

Introduction to Deep Learning

Lecture 17

15.076 | Analytics for a Better World



"Deep Learning, oil on canvas, masterpiece by Salvador Dalí", a creation by Léonard Boussioux with 30+ prompts in DALL-E

1

Plan

Broad Introduction of AI

Part I: Fundamental Concepts of Neural Networks

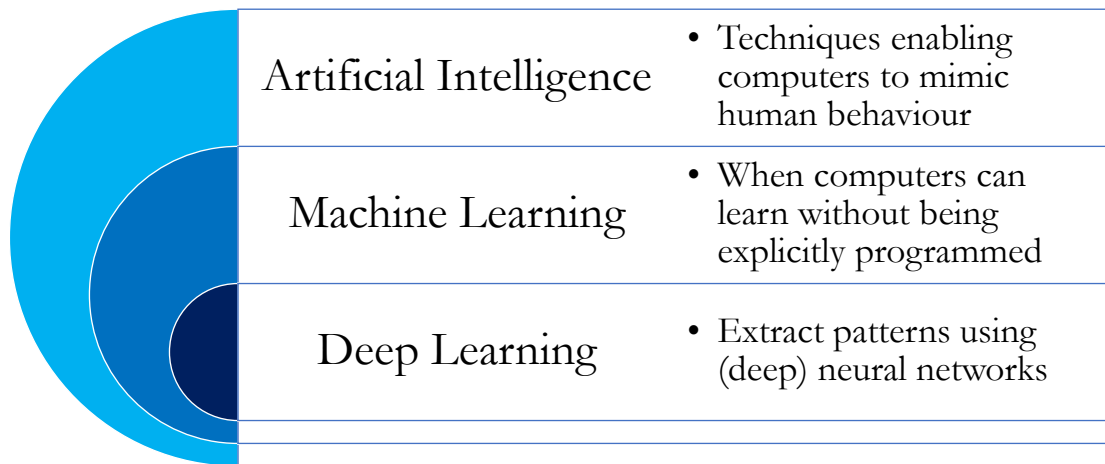
- The perceptron model
- Building neural networks
- Training neural networks
- Handle Overfitting

Part II: Convolutional Neural Networks

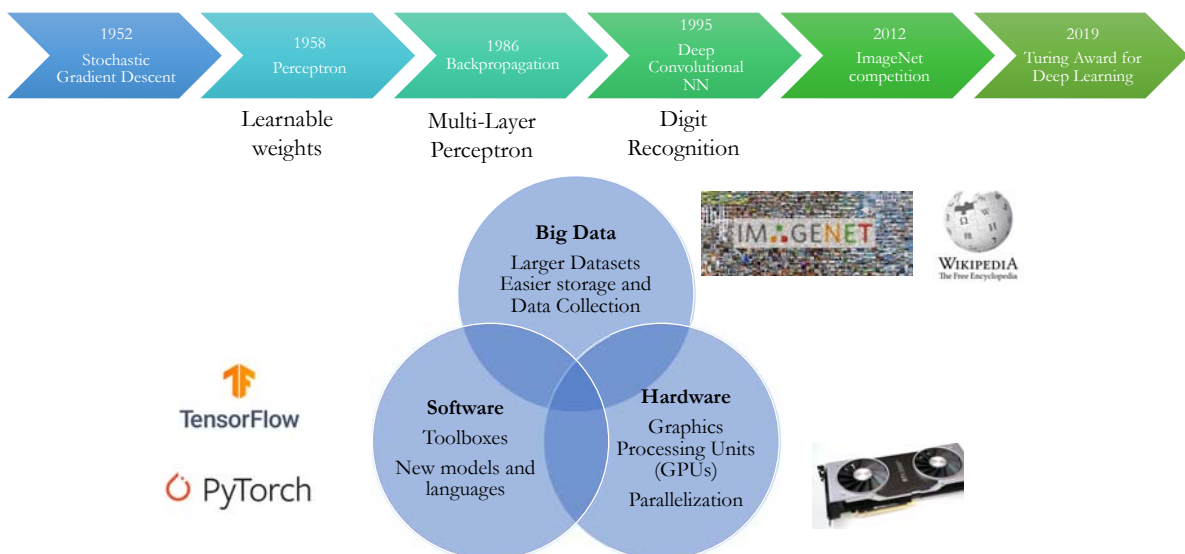


2

Deep Learning in the landscape



What is the history behind the current AI boom?



Deep Learning



What society thinks I do



What my friends think I do



What other computer scientists think I do



What mathematicians think I do



What I think I do

```
from keras import *
```

What I actually do

5

Deep Learning



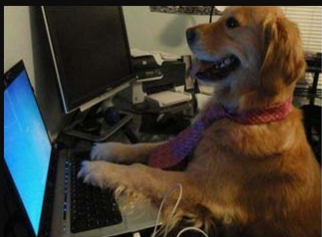
What society thinks I do



What my friends think I do



What other computer scientists think I do



What mathematicians think I do



What I think I do

```
from keras import *
```

So 2022...

What I actually do

6

Deep Learning



What society thinks I do



What my friends think I do



What other computer scientists think I do



What mathematicians think I do



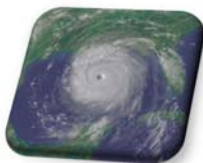
What I think I do

Hi ChatGPT, I am ~~lazy~~ efficient,
code this for me

What I actually do

7

During my PhD, I have used Deep Learning for:



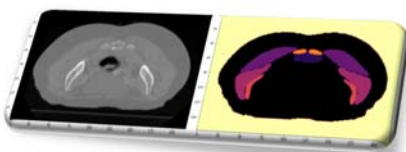
Tropical Cyclone
Forecasting



Healthcare
Operations



Wildlife
Identification



Tumor segmentation

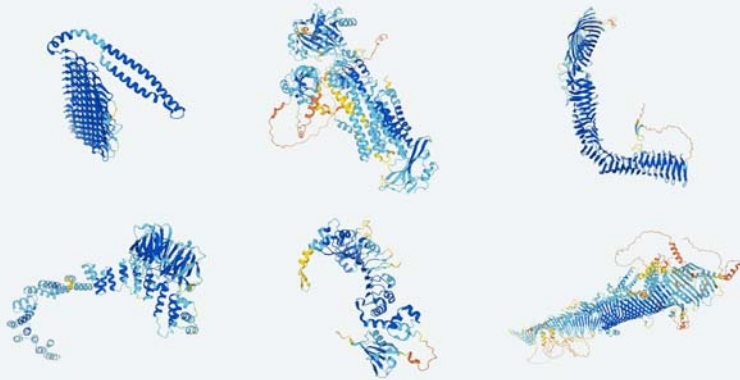


Stamp analytics

8

AlphaFold

An AI system developed by DeepMind that predicts a protein's 3D structure from its amino acid sequence.



9

DALL-E in 2022

<https://openai.com/blog/dall-e/>

BabyShark in an existential crisis after its overwhelming success on YouTube



DALL-E Prompt: an illustration of a baby shark in pajamas staring at its reflection in the mirror

Pikachu very proud for its first day at MIT

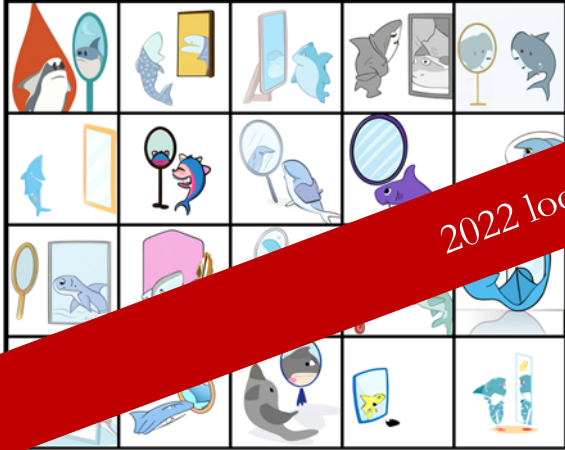


DALL-E Prompt: an illustration of a pikachu with a cape holding a calculator

DALL-E in 2022

<https://openai.com/blog/dall-e/>

BabyShark in an existential crisis after its overwhelming success on YouTube



DALL-E Prompt: an illustration of a baby shark in pajamas staring at its reflection in the mirror

Pikachu very proud for its first



DALL-E Prompt: an illustration of a pikachu with a cape holding a calculator

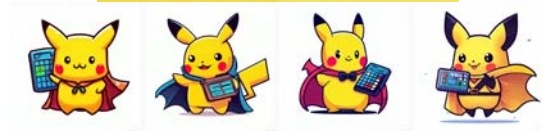
2022 looks so retro...

Same prompts in 2023...



DALL-E 2

Midjourney



DALL-E 2

12

I launched my
digital artist
career!
Exhibition
coming soon!



With great power comes great responsibility



AI that recognizes faces from the whole world



Real-time satellite analysis to get trading side information

Michelle Obama



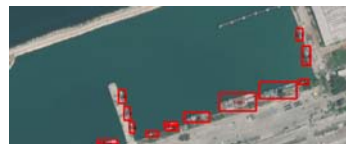
Bias, racism



Killer drones, AI-powered weapons

15

With great power comes great responsibility



BUSINESS • TECHNOLOGY

Exclusive: OpenAI Used Kenyan Workers on Less Than \$2 Per Hour to Make ChatGPT Less Toxic



AI

tion



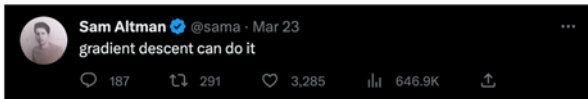
Bias, racism



Killer drones, AI-powered weapons

16

AGI doomers,
AI moratorium,
AI alignment,
Scale is all you need.




The New York Times

A.I. and Chatbots > Become an A.I. Expert > How Chatbots Work > Why Chatbots 'Hafunctor!' > How to Use Chatbots > What's the Future for A.I.?

Elon Musk and Others Call for Pause on A.I., Citing 'Profound Risks to Society'

More than 1,000 tech leaders, researchers and others signed an open letter urging a moratorium on the development of the most powerful artificial intelligence systems.

Give this article > > > 283

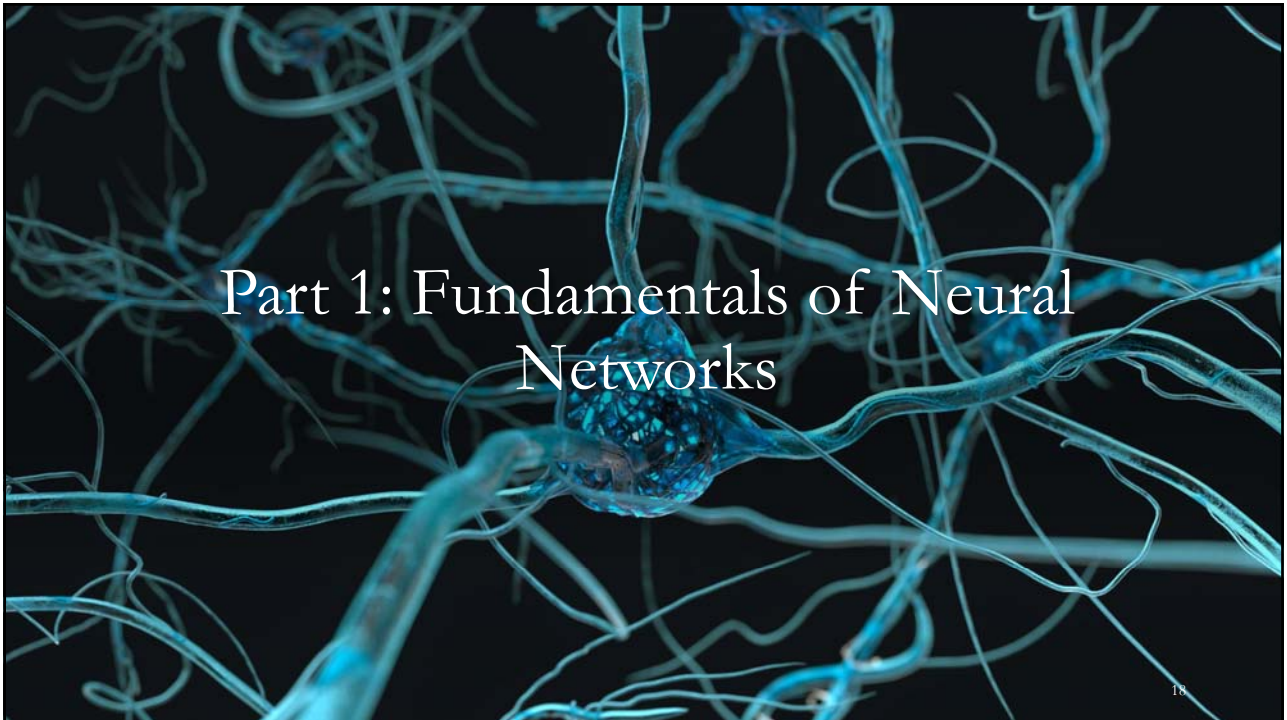


Elon Musk, the chief executive of Twitter and Tesla, and other tech leaders have criticized an "out-of-control race" to develop more advanced artificial intelligence. Benjamin Taniguchi/Associated Press

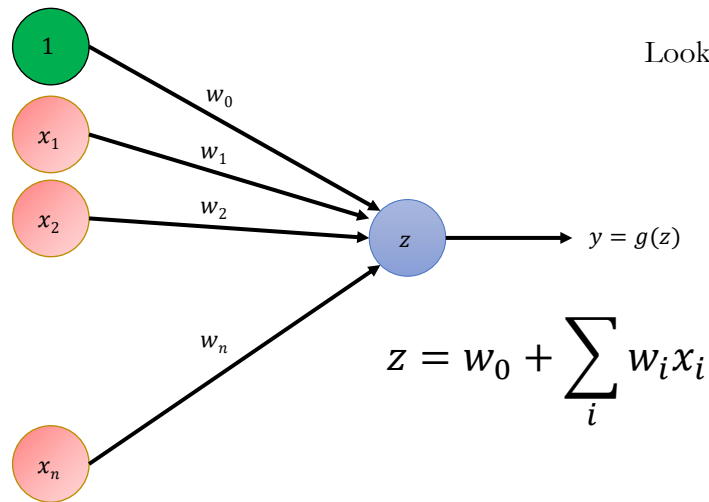
By Cade Metz and Gregory Schmidt
March 28, 2023
閱讀繁體中文版 閱讀簡體中文版

More than 1,000 technology leaders and researchers, including Elon Musk, have urged artificial intelligence labs to pause development of the most advanced systems, warning in [an open](#)

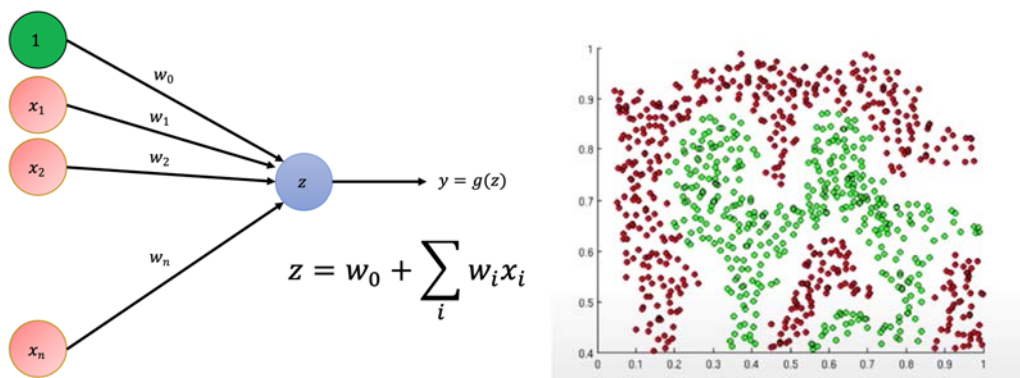
17



A simplified perceptron

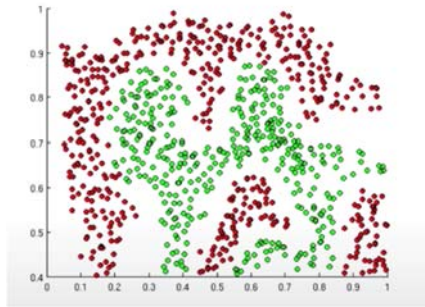


How can we separate the red and green points?

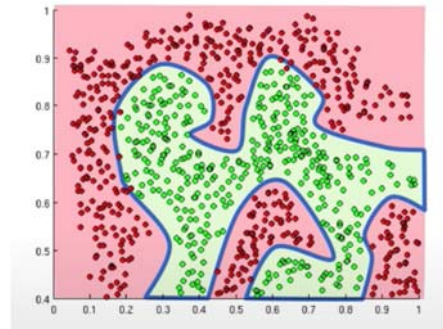


Using activation functions!

$$z = w_0 + \sum_i w_i x_i \longrightarrow y = g(z), \text{ where } g \text{ is a non-linear activation function.}$$



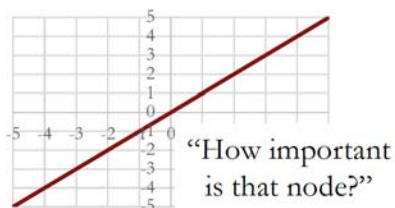
Linear activation functions produce linear decisions no matter the network size.



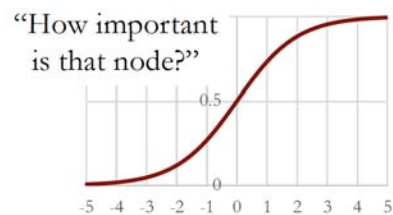
Non-linear activations can help approximate arbitrarily complex functions.

Classic Activation Functions

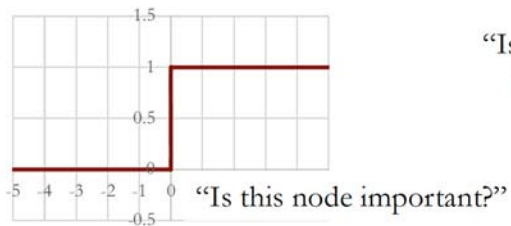
Linear activation



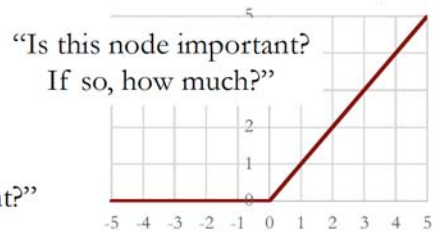
Sigmoid activation



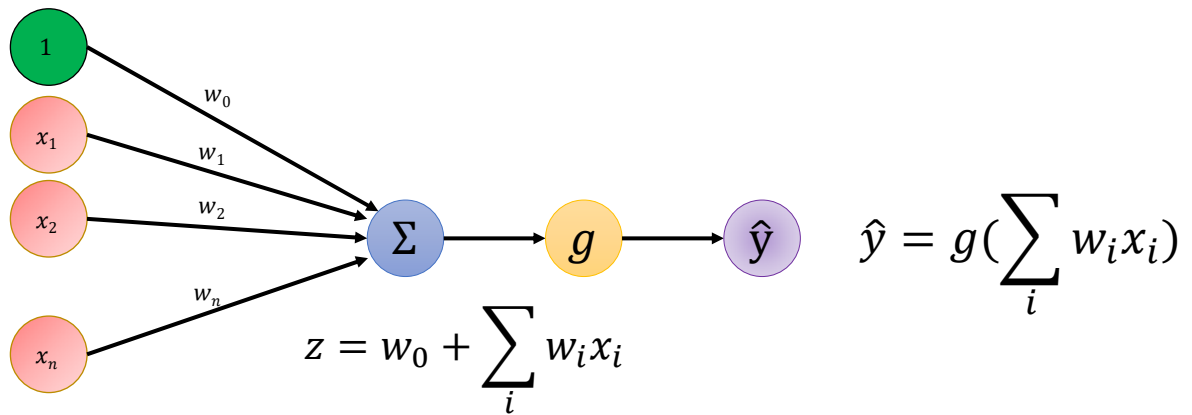
Threshold activation



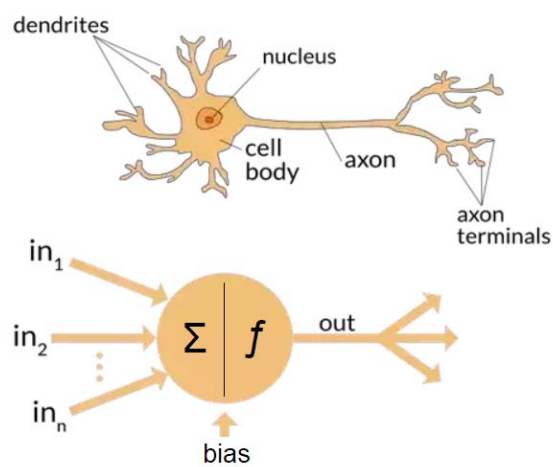
Rectified Linear Unit (ReLU)



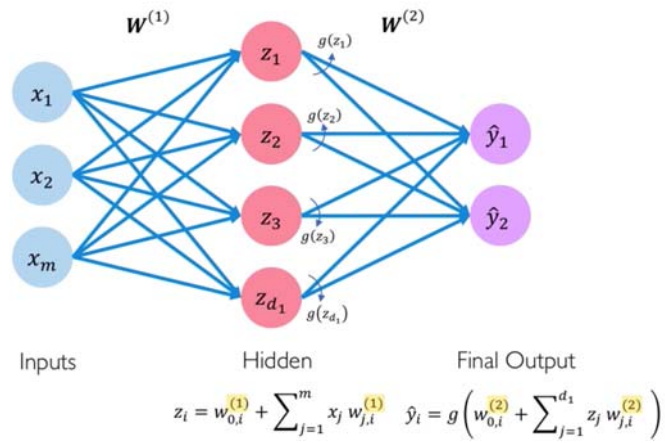
The full perceptron: forward propagation



Artificial Neuron vs Brain Neuron

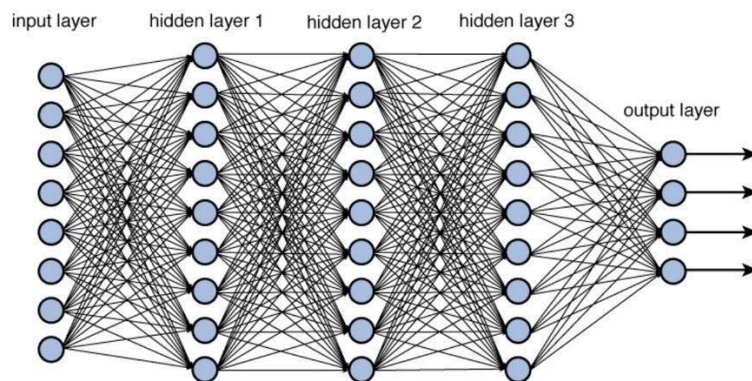


One-layer Neural Network



Picture source: <http://introtodeeplearning.com/>

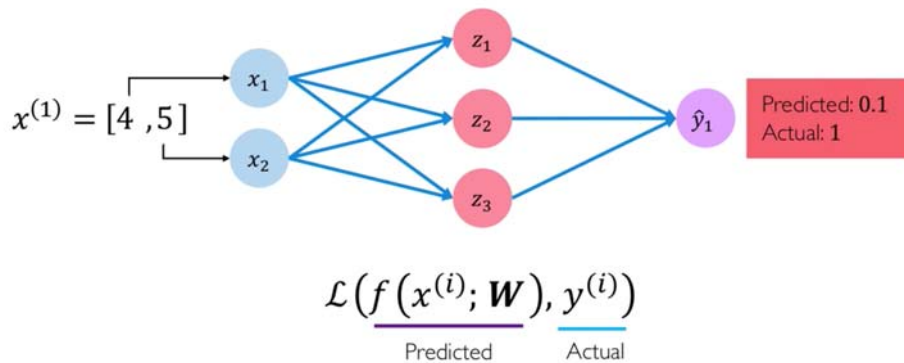
Deep Neural Network



26

Loss function

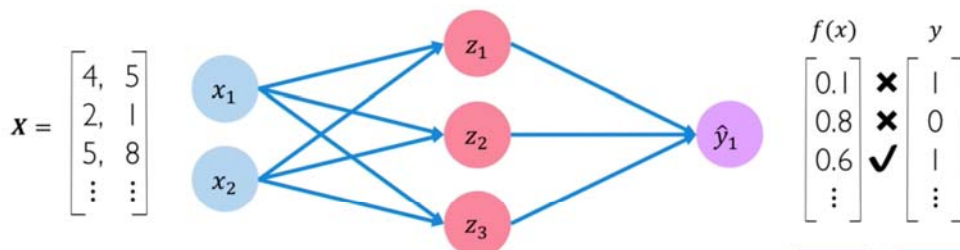
The **loss** of our network measures the cost incurred from incorrect predictions.



Picture source: <http://introtodeeplearning.com/>

Empirical Loss

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

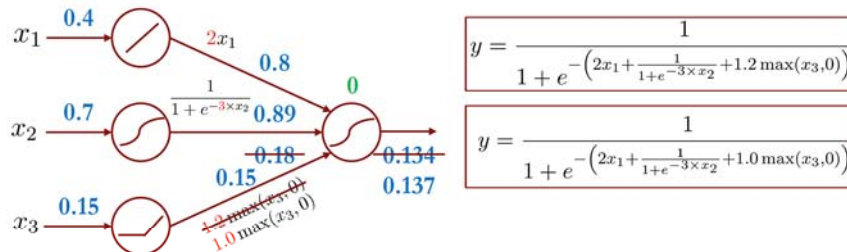


- Also known as:
- Objective function
 - Cost function
 - Empirical Risk

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

Picture source: <http://introtodeeplearning.com/>

Finding the Best Model Fit (explanation on the board)



- The parameters (“weights”) are chosen to maximize the model’s quality of fit—so the predictions \hat{y}_i match the observations y_i

regression

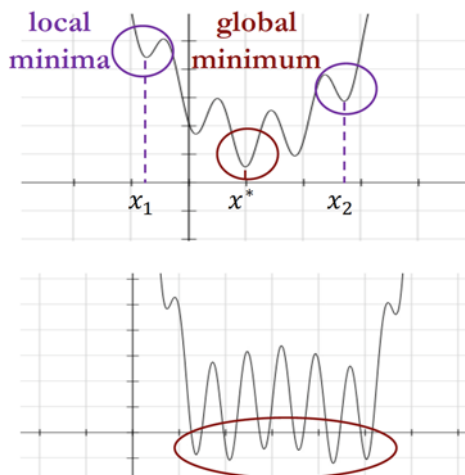
$$\min \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

classification

$$\min - \sum_{i=1}^n y_i \log(\hat{y}_i)$$

29

Nonlinear Optimization



Deep learning loss functions often have many local minima.

→ Deep learning outputs may not be the best possible ones.

However, many local minima often lead to similar solutions.

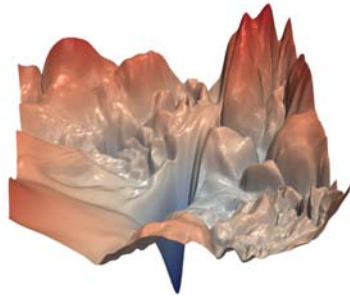
→ Deep learning outputs are often surprisingly “good enough”!

30

Training Neural Networks looks Hard!

Problem 1: Need to optimize a huge number of weights \mathbf{W} .

Problem 2: Loss landscape not smooth at all.



31

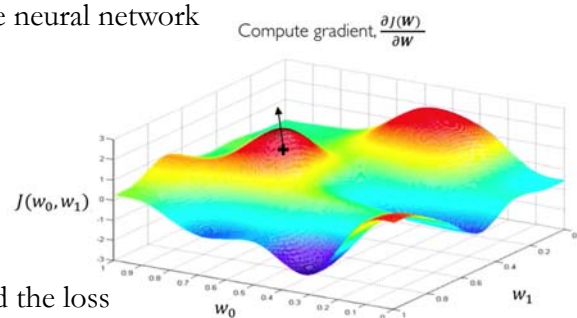
Fortunately there is Gradient Descent!

Objective: Minimizing average loss function of the neural network

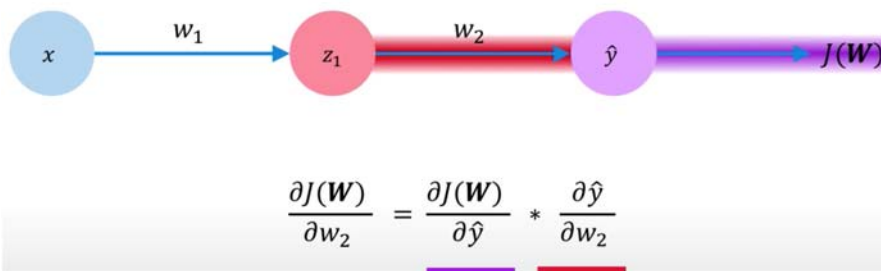
$$\mathcal{J}(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i, \mathbf{W}), y_i)$$

Algorithm:

1. Initialize weight matrix \mathbf{W}
2. **Backpropagation:** Loop until convergence...
 - a. **Forward pass:** Compute the predictions and the loss
 - b. **Backward pass:** Compute the gradient $\frac{\partial \mathcal{J}}{\partial \mathbf{W}}$ (i.e., what happens if we change the weights just a tiny bit)
 - c. **Update:** $w \leftarrow w - \alpha \frac{\partial \mathcal{J}}{\partial w}$ **α : "learning rate"**
3. Return the weight matrix \mathbf{W} and the predictions.



And backpropagation with chain rule



Picture source: <http://introtodeeplearning.com/>

But how do you choose the learning rate?

Remember:

Optimization through gradient descent

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

How can we set
the learning rate?

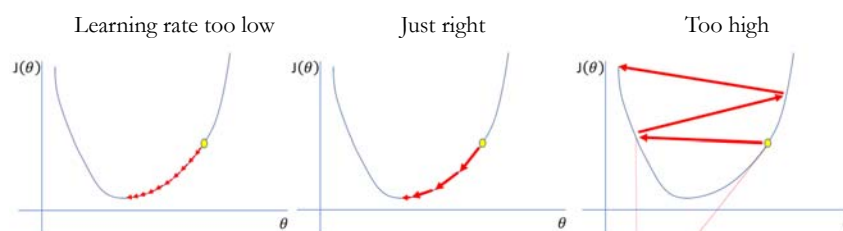
Ideas?

Adaptive Learning Rate

Learning rate is not fixed anymore, instead we can adapt it depending on:

- How large gradients are,
- How fast the learning is happening (momentum),

In Deep Learning, **Adam** is an extremely popular method and does all this automatically. You still need to choose a **base value** (typically between $10^{-3}/5.10^{-5}$ that depends on the task and architecture.)



Another problem: computing gradients is expensive!

Objective: Minimizing average loss function of the neural network

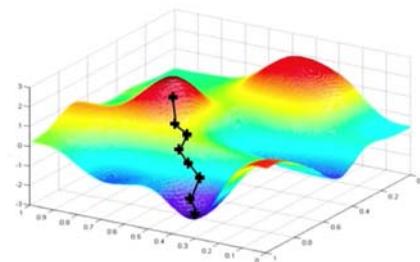
$$\mathcal{J}(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i, \mathbf{W}), y_i)$$

Algorithm:

1. Initialize weight matrix \mathbf{W}
2. **Backpropagation:** Loop until convergence...
 - a. **Forward pass:** Compute the predictions and the loss
 - b. **Backward pass:** Compute the gradient $\frac{\partial \mathcal{J}}{\partial \mathbf{W}}$
 - c. **Update:**

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \frac{\partial \mathcal{J}}{\partial \mathbf{W}}$$

3. Return the weight matrix \mathbf{W} and the predictions.



$$\frac{\partial \mathcal{J}}{\partial \mathbf{W}}$$

Expensive to compute

Solution: Stochastic Gradient Descent

Objective: Minimizing average loss function of the neural network

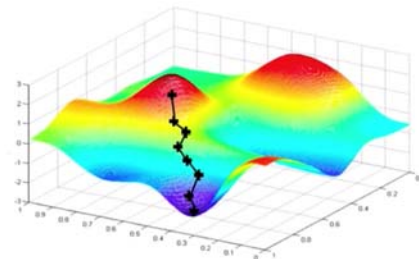
$$\mathcal{J}(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i, \mathbf{W}), y_i)$$

Algorithm:

1. Initialize weight matrix \mathbf{W}
2. **Backpropagation:** Loop until convergence...
 - a. **Forward pass:** Compute the predictions and the loss
 - b. **Backward pass:** Compute the gradient $\frac{\partial \mathcal{J}}{\partial \mathbf{W}}$
 - c. **Update:**

$$\mathbf{W} \leftarrow \mathbf{W} - \alpha \frac{\partial \mathcal{J}}{\partial \mathbf{W}}$$

3. Return the weight matrix \mathbf{W} and the predictions.



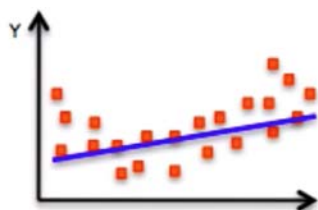
Expensive to compute

Batch Gradient Descent

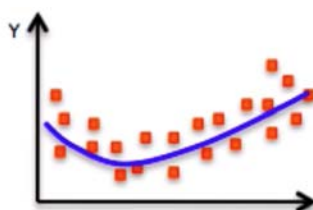
$$\frac{\partial \mathcal{J}(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial \mathcal{J}_k(\mathbf{W})}{\partial \mathbf{W}}$$

We use **mini-batches** while training allowing for a smoother convergence and larger learning rates.

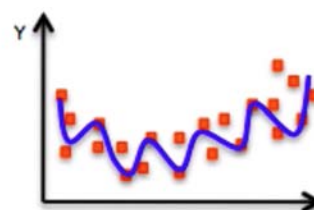
Another classic ML problem: overfitting



Underfitting
Model does not have capacity to fully learn the data



Ideal Fit



Overfitting
Too complex, does not generalize well

Overfitting in Neural Networks 1/3

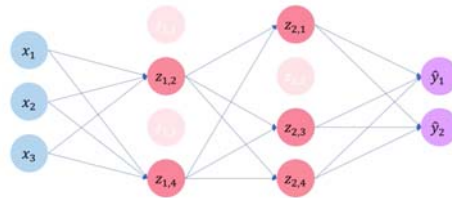
Overfitting #1: The neural network is “too complex”

- Too many hidden layers,
- Too many nodes in each layer.

→ Some neurons learn all the signal and others become useless, which is bad for generalization.

Regularization idea #1: **Drop-out**

Randomly dropping (i.e., setting to 0) some nodes in the network during training to force the network to not rely on just a few nodes.



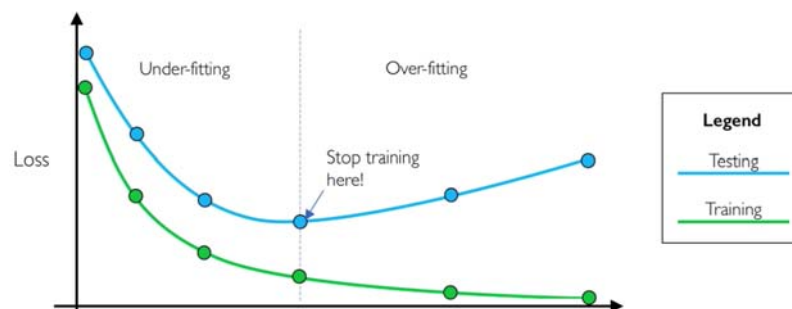
39

Overfitting in Neural Networks 2/3

Overfitting #2: The neural network is “too well trained”.

Regularization idea #2: **Early Stopping**

Stop the training early—after a moderate number of iterations—before the in-sample loss is completely minimized.



40

Overfitting in Neural Networks 3/3

Regularization idea #3: **Embed regularization** ideas in the training of neural networks, in the objective of the non-linear optimization.

- Recall the objective of minimizing the empirical loss:

$$\mathcal{J}(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f(\mathbf{x}_i, \mathbf{W}), y_i)$$

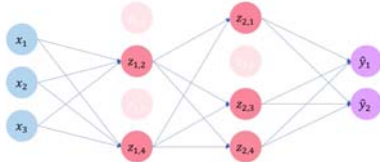
- We can penalize this loss-minimization objective as follows (analog of ridge regression in the context of neural networks)

$$\min \left(\mathcal{J}(\mathbf{W}) + \lambda \times \sum_{i \in \text{layers}} \sum_{k \in \text{nodes of } i} w_{ik}^2 \right)$$

41

Summary: The 3 Regularization Techniques

Regularization 1: Dropout



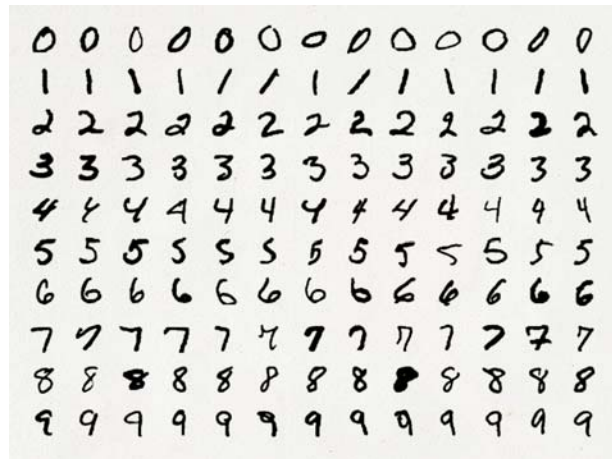
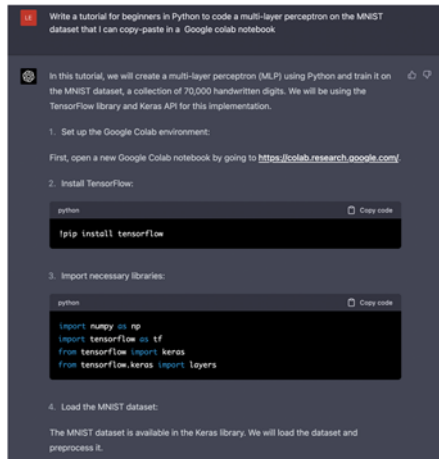
Regularization 2: Early Stopping



Regularization 3: L1 or L2 penalty on the NN weights.

Picture sources: <http://introtodeeplearning.com/>

Quick ChatGPT demo: Classifying MNIST digits



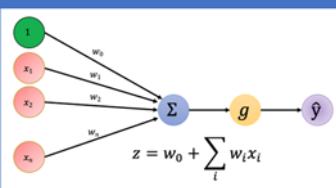
MNIST dataset

43

Summary of the neural networks fundamentals

The Perceptron

- Structure
- Nonlinear activation functions



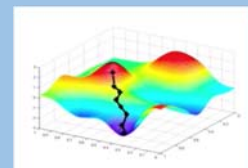
Neural Networks

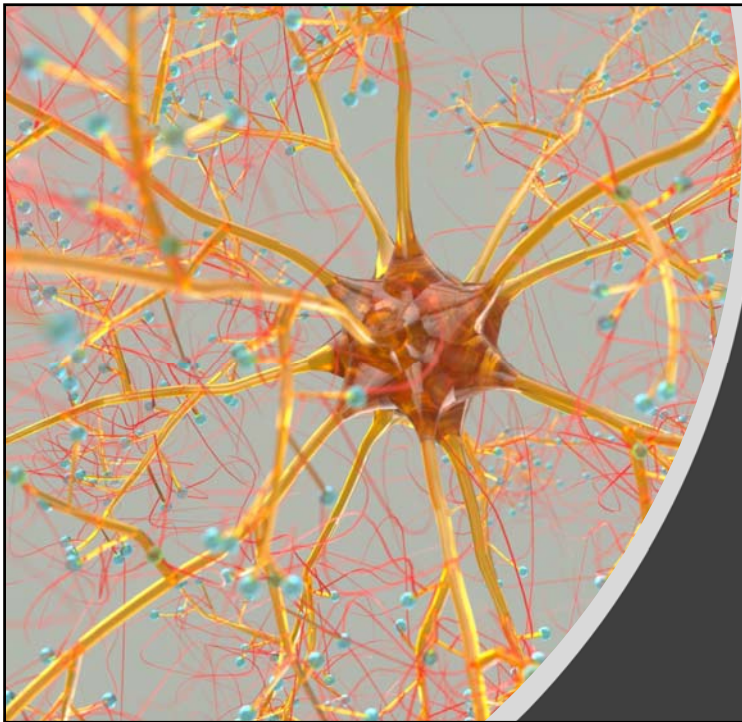
- Stacking layers of perceptrons
- Backpropagation and gradient descent



Training in Practice

- Adaptive learning
- Batching
- Regularization

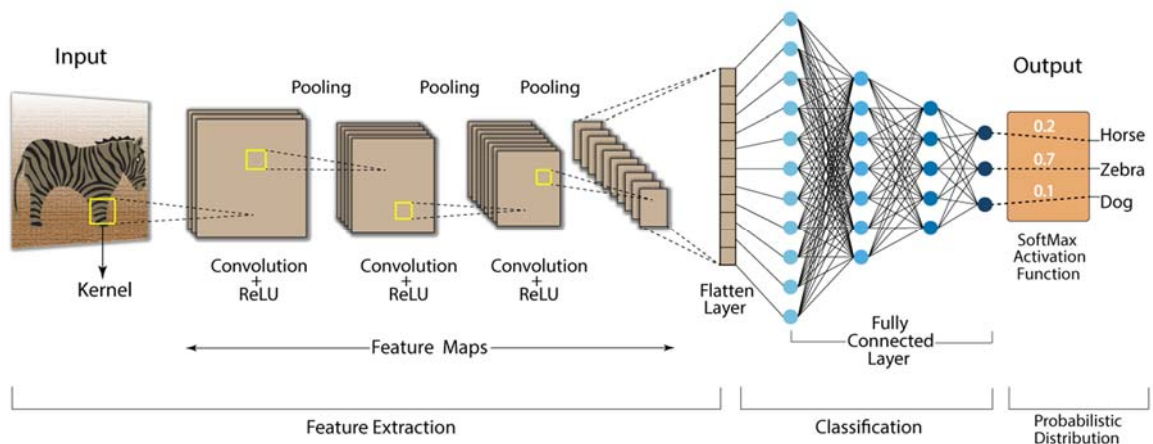




Convolutional Neural Networks

45

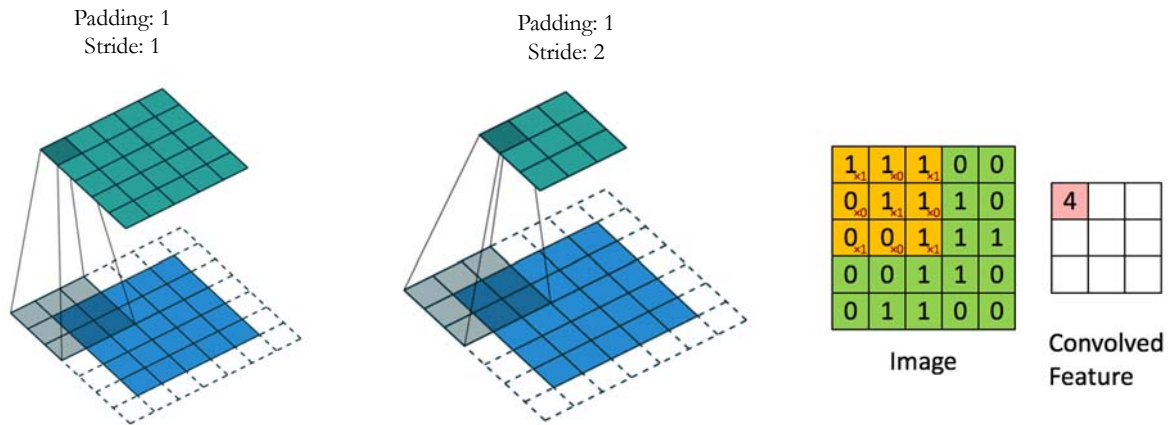
Overall Architecture of a CNN



46

Images source: <https://www.analyticsvidhya.com/blog/2021/05/20-questions-to-test-your-skills-on-cnn-convolutional-neural-networks/>

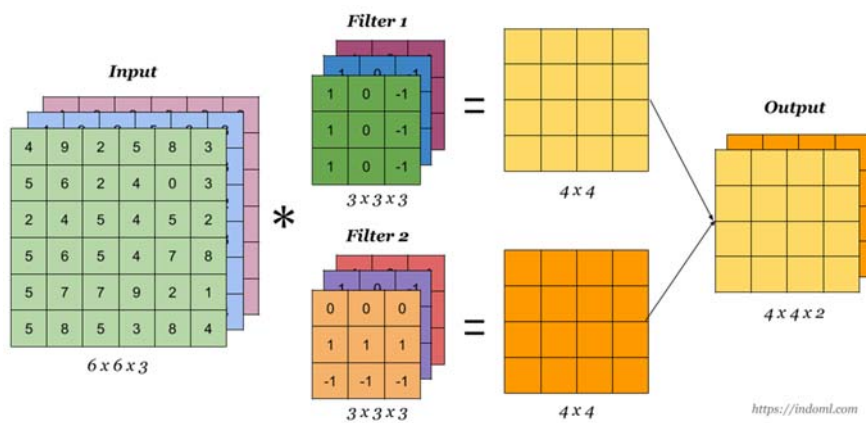
Convolution Operation



47

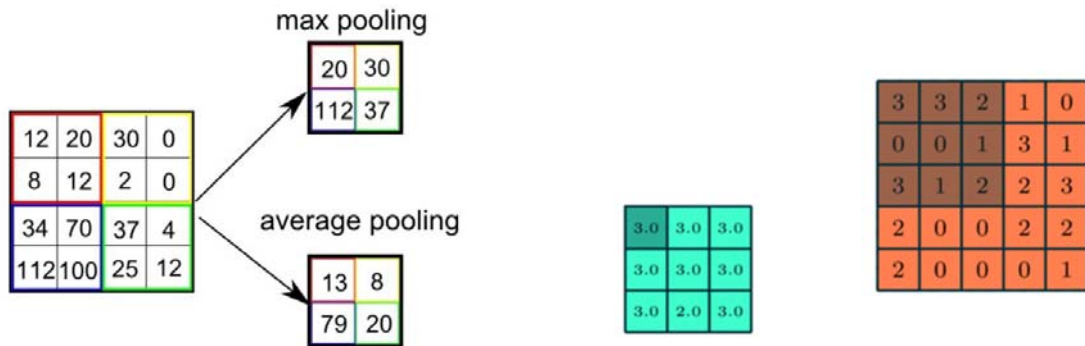
Images source: <https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-easy-way-3bd2b164a53>

3-dimensional data and Conv2D (optional)



48

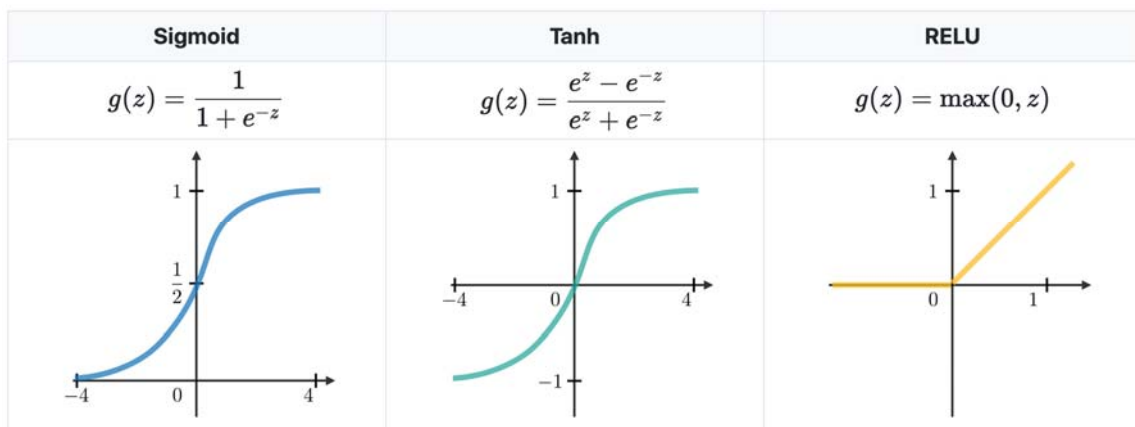
Pooling



49

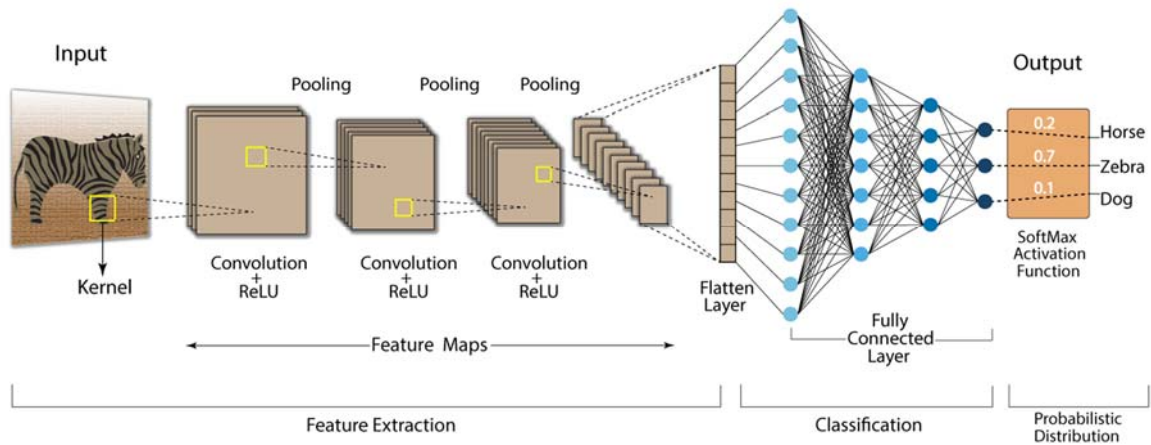
Images source: <https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-eli5-way-3bd2b164a53>

Activation functions



50

Now we can deal with images!

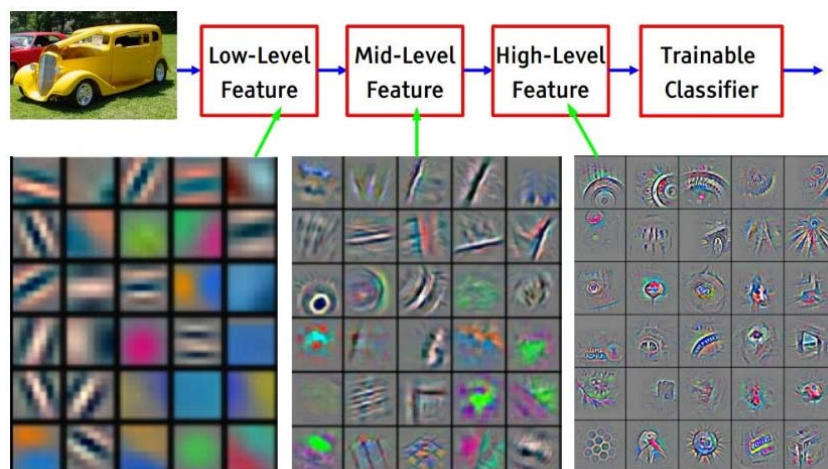


Check the fancy visualization of features learned with ConvNets: <https://microscope.openai.com/models>

Image source: <https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-e85-way-3bd2b1164d53>

51

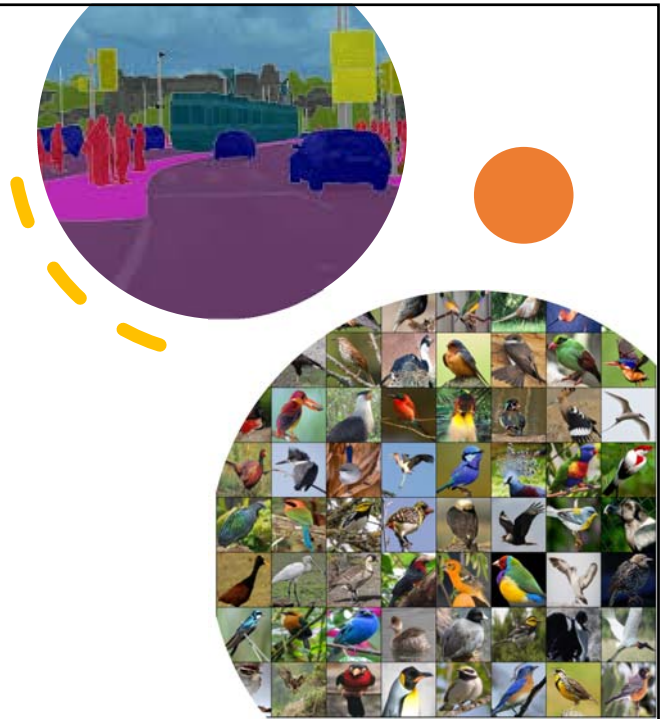
Visualizing CNN filters



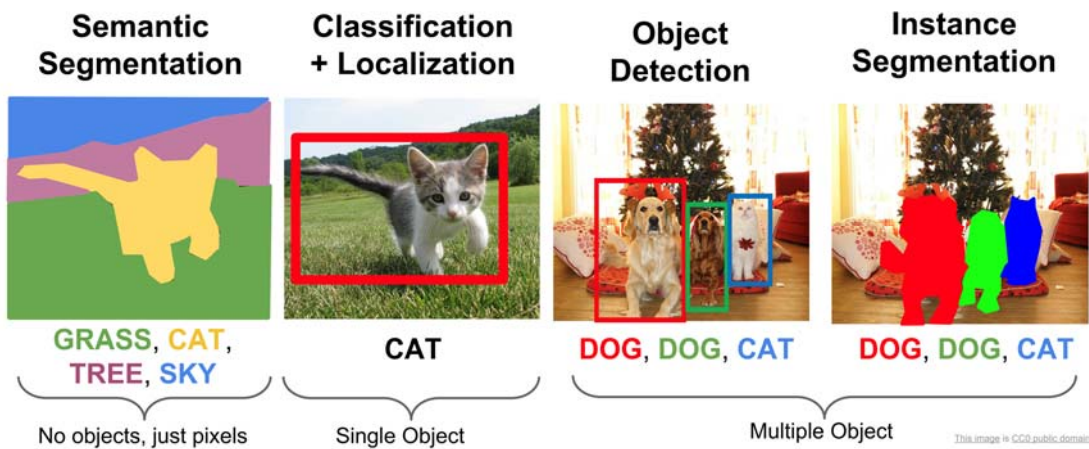
52

What kind of tasks can we perform?

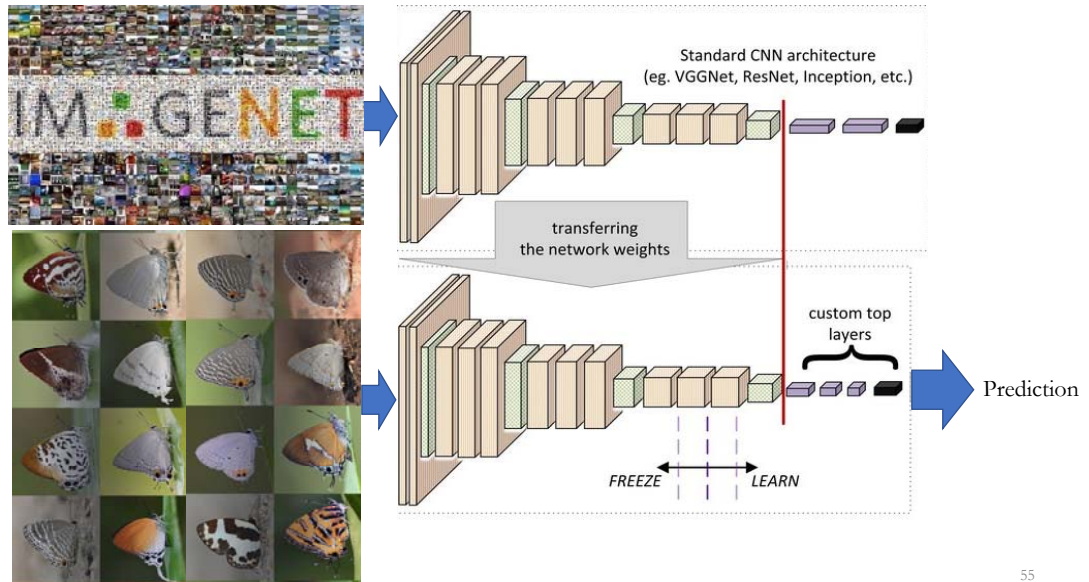
- Classification, Regression
- Segmentation
- Prediction of next image
- Content generation
- Feature extraction
- Descriptions
- ...



Other Computer Vision Tasks

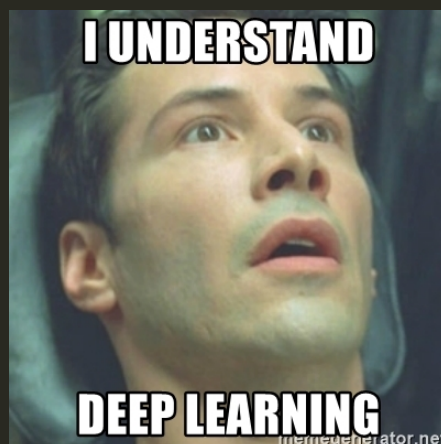


Transfer Learning (optional)



55

Please ask any questions!



56



Resources

MIT Spring Courses that may be offered again in the future:

6.S985 Artificial Intelligence for Business

6.S986 Large Language Models and Beyond

15.S04 Special Seminar in Management (Intro to Deep Learning)

The famous Introduction to Deep Learning class of MIT covers more because this is a full week of content. I highly recommend checking their slides in the topics of interest! <http://introtodeeplearning.com/>

The Computer Vision class of MIT (you have a quick glimpse at the subjects of interest): <http://6.869.csail.mit.edu/sp22/>

The NLP course of MIT: <https://www.mit.edu/~jda/teaching/6.864/>

Berkeley course about DL. Excellent too! I have looked at the slides about Transformers there several times in the past! <https://cs182sp21.github.io/>

Stanford course about CNNs, excellent as well: <http://cs231n.stanford.edu/>