

# Understanding Machine Learning & Neural Network

A Hardware Perspective

Chixiao Chen

# Overview

- Linear Classifier
  - Score function
  - Loss function and softmax classifier
  - Gradient Descent
- Neural Network

# Discussion on HW 2

- Problem 2

Assuming you are asked to design a specific processor to implementing matrix-vector multiplication based on the basic RV32 architecture. In other words, preform the following computing in a single cycle processor,

$$\underbrace{\begin{bmatrix} 1 & 2 & 3 \end{bmatrix}}_{1 \times 3} \cdot \underbrace{\begin{bmatrix} 2 & 1 & 3 \\ 3 & 3 & 2 \\ 4 & 1 & 2 \end{bmatrix}}_{3 \times 3} = \underbrace{\begin{bmatrix} 1 \cdot 2 + 2 \cdot 3 + 3 \cdot 4 \\ 1 \cdot 1 + 2 \cdot 3 + 3 \cdot 1 \\ 1 \cdot 3 + 2 \cdot 2 + 3 \cdot 2 \end{bmatrix}}_{1 \times 3} = \begin{bmatrix} 20 \\ 10 \\ 13 \end{bmatrix}$$

Figure 2: A matrix-vector product example

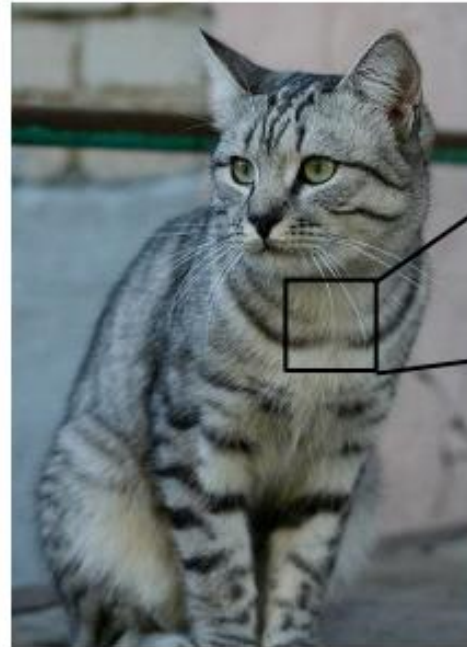
Register `s0` is the base pointer to save the following sequence  $\{1, 2, 3, 2, 3, 4, 1, 3, 1, 3, 2, 2\}$  in words (32bit). Register `t0` is the base pointer to store the output results. For the multiplication, please refer to the RV32M extension in appendix I. For simplicity, we only use the unsigned multiplication instruction `MUL` here.

- (a) Write assembly codes to implement the matrix-vector product based on RV32I/RV32M instruction set. (hint: please use branch to minimize the code length.)
- (b) Tailor a specific instruction set which only needs to support the above code. Use a table to illustrate all the instructions and their opcode/operand fields.
- (c) Write an RTL code to realize all the instructions used in the above table. (hint: no need for pipeline here, and reuse the code in HW#1.)

# Linear Classifier

# Classification

- The most common and simple task in ML
- Problem: from human to data



This image by Nikita is  
licensed under [CC-BY 2.0](#)

(assume given set of discrete labels)  
{dog, cat, truck, plane, ...}

[	105	112	100	111	104	90	106	99	96	103	112	119	104	97	93	87]
[	91	98	102	106	104	79	98	103	99	105	123	136	110	105	94	85]
[	76	85	90	105	128	105	87	96	95	99	115	112	106	103	99	85]
[	99	81	81	93	120	131	127	100	95	98	102	99	96	93	101	94]
[	106	91	61	64	69	91	88	85	101	107	109	98	75	84	96	95]
[	114	108	85	55	55	69	64	54	64	87	112	129	98	74	84	91]
[	133	137	147	103	65	81	80	65	52	54	74	84	102	93	85	82]
[	128	137	144	140	109	95	86	70	62	65	63	63	60	73	86	101]
[	125	133	140	137	119	121	117	94	65	79	80	65	54	64	72	98]
[	127	125	131	147	133	127	126	131	111	96	89	75	61	64	72	84]
[	115	114	109	123	150	140	131	118	113	109	100	92	74	65	72	78]
[	89	93	90	97	108	147	131	118	113	114	113	100	106	95	77	80]
[	63	77	86	81	77	79	102	123	117	115	117	125	125	130	115	87]
[	62	65	82	89	78	71	80	101	124	126	119	101	107	114	131	110]
[	63	65	75	88	89	71	62	81	120	138	135	105	81	90	110	118]
[	87	65	71	87	106	95	69	45	76	130	126	107	92	94	105	112]
[	118	97	82	86	117	123	116	66	41	51	95	93	89	95	102	107]
[	104	146	112	80	82	120	124	104	76	48	45	66	88	101	102	109]
[	157	170	157	120	93	86	114	132	112	97	69	55	70	82	99	94]
[	130	120	134	161	139	100	109	118	121	134	114	87	65	53	69	86]
[	120	112	96	117	150	144	120	115	104	107	102	93	87	81	72	79]
[	123	107	96	86	83	112	153	149	122	109	104	75	80	107	112	99]
[	122	121	102	80	82	86	94	117	145	148	153	102	50	78	92	107]
[	122	164	148	103	71	56	78	83	93	103	119	139	102	61	60	84]]

What the computer sees

An image is just a big grid of  
numbers between [0, 255]:

e.g. 800 x 600 x 3  
(3 channels RGB)

# Why Hard ?

- Hard to get criteria



This image is CC0 1.0 public domain



This image is CC0 1.0 public domain



This image is CC0 1.0 public domain



This image is CC0 1.0 public domain

- Harder if occlusion



This image is CC0 1.0 public domain



This image is CC0 1.0 public domain

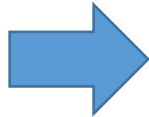


This image by jonsson is licensed under CC-BY 2.0

# An Image Classifier

- Method 1: Define explicit feature (Unfortunately not working)
- Method 2: Machine Learning / Data-driven approach
  - Collect data set and train a model

```
def classify_image(image):  
    # Some magic here?  
    return class_label
```



```
def train(images, labels):  
    # Machine learning!  
    return model
```



Memorize all  
data and labels

```
def predict(model, test_images):  
    # Use model to predict labels  
    return test_labels
```



Predict the label  
of the most similar  
training image



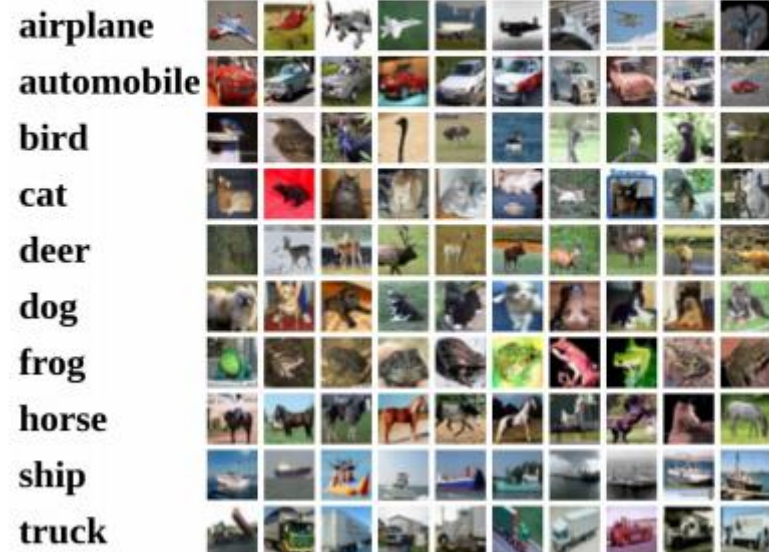
# Common Data set

- MNIST
  - Handwrite digits
  - 10 classes



- CIFAR 10 (100)
  - Small but real pictures

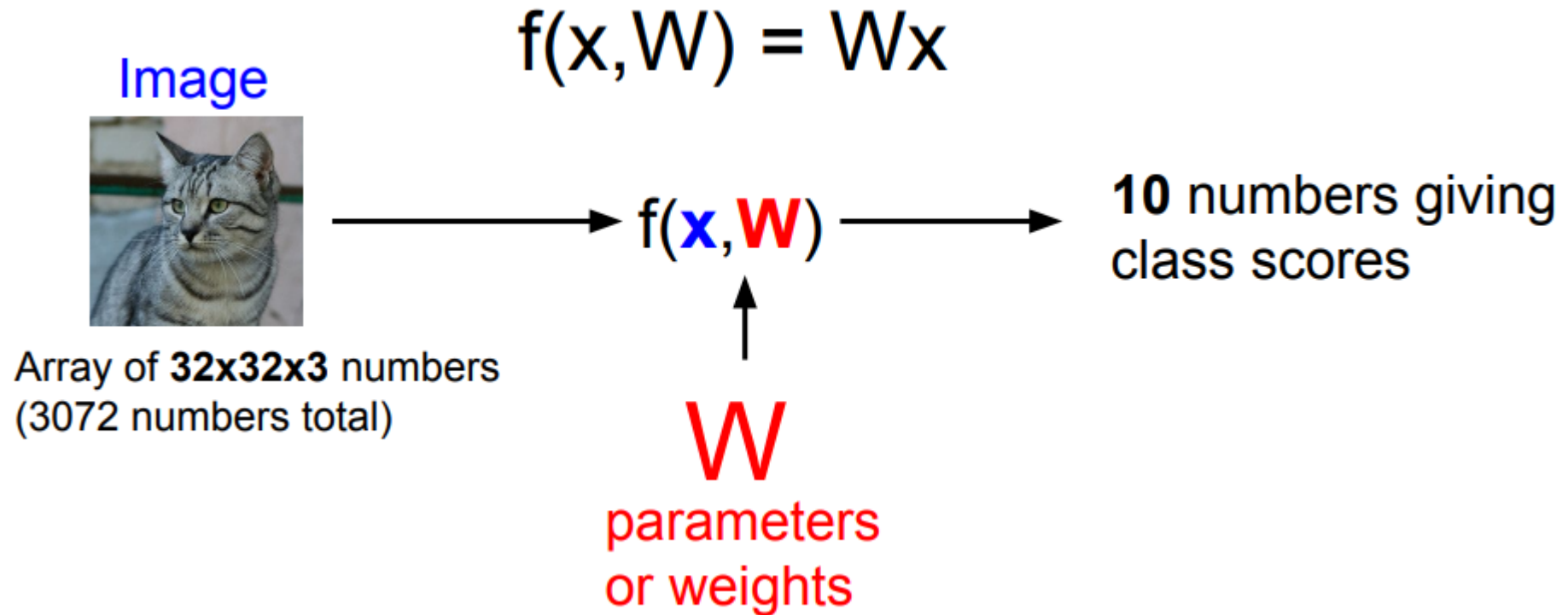
**50,000** training images  
**10,000** testing images





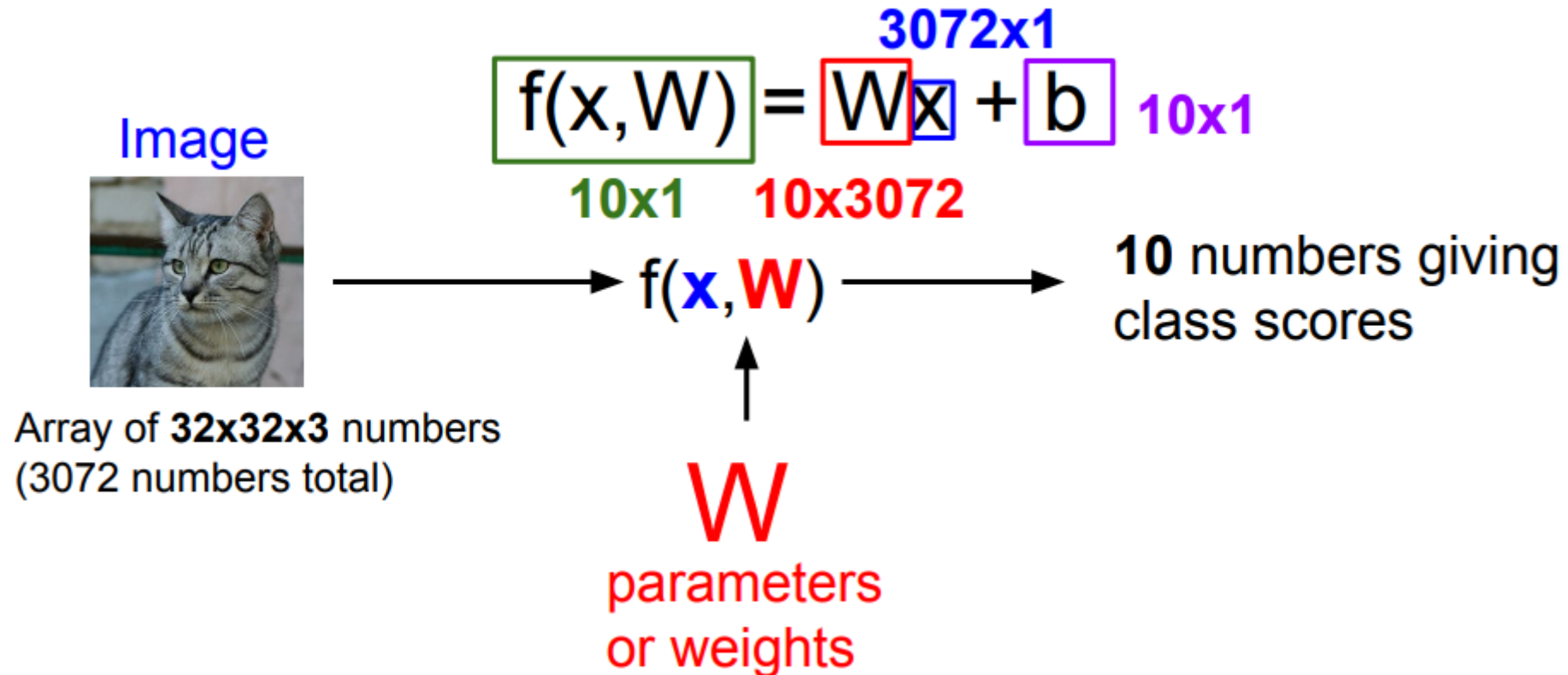
# Linear Classifier

- Use CIFAR 10 as examples
- What's the size of  $x$ ,  $w$  and  $f$  ?



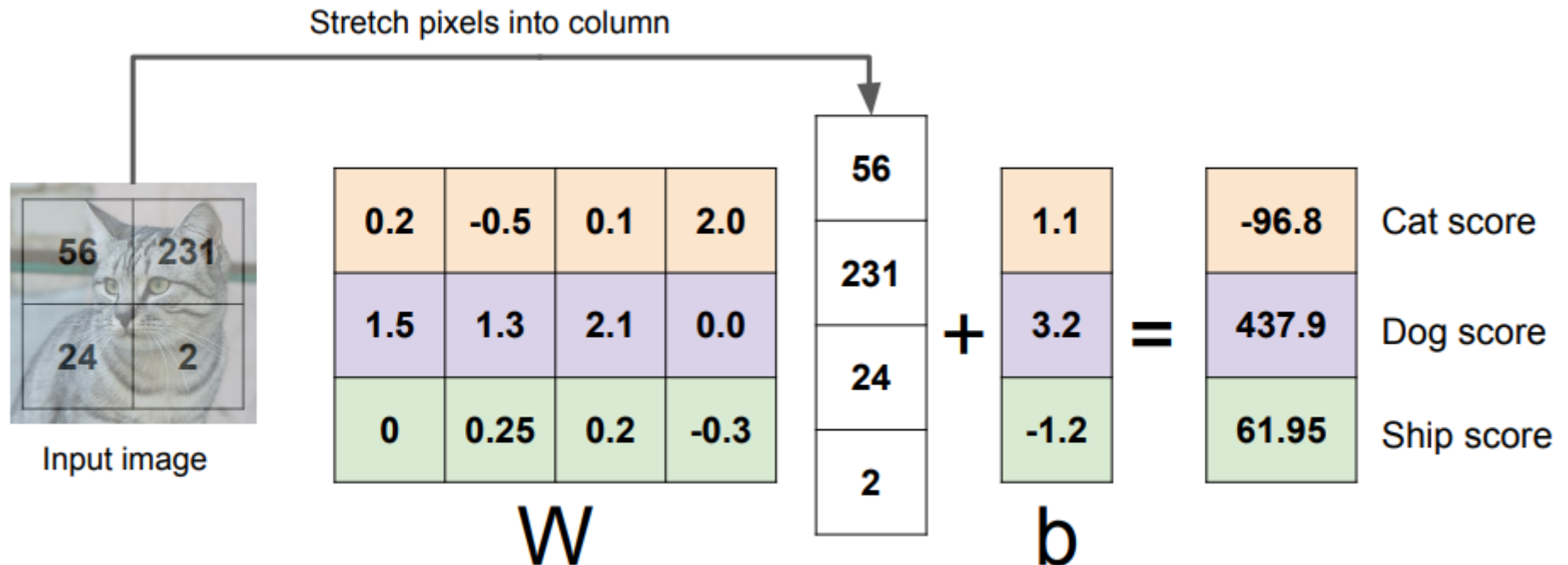
# Score Function

- In practice, we may need bias
- $f()$  is also called score function



# A simple example

- Assuming we have an image with 4 pixels, and 3 classes (cat/dog/ship)
- Pick up the maximum.



# Loss function

- A loss function tells how good a classifier is

Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where  $x_i$  is image and  
 $y_i$  is (integer) label

Loss over the dataset is a  
sum of loss over examples:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

Suppose: 3 training examples, 3 classes.  
With some  $W$  the scores  $f(x, W) = Wx$  are:



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>

# Softmax Classifier

- Other loss function, such as multiclass SVM, is not included in the class.
- Motivation: interpret raw classifier scores as probabilities
- Also known as logistic regression classifier
- Softmax uses a cross-entropy loss

$$L_i = -\log \left( \frac{e^{f_{v_i}}}{\sum_j e^{f_j}} \right)$$

# Softmax Classifier Example - I



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{Softmax Function}$$

Probabilities  
must be  $\geq 0$

cat

**3.2**

car

5.1

frog

-1.7

exp



**24.5**

164.0

0.18

unnormalized  
probabilities

# Softmax Classifier Example - II



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k | X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{Softmax Function}$$

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

cat  
car  
frog

**3.2**  
**5.1**  
**-1.7**

exp

**24.5**  
**164.0**  
**0.18**

normalize

**0.13**  
**0.87**  
**0.00**

unnormalized  
probabilities

probabilities



# Softmax Classifier Example - III



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{Softmax Function}$$

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

$$L_i = -\log P(Y = y_i|X = x_i)$$

cat  
car  
frog

3.2

5.1

-1.7

exp

24.5

164.0

0.18

normalize

0.13

0.87

0.00

$$\rightarrow L_i = -\log(0.13) = 2.04$$

Unnormalized  
log-probabilities / logits

unnormalized  
probabilities

probabilities

**Maximum Likelihood Estimation**  
Choose probabilities to maximize  
the likelihood of the observed data

# Softmax Classifier Example - III



Want to interpret raw classifier scores as **probabilities**

$$s = f(x_i; W)$$

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}} \quad \text{Softmax Function}$$

Probabilities  
must be  $\geq 0$

Probabilities  
must sum to 1

$$L_i = -\log P(Y = y_i|X = x_i)$$

cat  
car  
frog

**3.2**  
**5.1**  
**-1.7**

Unnormalized  
log-probabilities / logits

exp

**24.5**  
**164.0**  
**0.18**

unnormalized  
probabilities

normalize

**0.13**  
**0.87**  
**0.00**

probabilities

compare

**1.00**  
**0.00**  
**0.00**

Correct  
probs

Cross Entropy

$$H(\underline{P}, \underline{Q}) = H(p) + D_{KL}(P||Q)$$

# A computing trick

- Exponential computing might generate very big numbers, exceeding the integer range

- Scaling is before exponential domain

$$\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} = \frac{C e^{f_{y_i}}}{C \sum_j e^{f_j}} = \frac{e^{f_{y_i} + \log C}}{\sum_j e^{f_j + \log C}}$$

- Sample code:  
also work in  
HW!

```
f = np.array([123, 456, 789]) # example with 3 classes and each having large scores
p = np.exp(f) / np.sum(np.exp(f)) # Bad: Numeric problem, potential blowup

# instead: first shift the values of f so that the highest number is 0:
f -= np.max(f) # f becomes [-666, -333, 0]
p = np.exp(f) / np.sum(np.exp(f)) # safe to do, gives the correct answer
```

# How to find the best weight ?

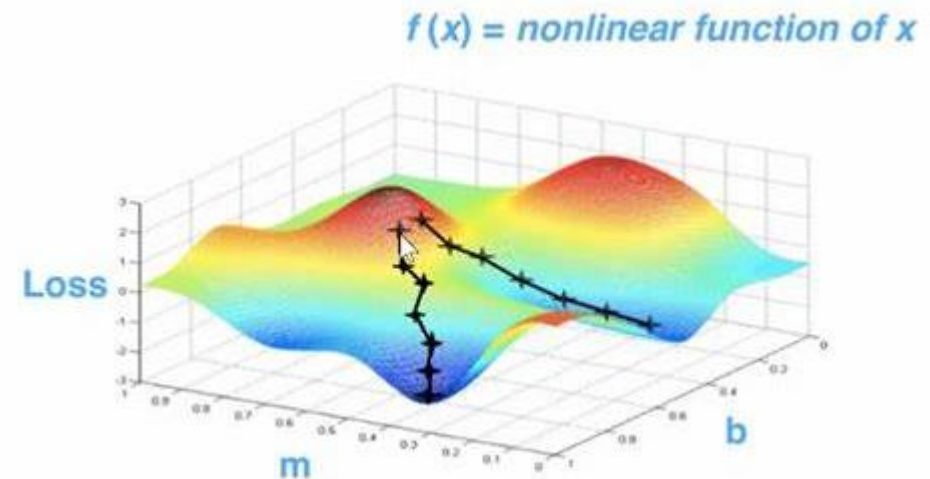
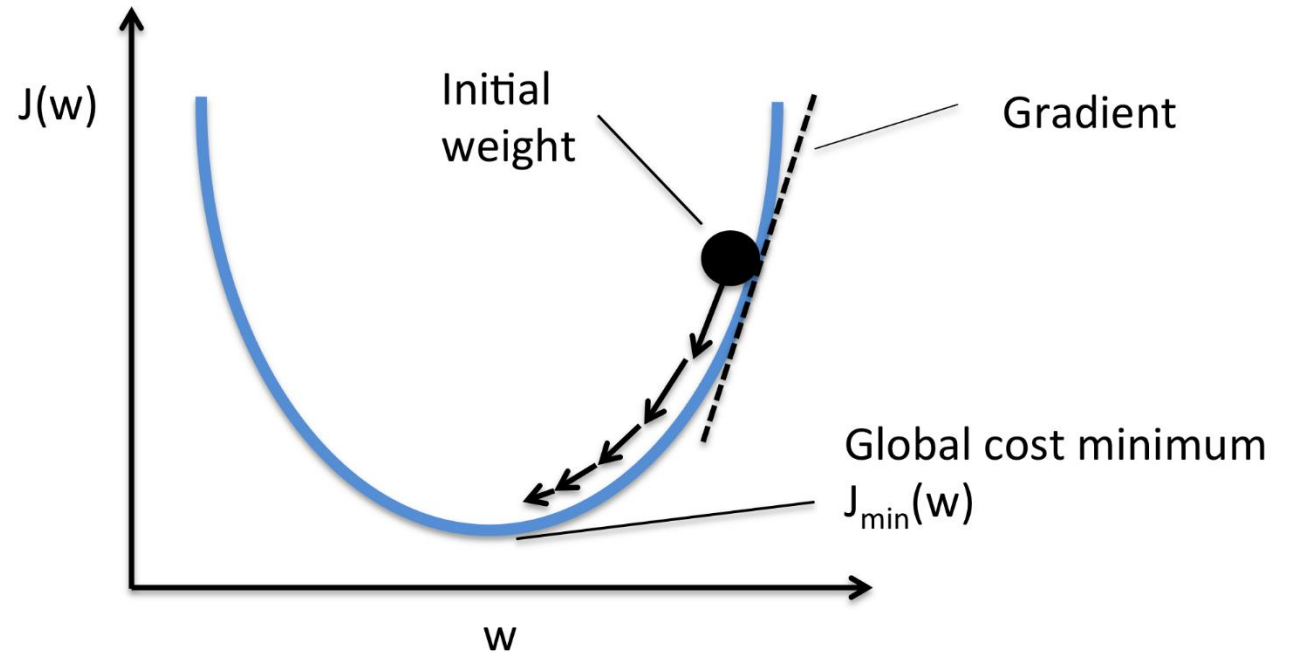
- Recap: dataset, score function, loss function
- How to find the best score function / weights?
- Normally, random initiate provide 10%~15% accuracy for 10 class, like guess
- Find the min. (            ) ?

# Gradient Descent

- Derivative / Gradient

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

- The slope in any direction is the dot product of the direction with the gradient. The direction of steepest descent is the negative gradient.



# Numerical Gradient Computing

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (first dim):**

[0.34 + **0.0001**,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25322**

**gradient dW:**

**[-2.5,**  
?,  
?,

$$(1.25322 - 1.25347)/0.0001 = -2.5$$

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

?,  
?,...]

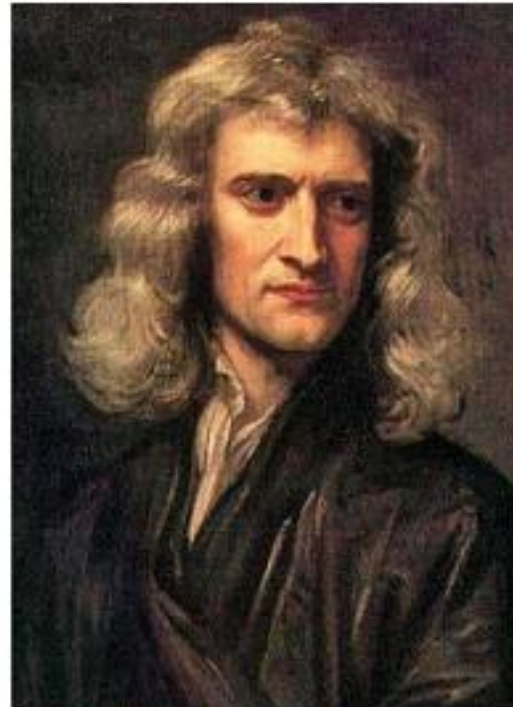


# Analytical Gradient Computing

- We all learned calculus.
- Linear derivative is super easy:

$$f(a)=ax+b, f'(a)=x$$

- What about softmax?



[This image](#) is in the public domain



[This image](#) is in the public domain



# Derivative of softmax classifier

$$P_k = \frac{e^{f_k}}{\sum_j e^{f_j}} \quad L_i = -\sum_k p_{i,k} \log P_k \quad f_m = (x_i W)_m$$

when  $k = m$ ,

$$\frac{\partial P_k}{\partial f_m} = \frac{e^{f_k} \sum_j e^{f_j} - e^{f_k} \cdot e^{f_k}}{(\sum_j e^{f_j})^2} = P_k(1 - P_k)$$

when  $k \neq m$ ,

$$\frac{\partial P_k}{\partial f_m} = -\frac{e^{f_k} e^{f_m}}{(\sum_j e^{f_j})^2} = -P_k P_m$$

then:

$$\frac{\partial L_i}{\partial f_m} = -\sum_k p_{i,k} \frac{\partial \log P_k}{\partial f_m}$$

$$\begin{aligned} &= -\sum_k p_{i,k} \frac{1}{P_k} \frac{\partial P_k}{\partial f_m} \\ &= -\sum_{k=m} p_{i,k} \frac{1}{P_k} P_k (1 - P_k) + \sum_{k \neq m} p_{i,k} \frac{1}{P_k} P_k P_m \\ &= \sum_{k \neq m} p_{i,k} P_m - \sum_{k=m} p_{i,k} (1 - P_k) \\ &= \begin{cases} P_m & , \quad m \neq y_i \\ P_m - 1 & , \quad m = y_i \end{cases} \\ &= P_m - p_{i,m} \end{aligned}$$

Last:

$$\frac{\partial L_i}{\partial W_k} = \frac{\partial L_i}{\partial f_m} \frac{\partial f_m}{\partial W_k} = x_i^T (P_m - p_{i,m})$$

$$\nabla_{W_k} L = -\frac{1}{N} \sum_i x_i^T (p_{i,m} - P_m) + 2\lambda W_k$$

# Numerical vs. Gradient

- Numerical gradient: approximate, slow, easy to write
- Analytic gradient: exact, fast, error-prone

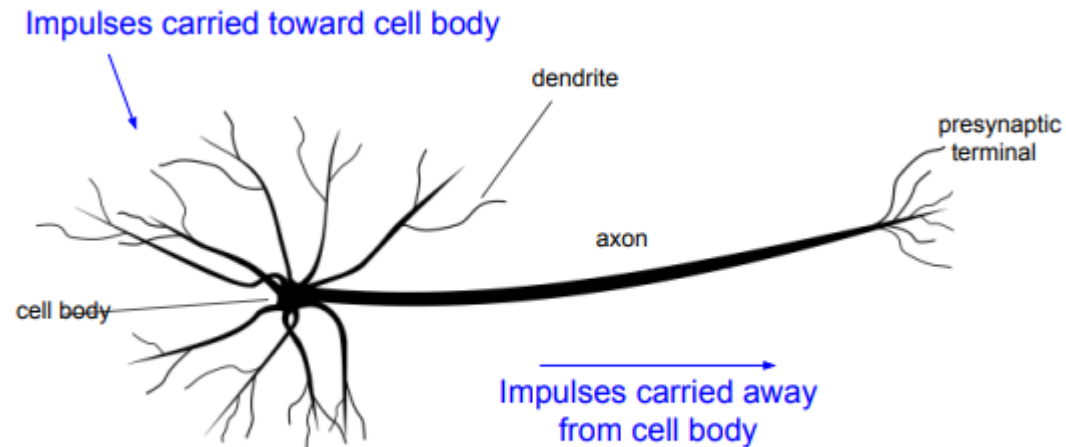
---- When you only have CPUs

---- Numerical gradient is NOT slow when you have accelerator!

Neuron Network

# Neuron

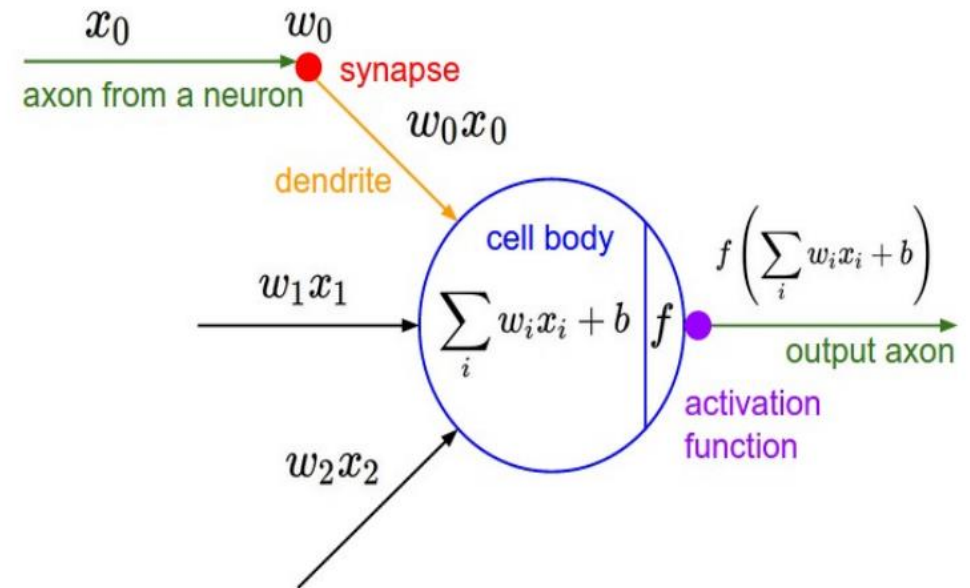
- Neuron Structure



This image by Felipe Perucho is licensed under [CC-BY 3.0](#)

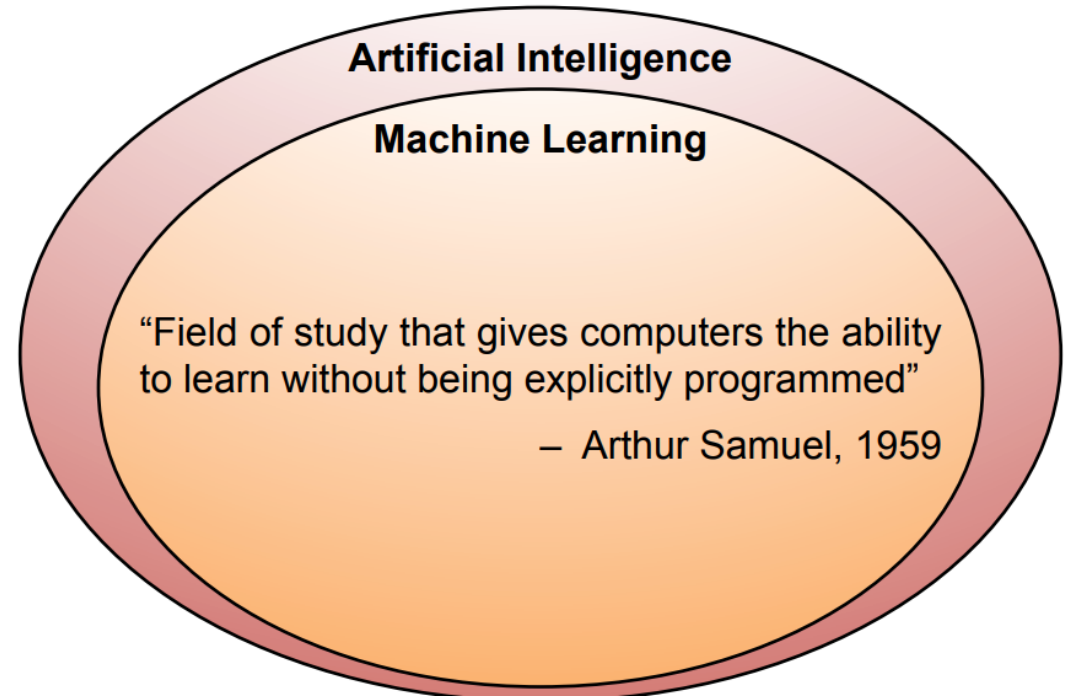
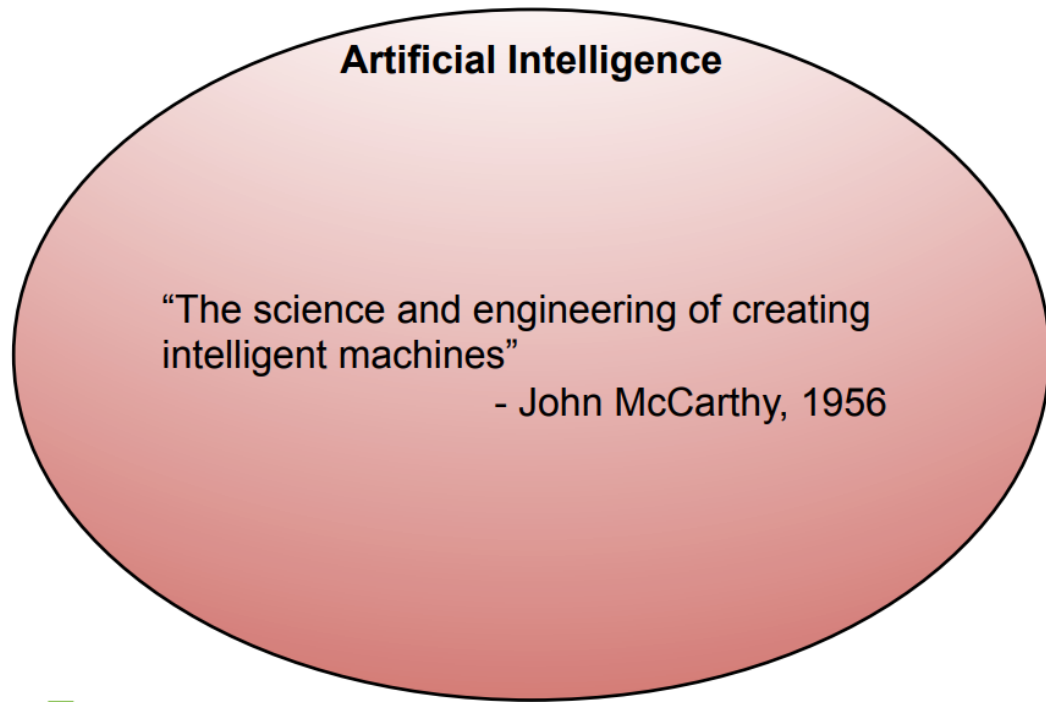
- It has many layers.

- Its math model

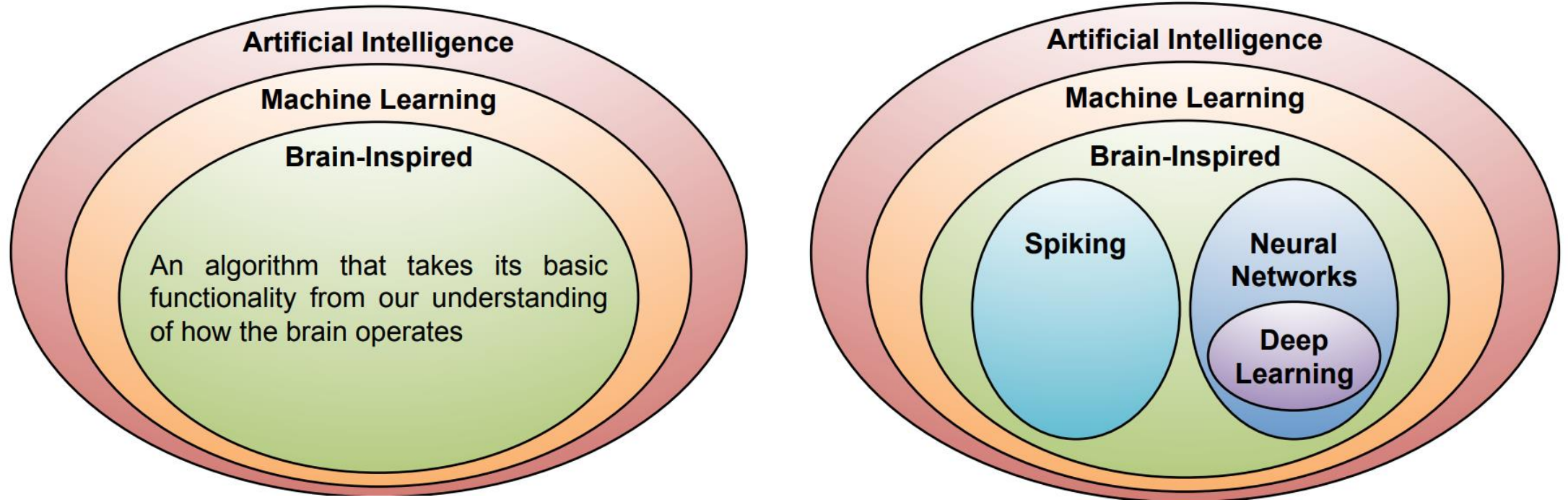


Very similar to linear classifier, **except** ...

# Artificial Intelligence and Machine Learning

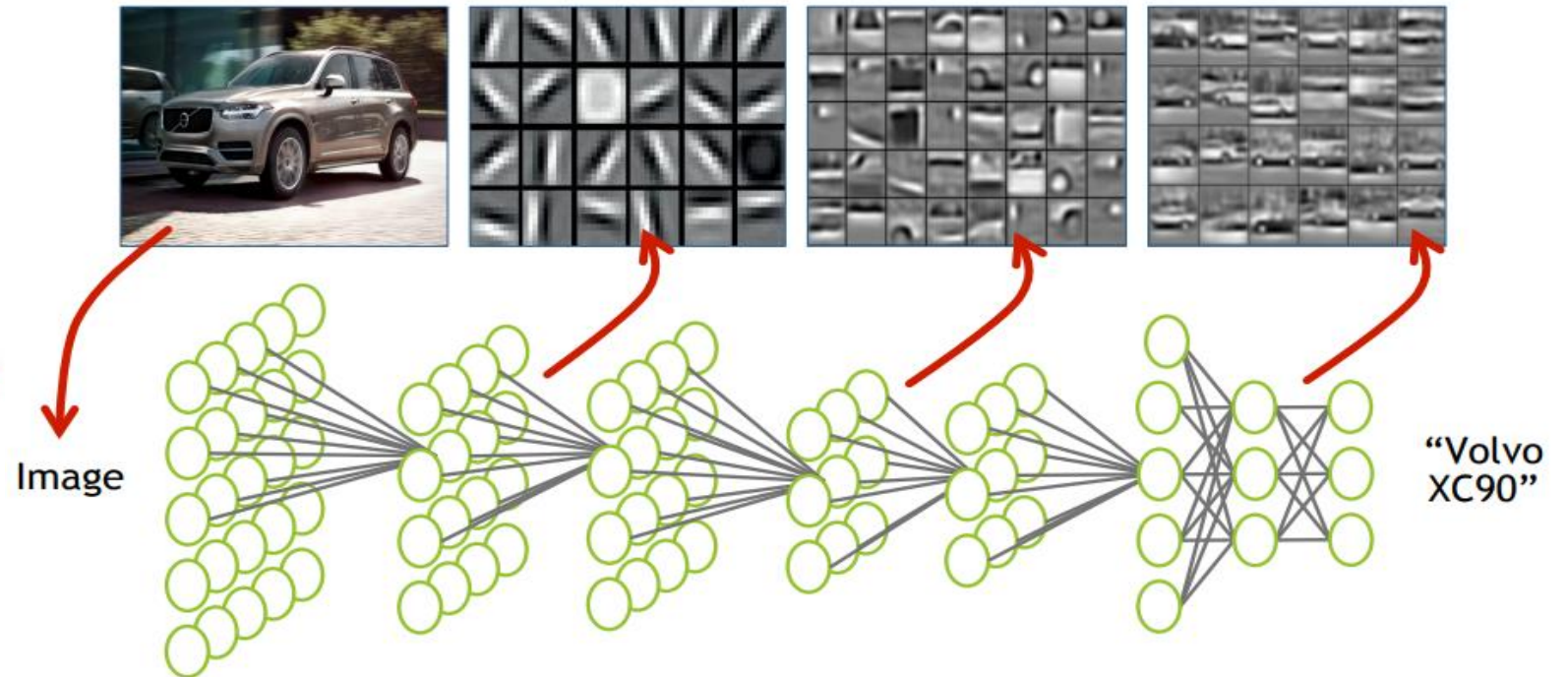
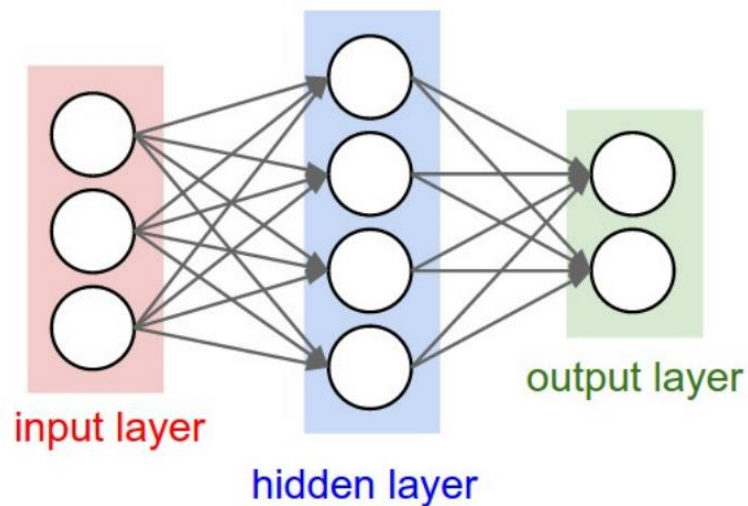


# Brain Inspired and Neural Network



# MLP and Deep Learning

- Previously, we have a name called: MLP: multi-level perceptron
- Now: Deep Learning





# Back-propagation

- How can we find the derivative for multi-layer structure?

Backpropagation: a simple example

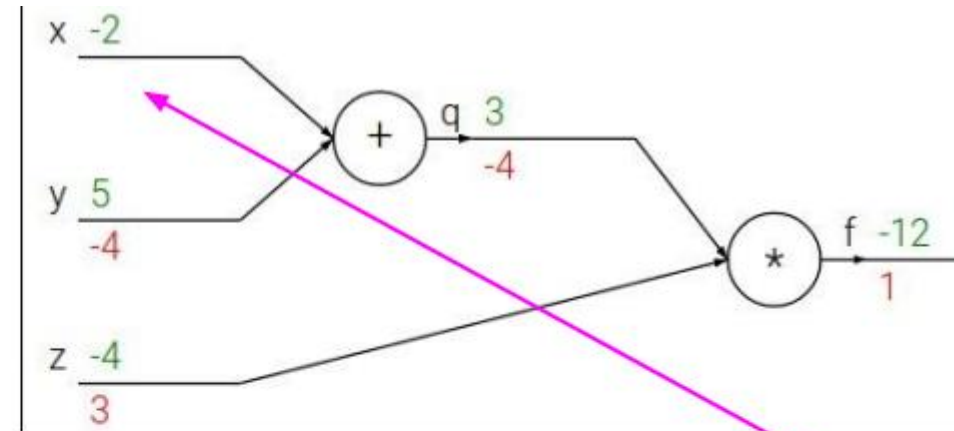
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Chain rule:

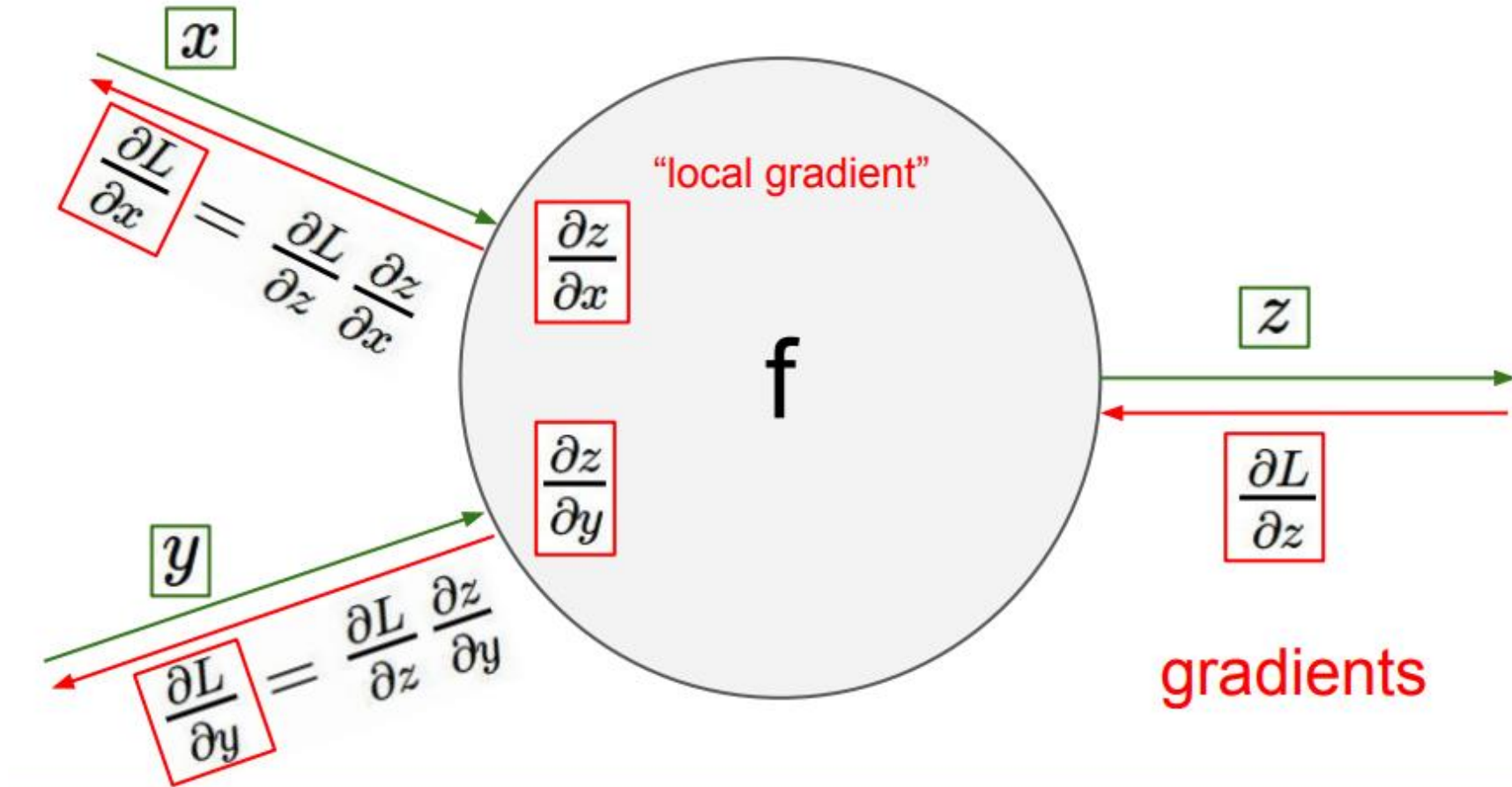
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Upstream  
gradient

Local  
gradient

$$\frac{\partial f}{\partial x}$$

# Chain Rule



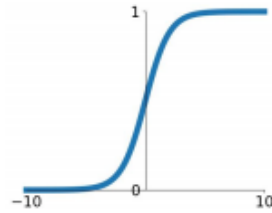
# Activation Functions

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1+e^{-x})^2} = \left( \frac{1+e^{-x}-1}{1+e^{-x}} \right) \left( \frac{1}{1+e^{-x}} \right) = (1-\sigma(x))\sigma(x)$$

**Sigmoid**

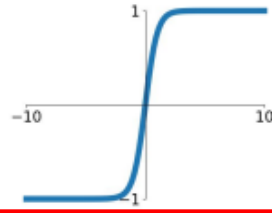
Sigmoid is easy  
for derivative

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



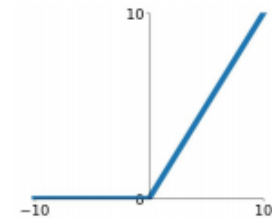
**tanh**

$$\tanh(x)$$



**ReLU**

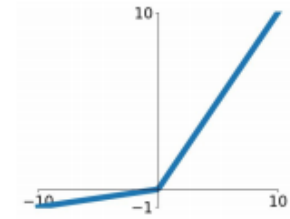
$$\max(0, x)$$



ReLU is the most common one  
Used in recent CNN

**Leaky ReLU**

$$\max(0.1x, x)$$

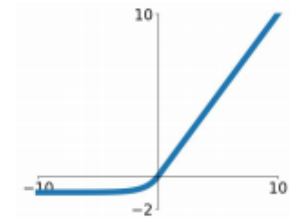


**Maxout**

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

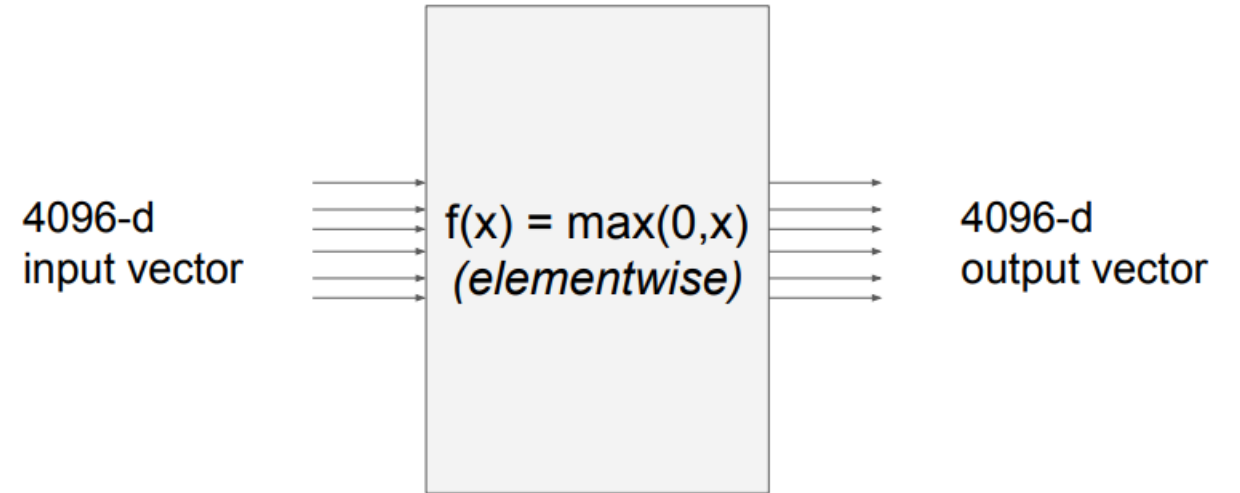
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



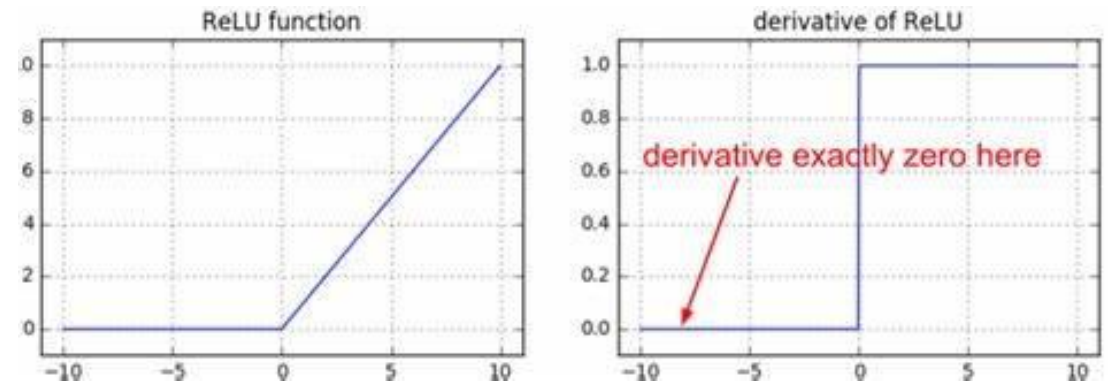
# Vectorized Gradient Computing

- Jacobian Matrix

$$\mathbf{J} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \cdots & \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$



- How big will be the Jacobian matrix?
- Mini-batch make it worse



# Using Numpy to train a 2-layer model

- For simplicity: use sigmoid as activation

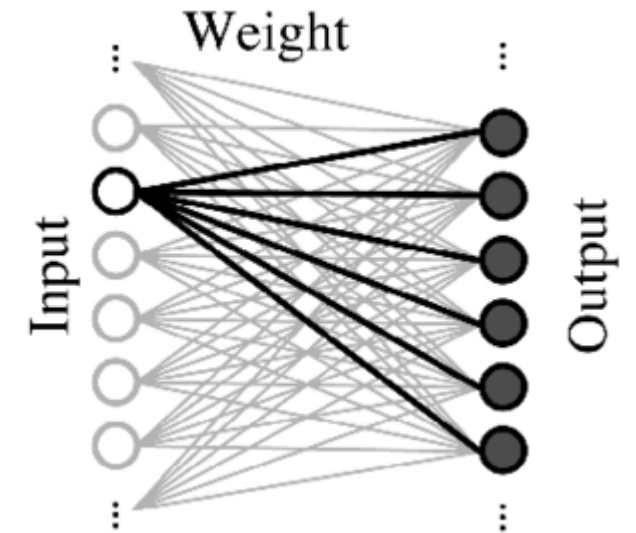
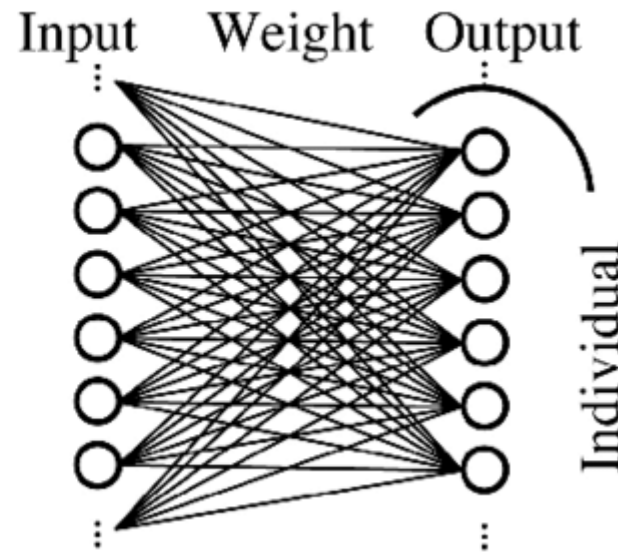
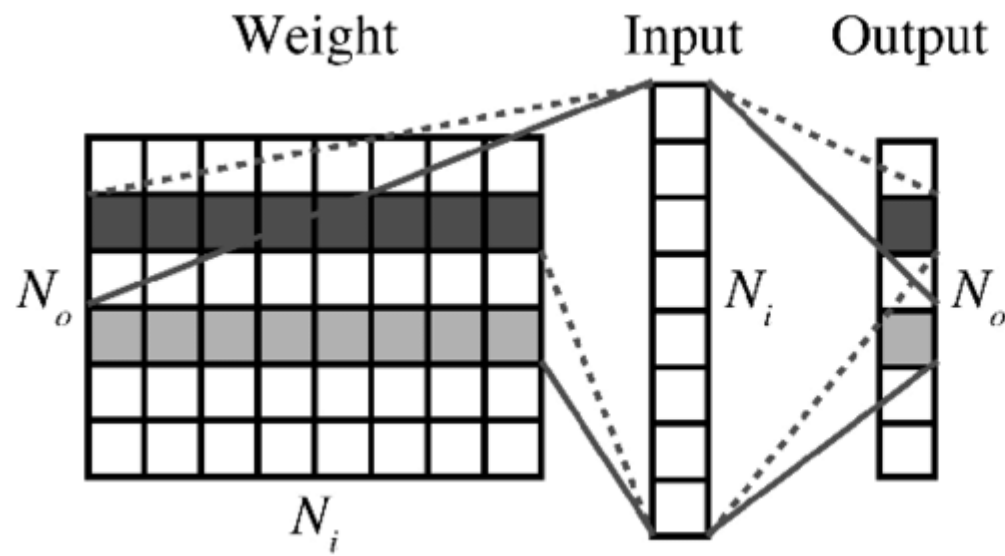
```
1  import numpy as np
2  from numpy.random import randn
3
4  N, D_in, H, D_out = 64, 1000, 100, 10
5  x, y = randn(N, D_in), randn(N, D_out)
6  w1, w2 = randn(D_in, H), randn(H, D_out)
7
8  for t in range(2000):
9      h = 1 / (1 + np.exp(-x.dot(w1)))
10     y_pred = h.dot(w2)
11     loss = np.square(y_pred - y).sum()
12     print(t, loss)
13
14     grad_y_pred = 2.0 * (y_pred - y)
15     grad_w2 = h.T.dot(grad_y_pred)
16     grad_h = grad_y_pred.dot(w2.T)
17     grad_w1 = x.T.dot(grad_h * h * (1 - h))
18
19     w1 -= 1e-4 * grad_w1
20     w2 -= 1e-4 * grad_w2
```

# Inference vs. Training

- Inference: Apply weights to determine output
- Training: Determine weights
  - Supervised: Training set has inputs and outputs, i.e., labeled
  - Unsupervised: Training set is unlabeled
  - Semi-supervised: Training set is partially labeled
  - Reinforcement: Output assessed via rewards and punishments

# Fully Connected Layer

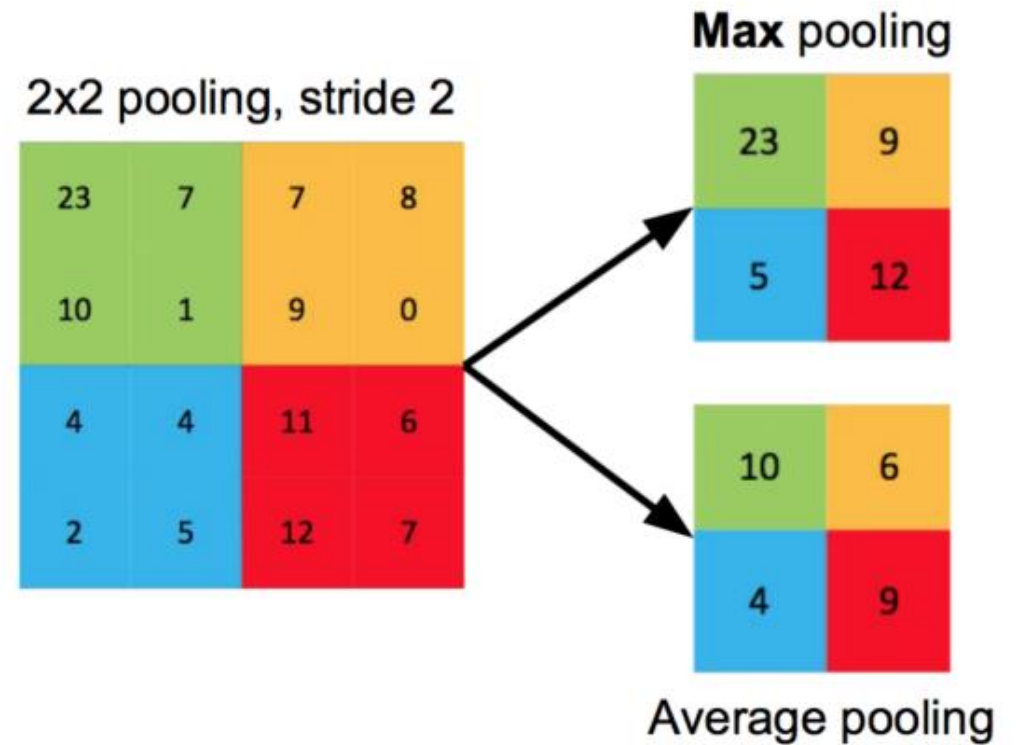
- Just as the linear classifier
- Even today, the last stage of DNNs are still FC layers





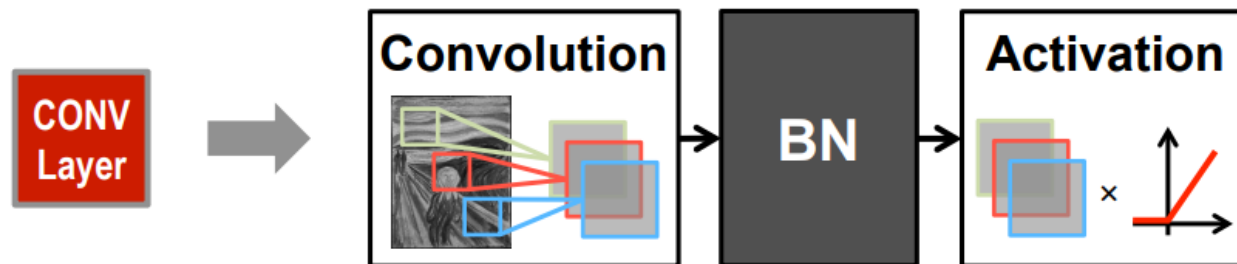
# Pooling Layer

- Reduce resolution of each channel independently
- Overlapping or non-overlapping  
→ depending on stride
- Increases translation-invariance and noise-resilience



# Normalization Layer

- Batch Normalization (BN)
  - Normalize activations towards mean=0 and std. dev.=1 based on the statistics of the training dataset
  - Believed to be key to getting high accuracy and faster training on very deep neural networks.
  - Makes training harder, but almost do not affect inference.



$$y = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} \gamma + \beta$$

Annotations for the equation:

- $\mu$ : data mean (red arrow)
- $\sigma^2$ : data std. dev. (red arrow)
- $\epsilon$ : small const. to avoid numerical problems (gray arrow)
- $\gamma$ : learned scale factor (blue arrow)
- $\beta$ : learned shift factor (blue arrow)

# Summary

Next time: CNN