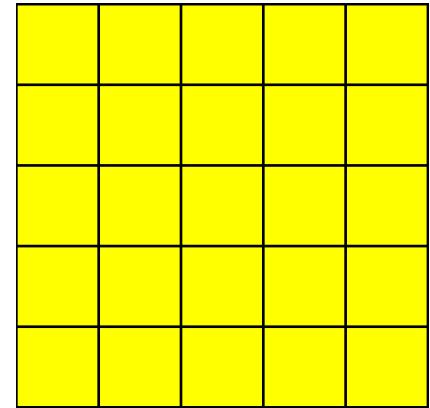
5x5 input image



5x5 filter



5x5 output



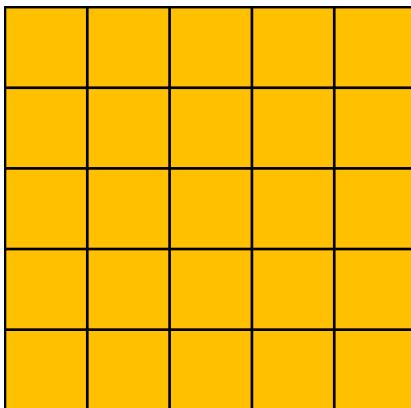
“same” convolution

$$\text{Padding} = (5-1)/2 = 2 \text{ px (per side)}$$

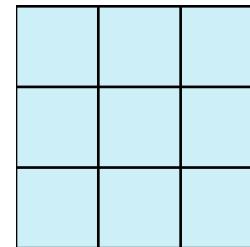
# “Valid” convolution

- How much do we need to pad?

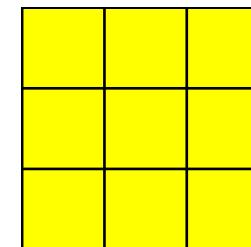
5x5 input image



3x3 filter



3x3 output

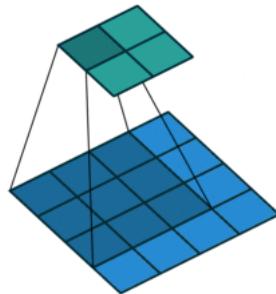


$$\text{Padding} = (3-1)/2 = 1 \text{ px}$$

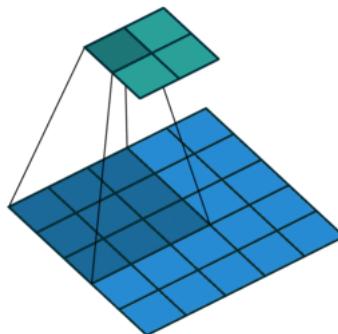
- **padding:** One of `"valid"`, `"causal"` or `"same"` (case-insensitive). `"valid"` means "no padding". `"same"` results in padding the input such that the output has the same length as the original input. `"causal"` results in causal (dilated) convolutions, e.g. `output[t]` does not depend on `input[t+1:]`. Useful when modeling temporal data where the model should not violate the temporal order. See [WaveNet: A Generative Model for Raw Audio, section 2.1](#).

# Stride

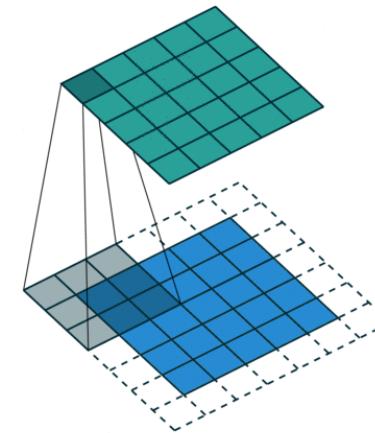
- Stride is the **step** that you make when **applying convolutions**
- By changing the stride you can:
  - change the size of the output by keeping the same filter size
  - cover a larger input area with the same convolutions



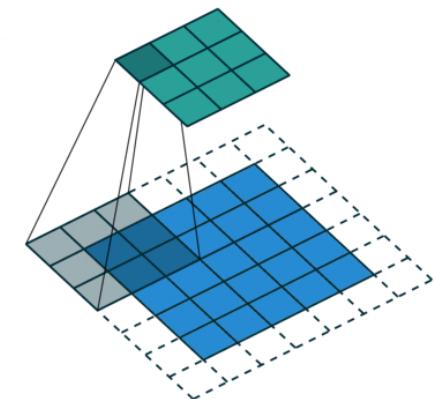
no padding, no stride



no padding, stride



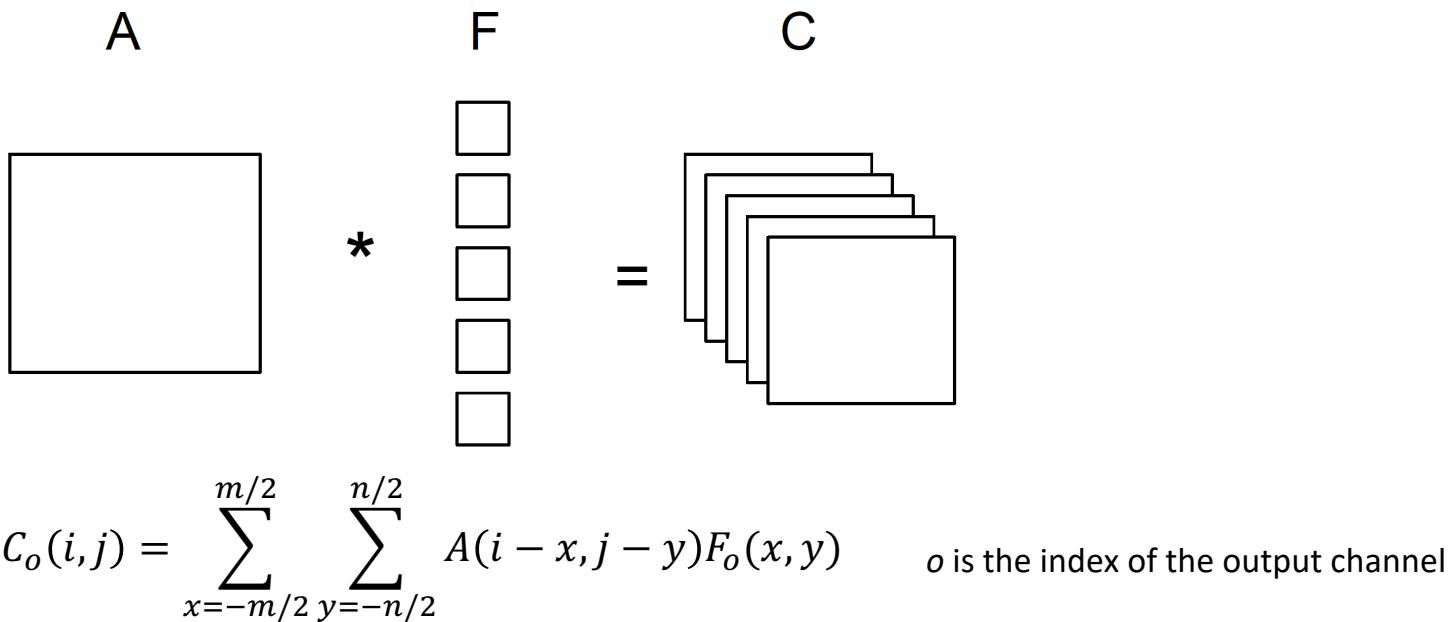
padding, no stride



padding, stride

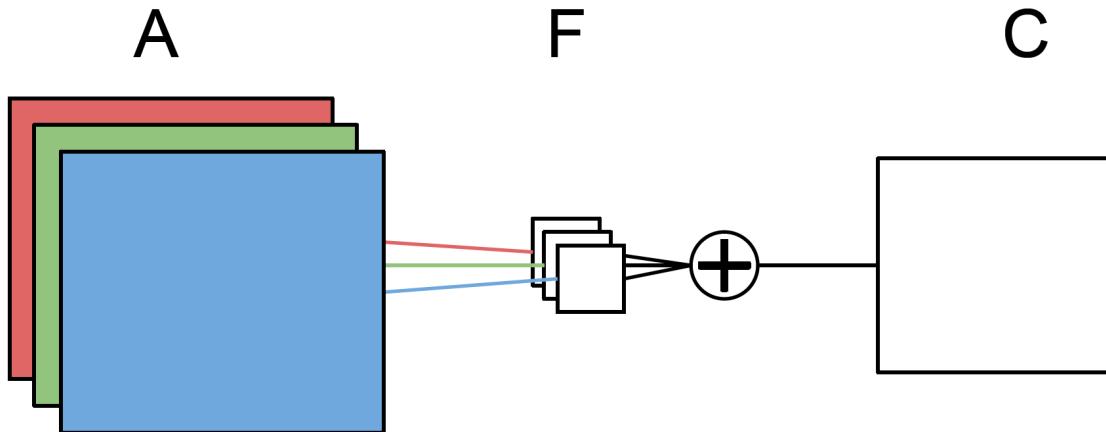
# Convolutional Neural Network

- Applying **multiple filters** we obtain **multiple output channels**



# Convolutional Neural Network

- Convolution with **multiple input channels** (example: RGB images)



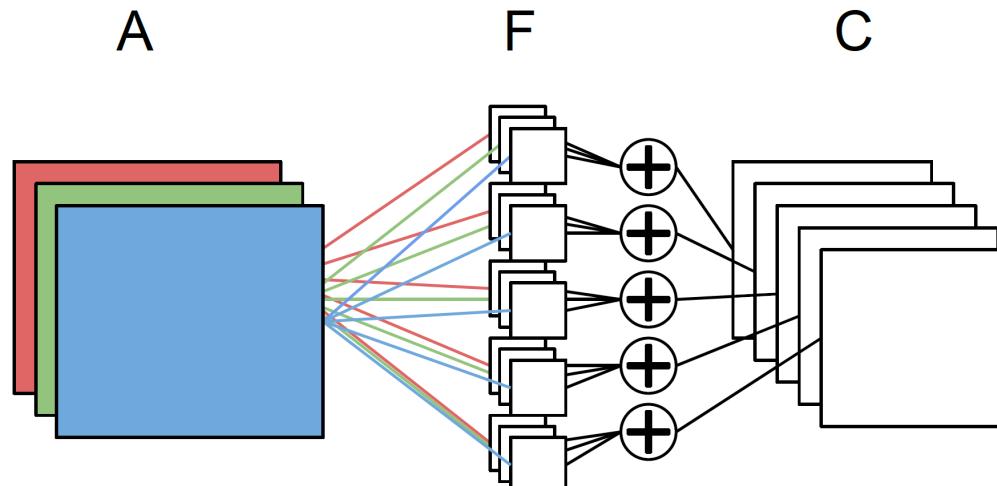
$$C(i, j) = \sum_{k=0}^l \sum_{x=-m/2}^{m/2} \sum_{y=-n/2}^{n/2} A_k(i - x, j - y) F_k(x, y)$$

F is now a set of filters  
k is the index of the input channel

One example of  
**multiple channels** is  
also the **output** of the  
previous slide

# Convolutional Neural Network

- Convolution with **multiple inputs** and **outputs**

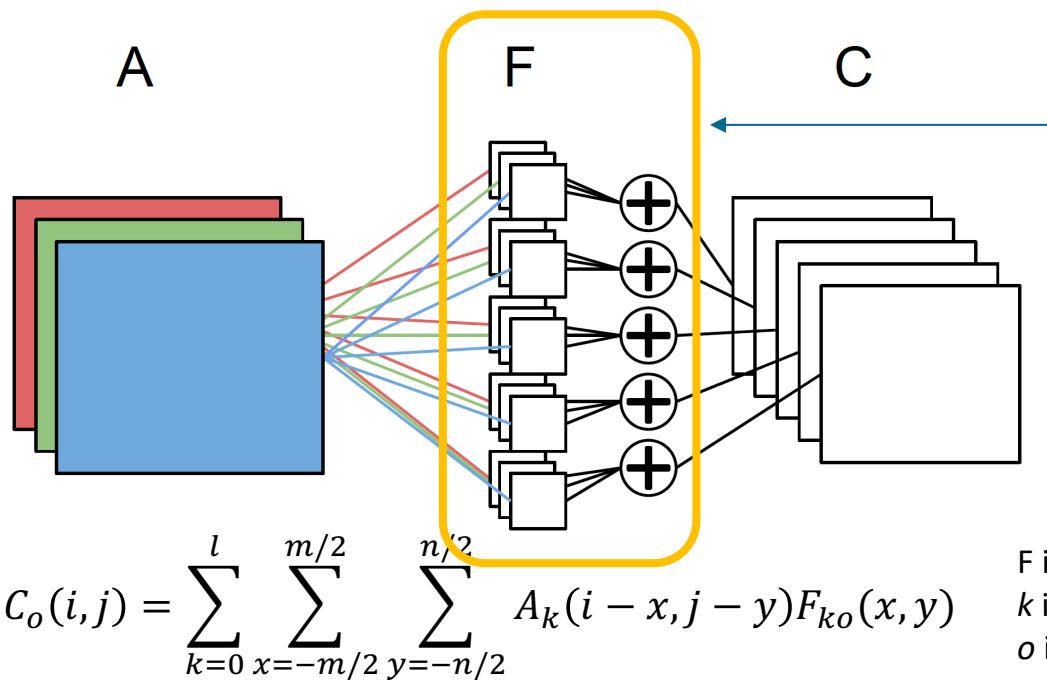


$$C_o(i, j) = \sum_{k=0}^l \sum_{x=-m/2}^{m/2} \sum_{y=-n/2}^{n/2} A_k(i - x, j - y) F_{ko}(x, y)$$

F is now a set of filters  
k is the index of the input channel  
o is the index of the output channel

# Convolutional Neural Network

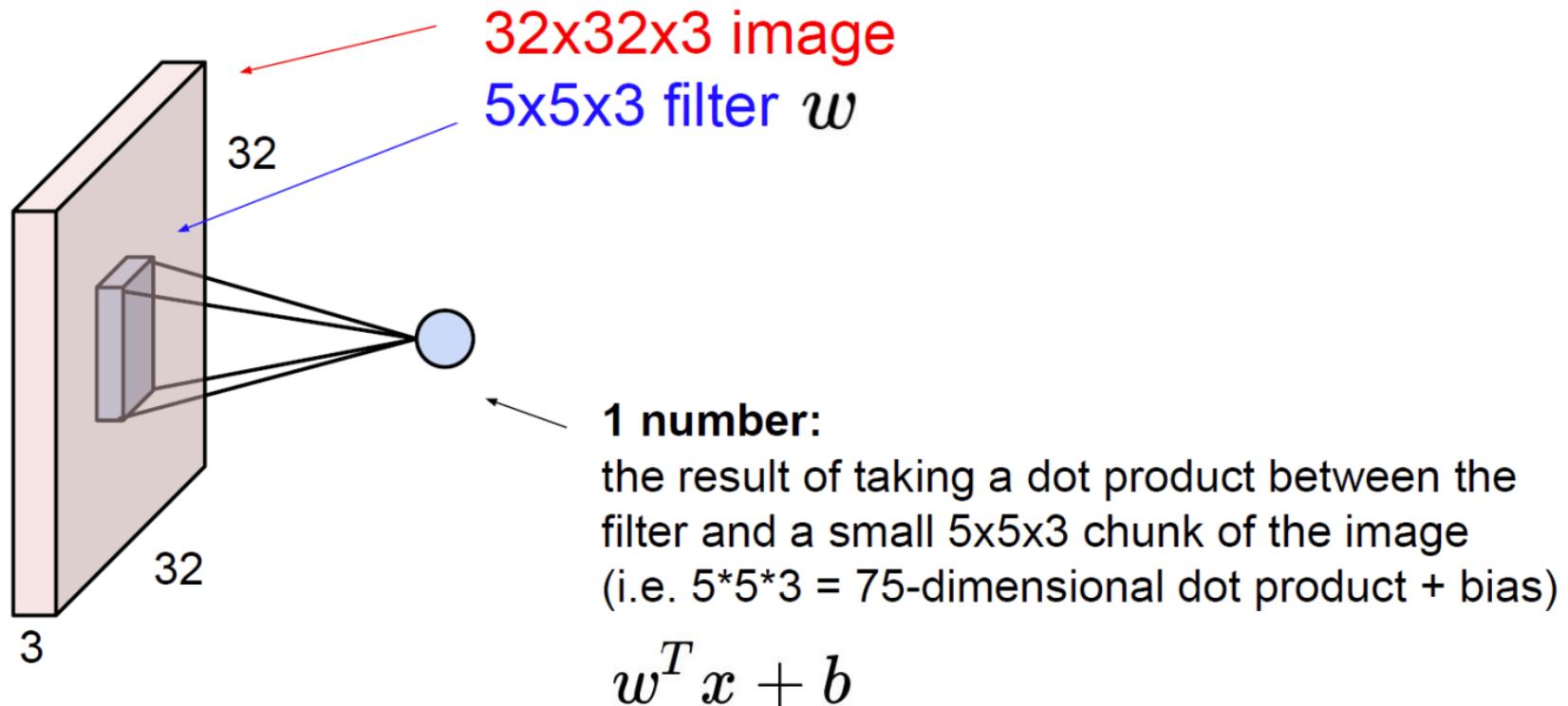
- Convolution with **multiple inputs** and **outputs**



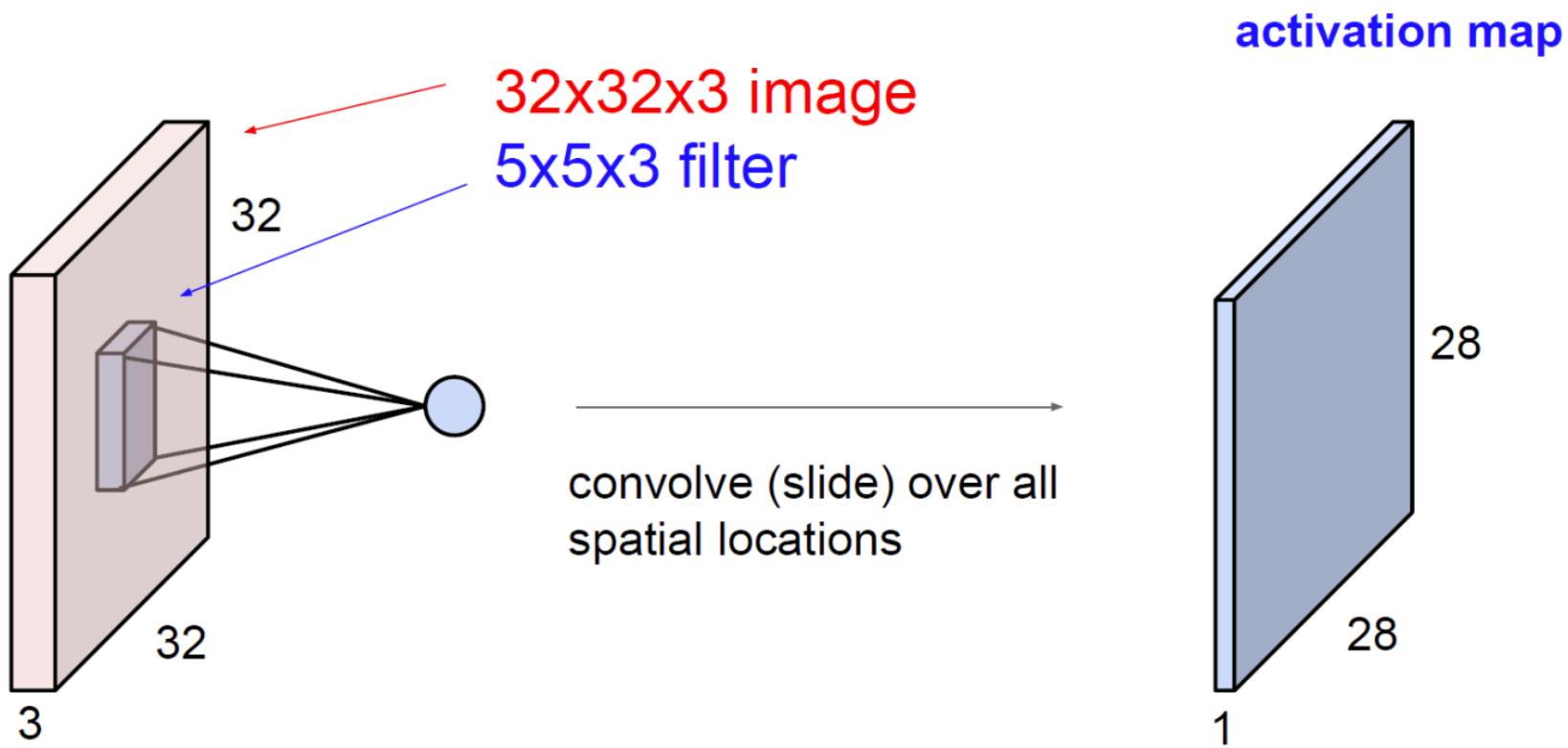
We can call this a  
**Convolutional layer**

$F$  is now a set of filters  
 $k$  is the index of the input channel  
 $o$  is the index of the output channel

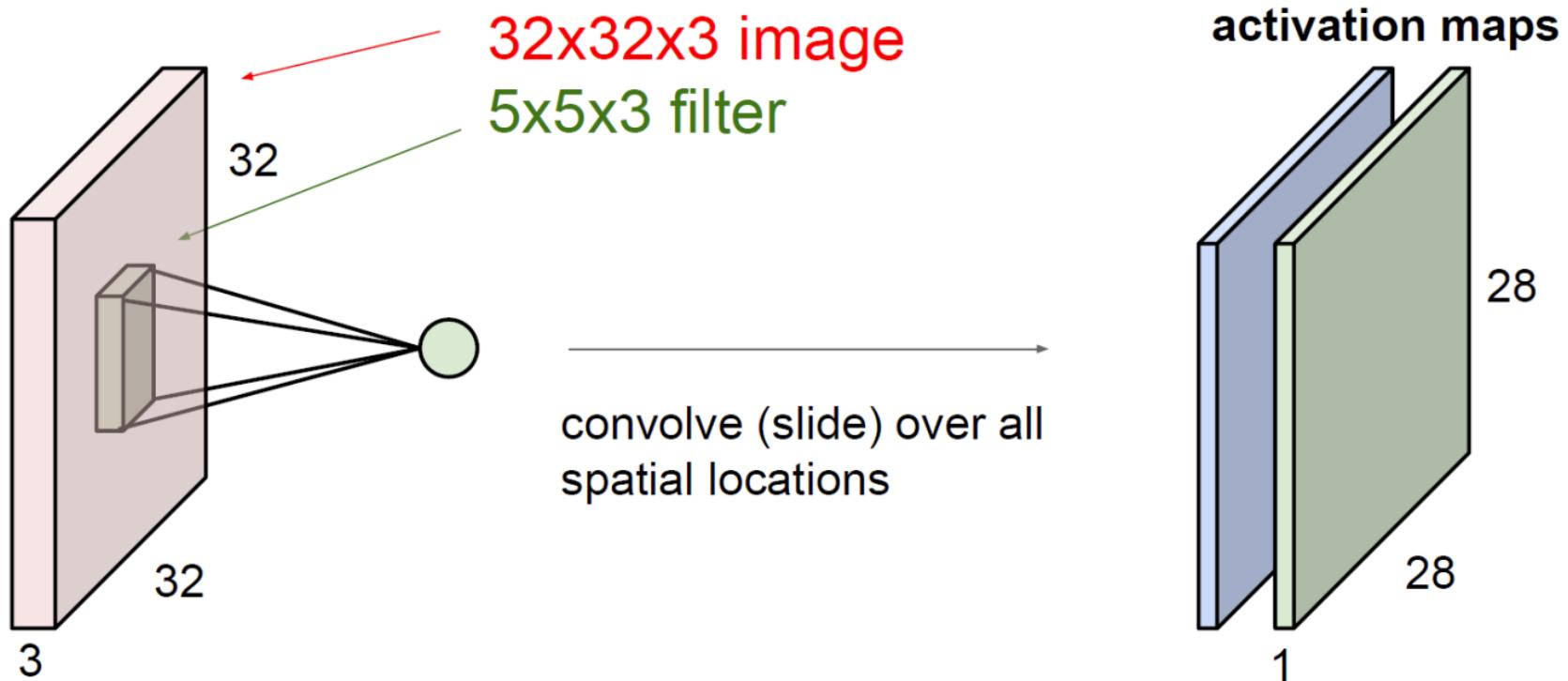
# Convolutional Neural Network



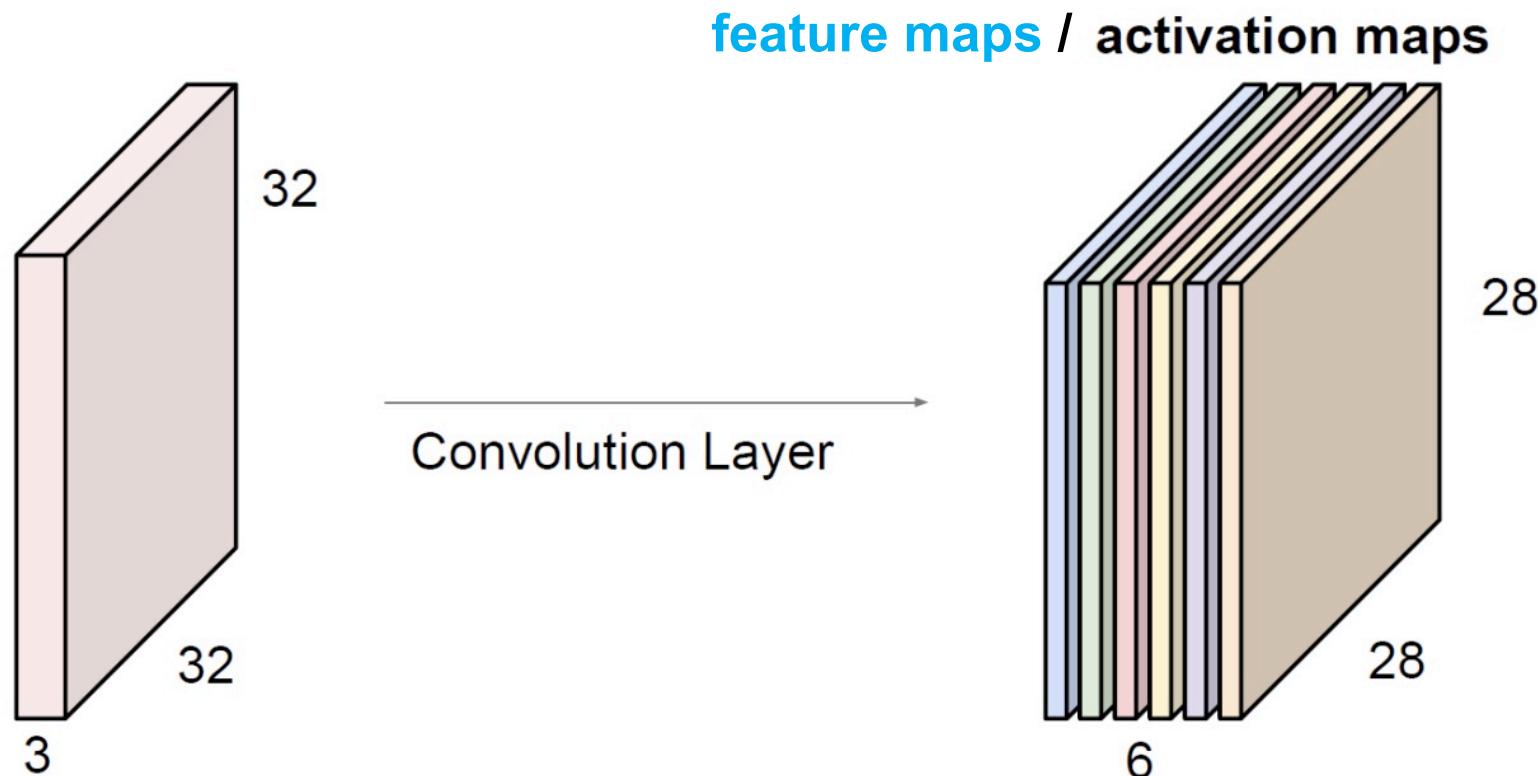
# Convolutional Neural Network



# Convolutional Neural Network



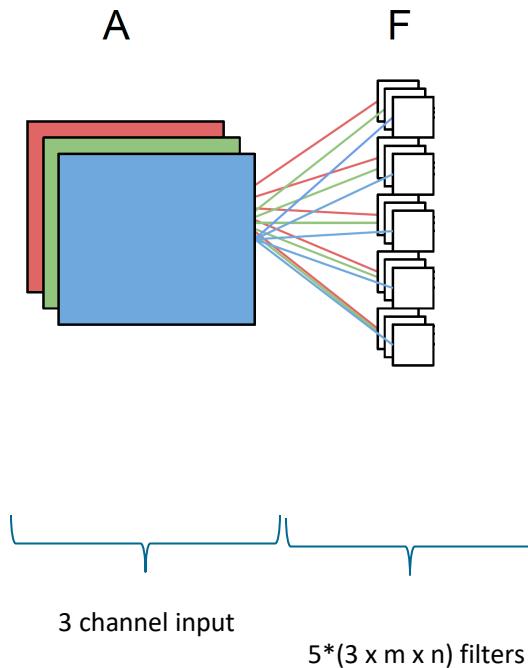
# Convolutional Neural Network



We stack these up to get a “new image” of size 28x28x6!

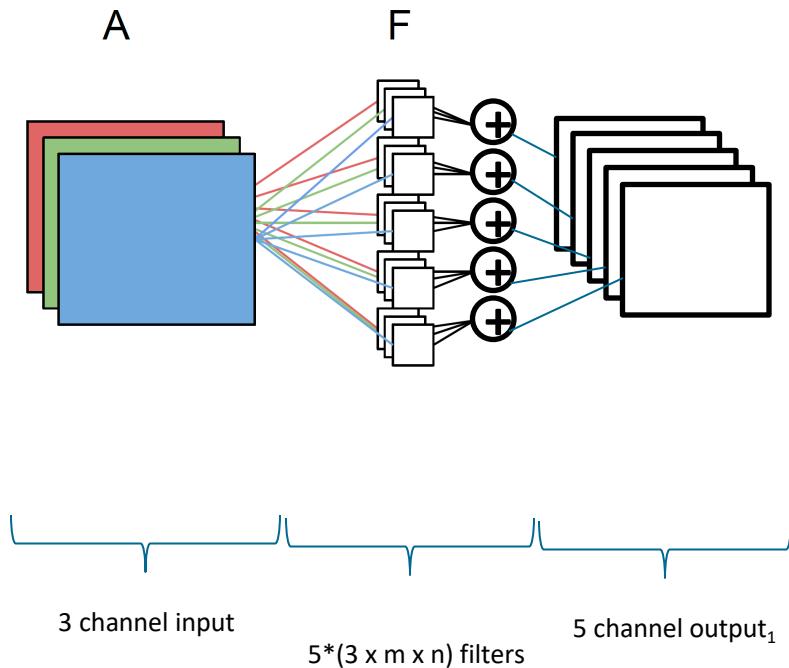
# Convolutional Neural Network

- A convolutional network is (basically) a stack of convolutional layers



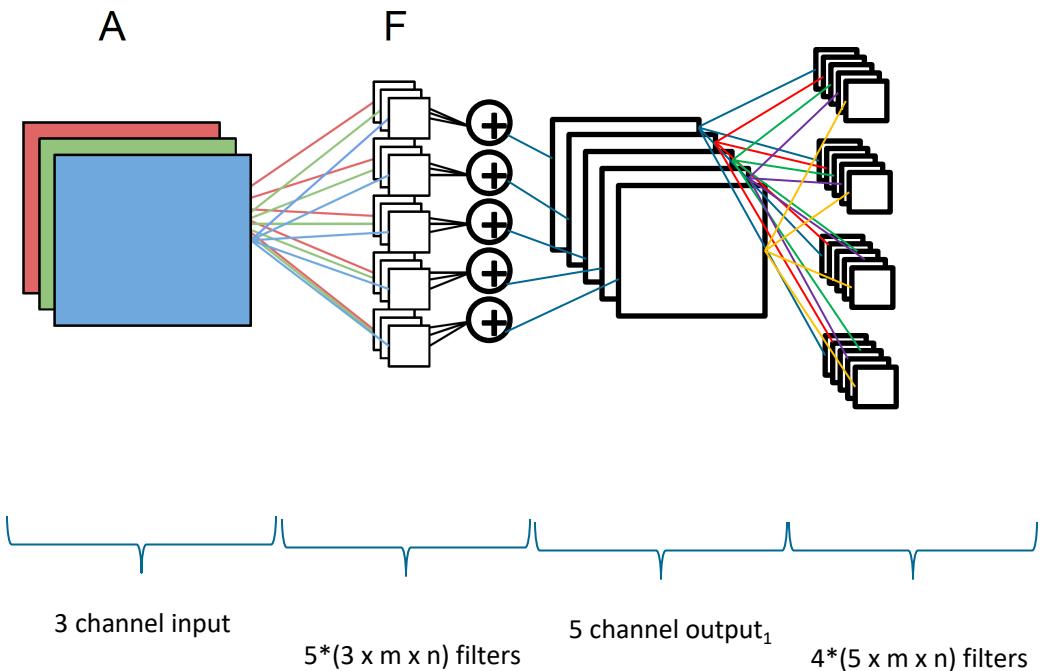
# Convolutional Neural Network

- A convolutional network is (basically) a stack of convolutional layers



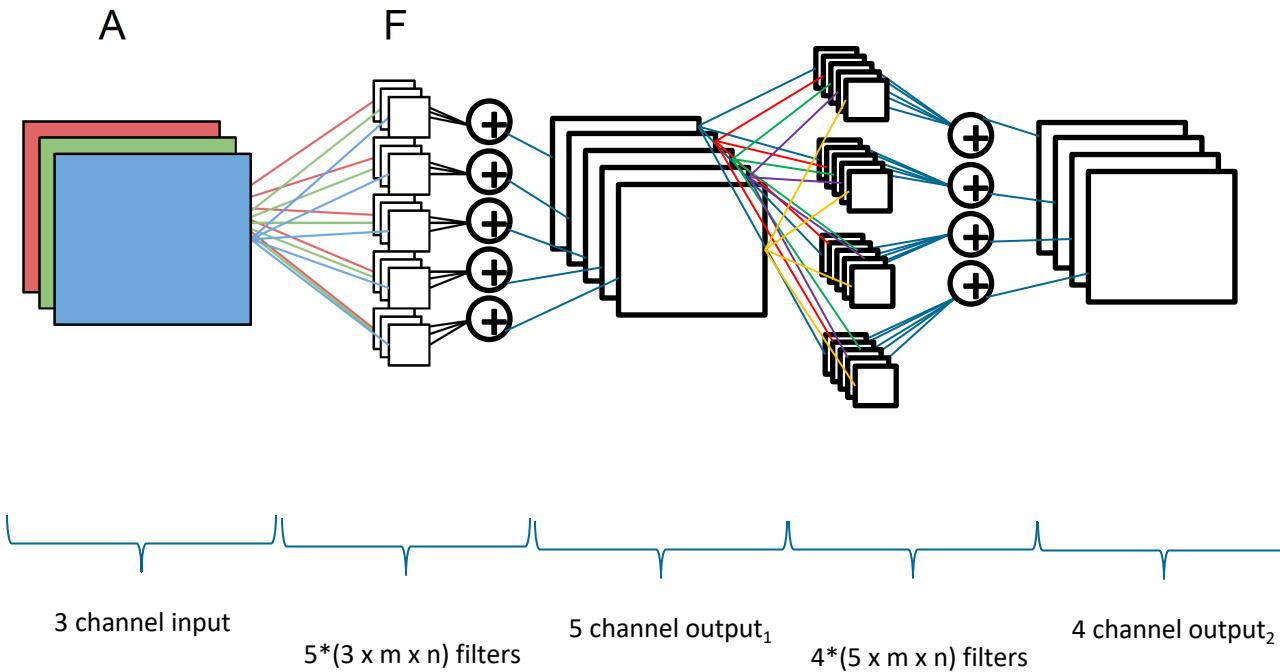
# Convolutional Neural Network

- A convolutional network is (basically) a stack of convolutional layers



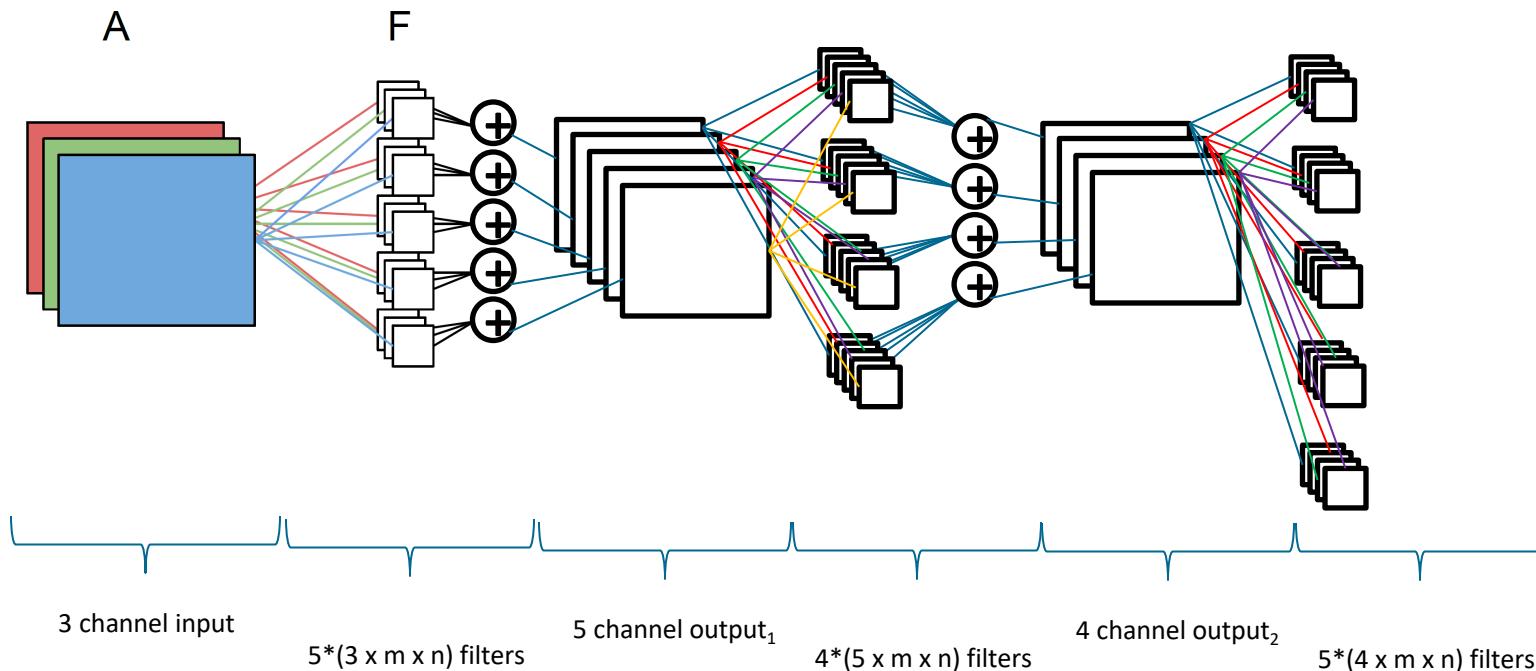
# Convolutional Neural Network

- A convolutional network is (basically) a stack of convolutional layers



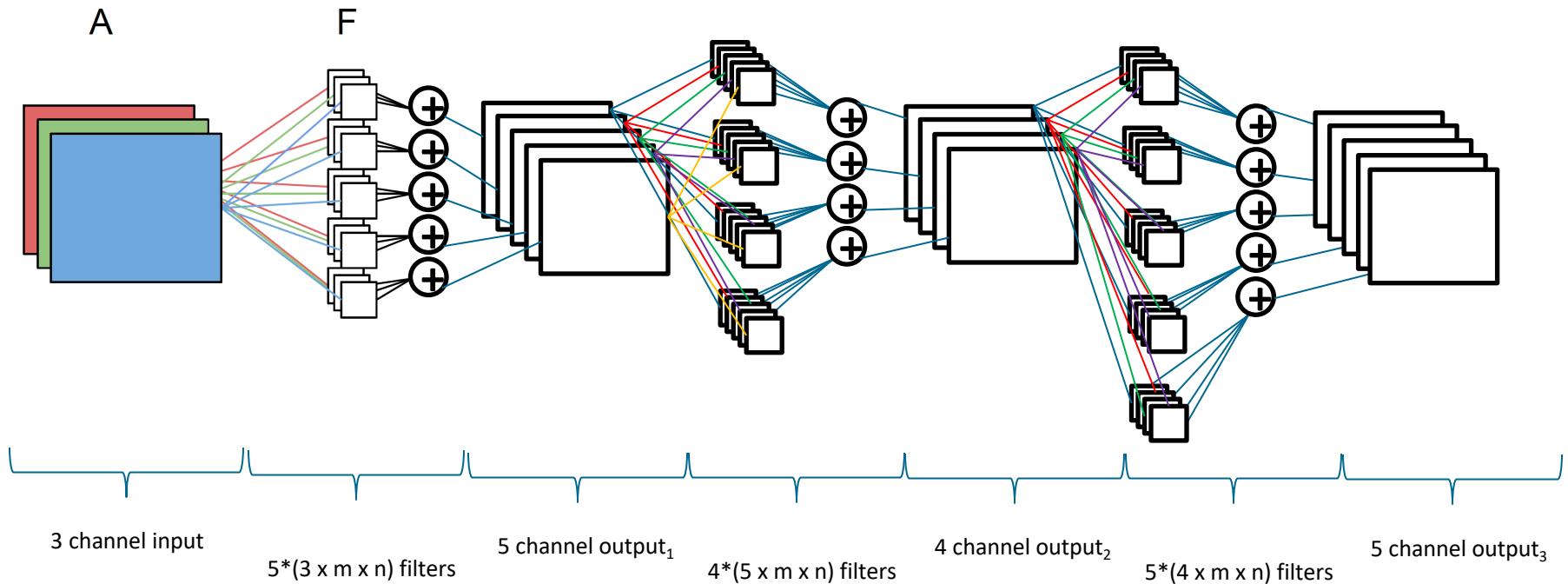
# Convolutional Neural Network

- A convolutional network is (basically) a stack of convolutional layers



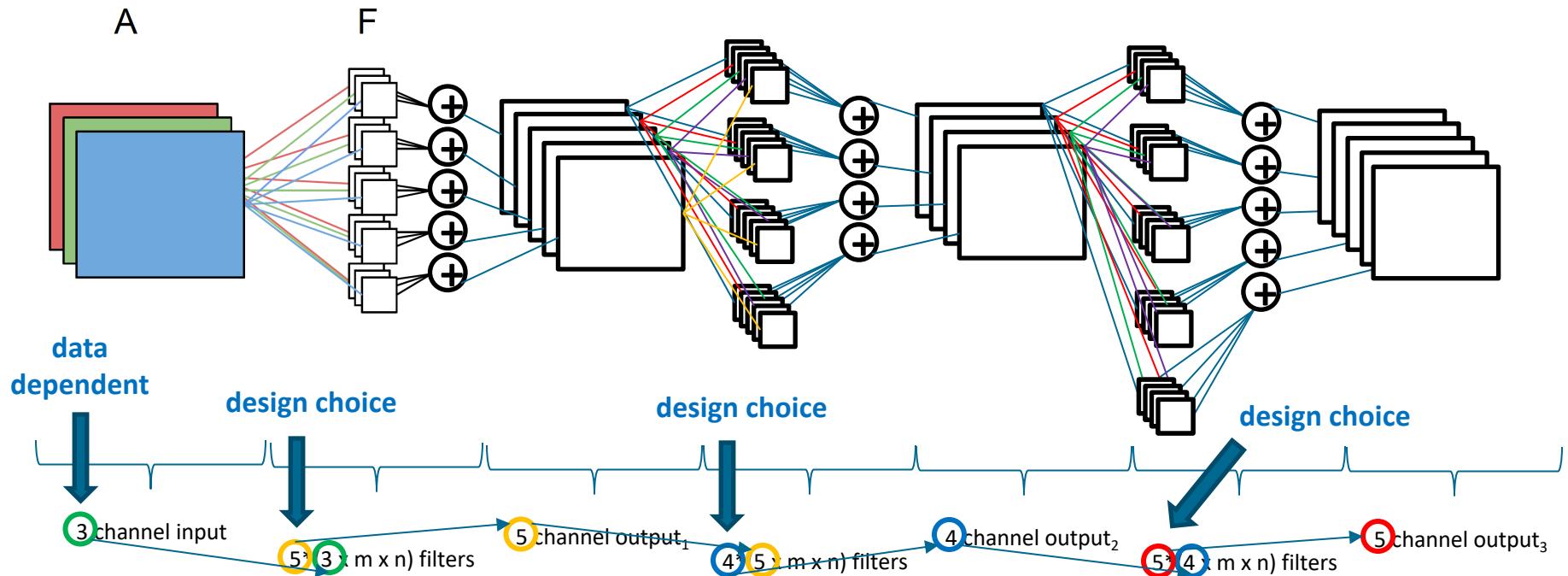
# Convolutional Neural Network

- A convolutional network is (basically) a stack of convolutional layers



# Convolutional Neural Network

- A convolutional network is (basically) a stack of convolutional layers



# VIBOT 2018

## Machine learning with Neural Networks

**Francesco Ciompi**

[francesco.ciompi@radboudumc.nl](mailto:francesco.ciompi@radboudumc.nl)