



Cornell Bowers C-IS

College of Computing and Information Science

CS 4782: Introduction to Deep Learning

Course Staff



Varsha Kishore
Instructor



Justin Lovelace
Instructor



Anissa Dallmann
TA



Luke Kulm
TA



Zach Ross
TA



Adhitya Polavaram
TA

ML/AI Courses at Cornell

CS 4780: Introduction to Machine Learning

CS 4756: Robot Learning

CS 4670: Introduction to Computer Vision

CS 4744: Computational Linguistics I

CS 4789: Introduction to Reinforcement Learning

CS 4775: Computational Genetics and Genomics

CS 4740: Natural Language Processing

...

Logistics

- All lectures will be held in person at Stocking Hall 202
- Lectures will be on Tuesdays and Thursdays from 2:55 to 4:10pm
- Aiming to have a small class
 - If you are on the waitlist, come talk to us after class
- Please participate!!

Logistics

- Course website: <https://www.cs.cornell.edu/courses/cs4782/2024sp/>
 - Tentative schedule, homework policies, grading policies, etc. are on the course page
- Office hours are on the course website!
- We also have a Canvas page
 - Links to the Ed Discussion
- No laptops/mobiles/smart devices and other devices in class please
- Notes will be printed

Grading

- Homework (30%)
 - There will be written assignments and coding projects
 - Google Cloud Credits for compute!
 - We recommend doing them in pairs!
 - 2-slip days for every assignment
- Mid-term exam (30%)
 - Will be similar to the homework assignments
- Project (20%)
 - Goal: familiarize yourself with deep learning libraries
 - Implement a method from a recent research paper and reproduce their results
- Participation (20%)
 - Attend classes!
 - Engage in class discussions
 - At the end of each module provide feedback

Academic Integrity

- Do not disclose exact solutions to members from other groups for assignments
 - High-level discussion is allowed
- Cite any external sources
- You can use ChatGPT/BARD/other AI assistants
 - But **add a note** explaining what you used it for and how you used it

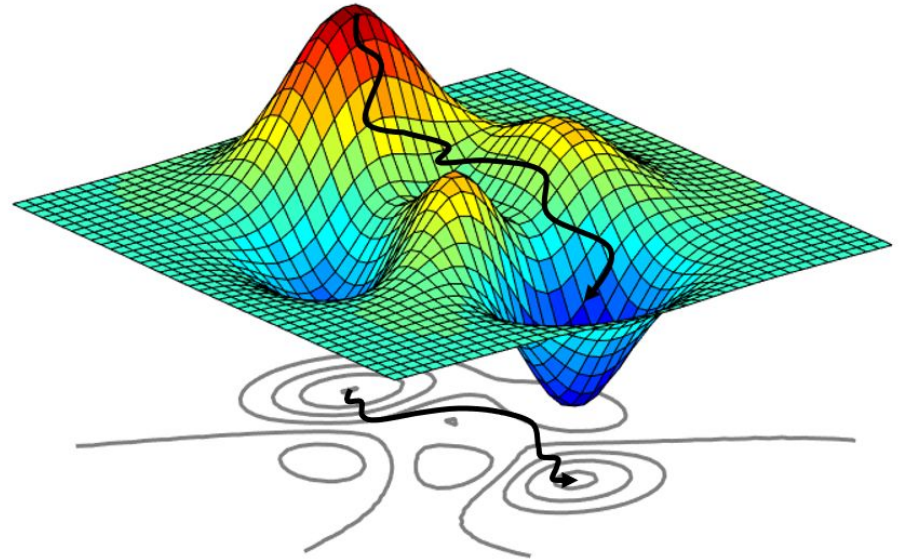
Course Objectives

By the end of the course you will be able to...

1. Design, train, and evaluate deep neural networks
2. Apply deep learning techniques to solve real-world problems in computer vision, natural language processing, and other complex domains
3. Critically evaluate pros/cons of different model architectures
4. Read and understand research in deep learning
5. Understand the core design principles behind leading deep learning systems like GPT-4, DALL-E 2/3, and Stable Diffusion

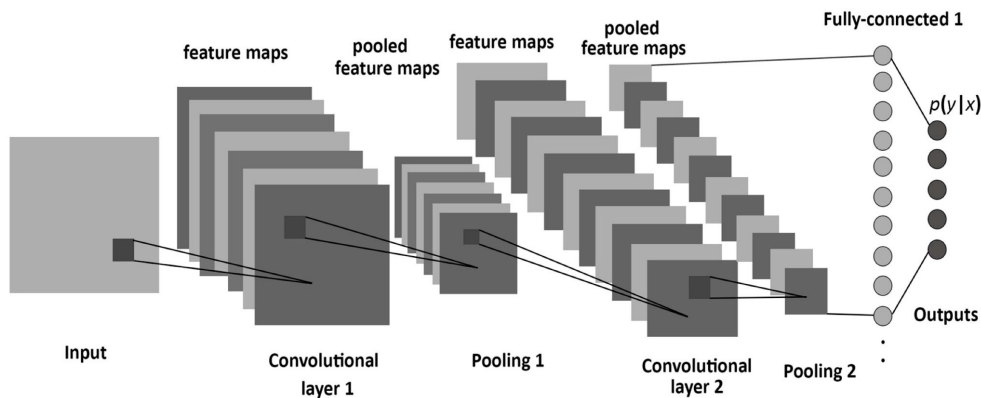
Training Neural Networks

- Optimization algorithms - gradient descent, SGD, AdaGrad, Adam
- Learning rate scheduling
- Hyperparameter Optimization
- Regularization



Computer Vision

- Convolutional neural networks
- Different convolutional architectures - vanilla CNN, LeNet, ResNet, DenseNets



<https://www.mdpi.com/1099-4300/19/6/242>



Natural Language Processing

- Word Embeddings
- Recurrent Neural Networks
 - RNNs/ LSTMs
- Attention and Transformers
- Large Language Models (LLMs)

Explaining a Joke

Input: Did you see that Google just hired an eloquent whale for their TPU team? It showed them how to communicate between two different pods!

<https://arxiv.org/abs/2204.02311>

Natural Language Processing

- Word Embeddings
- Recurrent Neural Networks
- RNNs/ LSTMs
- Attention and Transformers
- Large Language Models (LLMs)

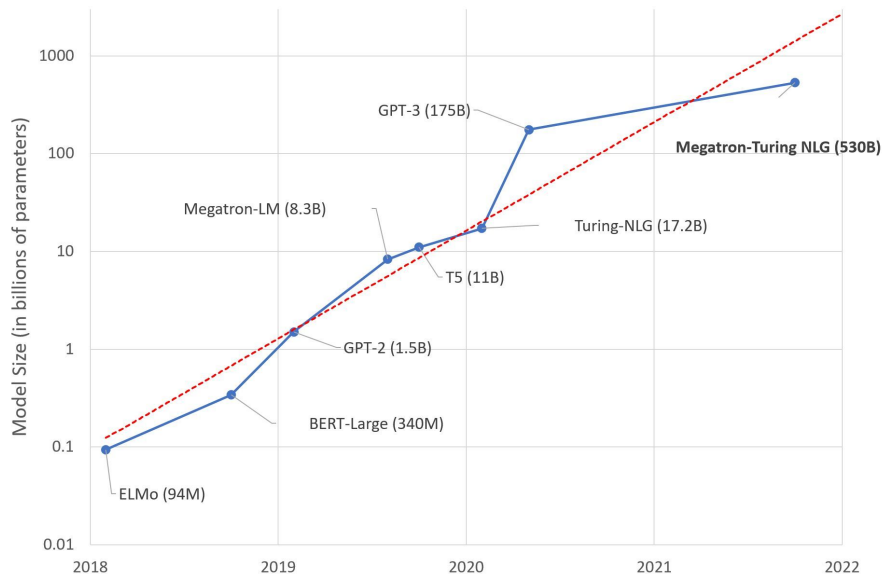
Explaining a Joke

Input: Did you see that Google just hired an eloquent whale for their TPU team? It showed them how to communicate between two different pods!

Model Output: TPUs are a type of computer chip that Google uses for deep learning. A "pod" is a group of TPUs. A "pod" is also a group of whales. The joke is that the whale is able to communicate between two groups of whales, but the speaker is pretending that the whale is able to communicate between two groups of TPUs.

Natural Language Processing

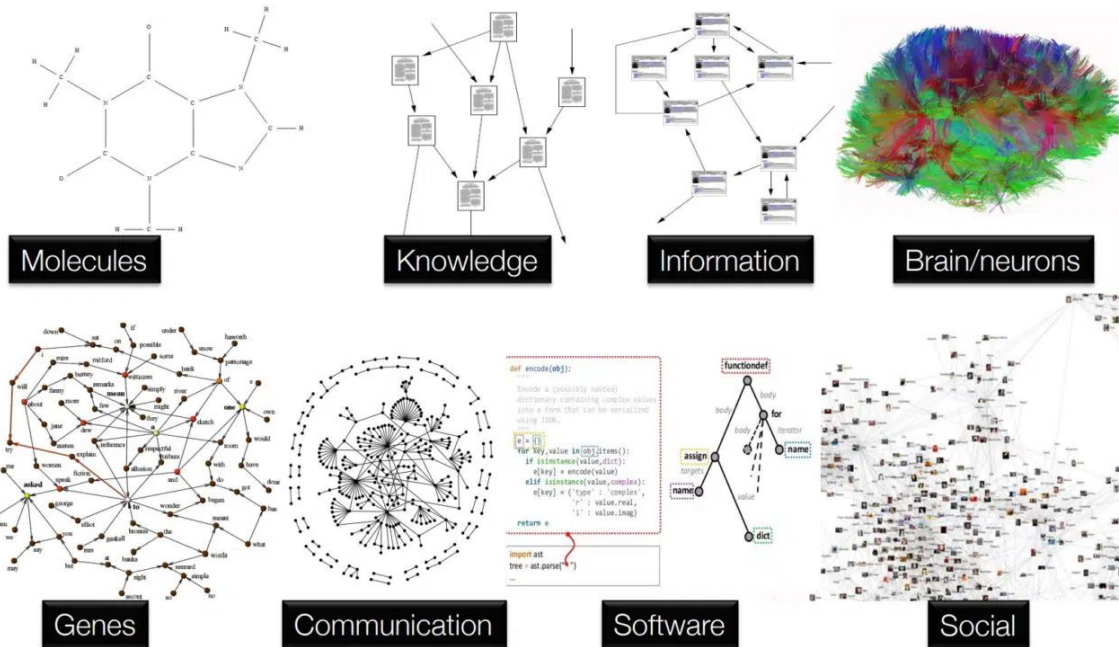
- Word Embeddings
- Recurrent Neural Networks
- RNNs/ LSTMs
- Attention and Transformers
- Large Language Models (LLMs)



<https://huggingface.co/blog/large-language-models>

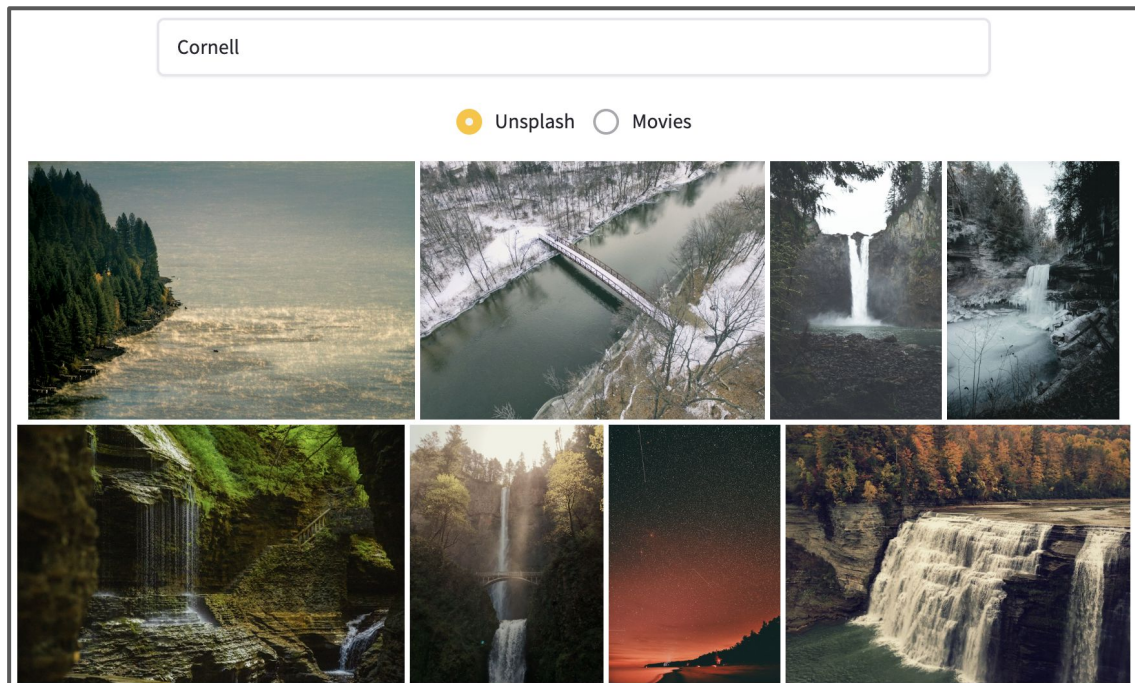
Graph Neural Networks

- Neural networks for data represented as graphs!

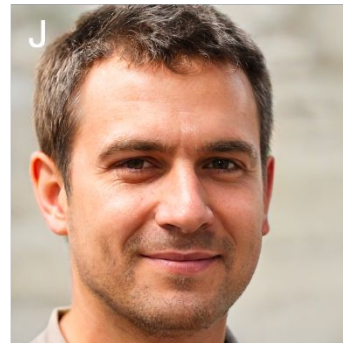
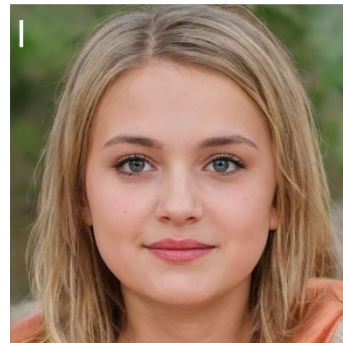
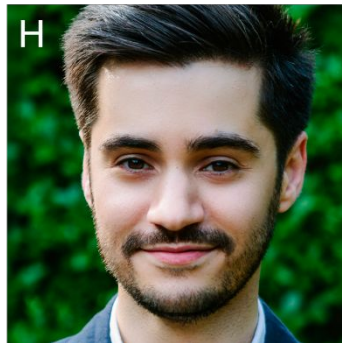
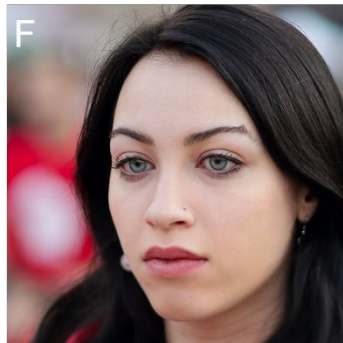
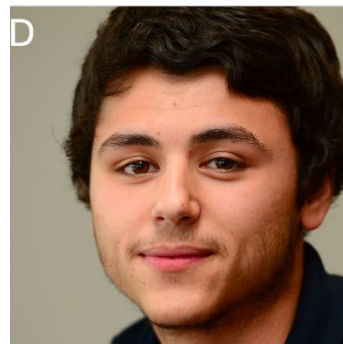
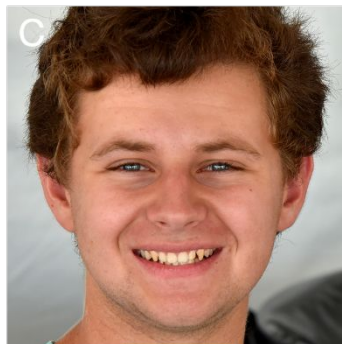
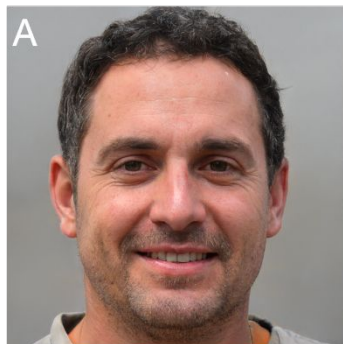


Modern Vision Networks

- Vision Transformers (ViTs)
- Vision Pre-Training
 - (Supervised, Self-supervised)
- Vision-Language Models

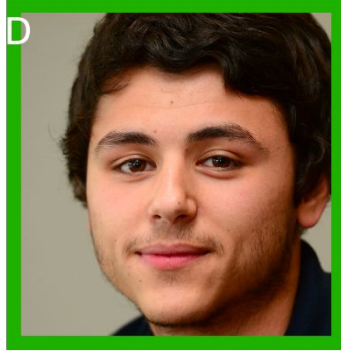


Real or Fake?



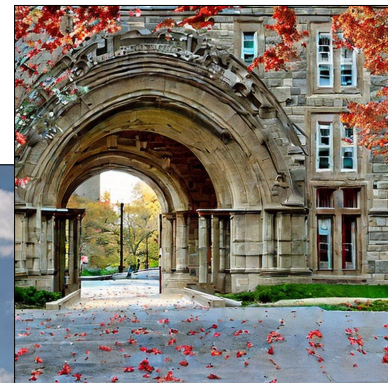
Cornell Bowers C-IS

Real or Fake?



Generative Models

- U-Nets
- Variational Autoencoders (VAEs)
- Generative Adversarial Networks (GANs)
- Diffusion Models
- Multi-Modal Diffusion



Reinforcement Learning

Technique for an agent to learn in an interactive environment by testing different actions and obtaining feedback from its experiences.

- Markov Decision Process
- Q-learning/Deep Q-learning
- Policy Gradients
- Exploration strategies
- RL from Human Feedback



AI in Human Society



ARTIFICIAL INTELLIGENCE

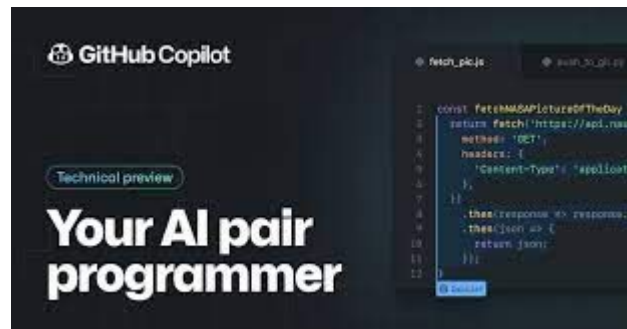
A program that can sense, reason, act, and adapt

MACHINE LEARNING

Algorithms whose performance improve as they are exposed to more data over time

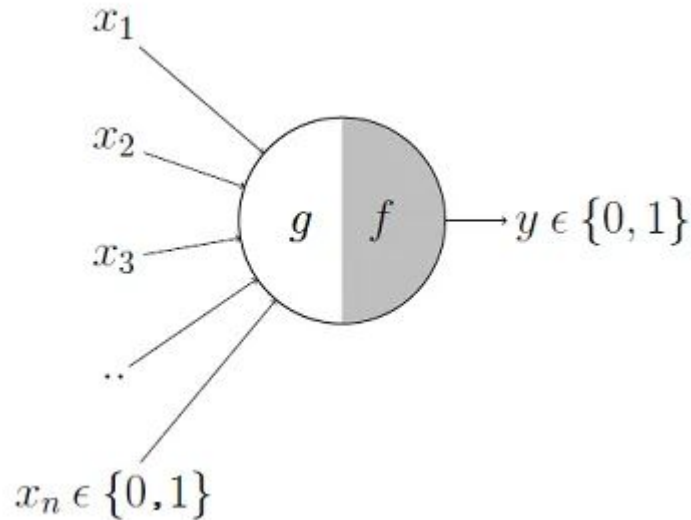
DEEP LEARNING

Subset of machine learning in which multilayered neural networks learn from vast amounts of data



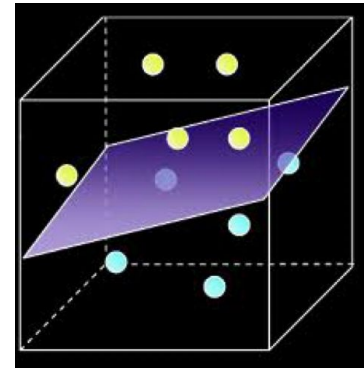
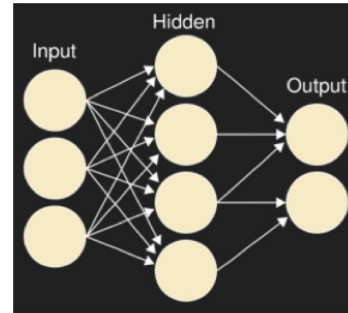
McCulloch-Pitts Neuron

Computational model of a neuron that was proposed by Warren McCulloch (neuroscientist) and Walter Pitts (logician) in 1943.

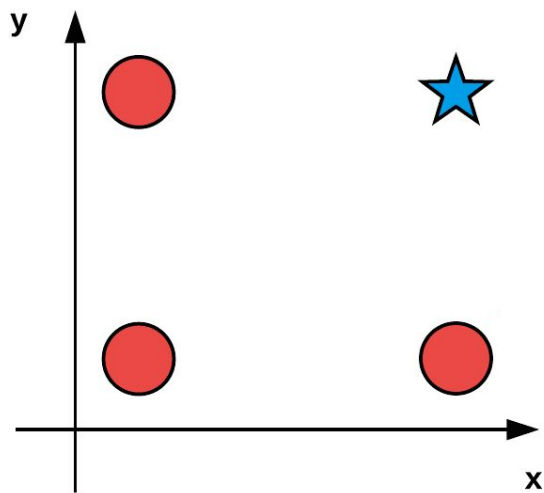


Perceptron (1957)

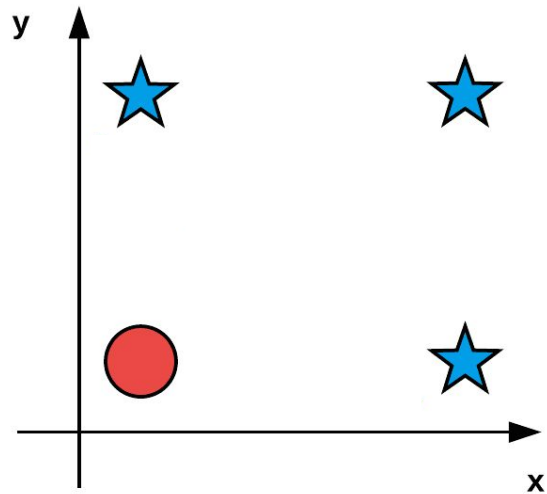
- Linear classifier, predecessor to Neural Network
- Trained with the perceptron update rule
- Invented @ Cornell University
 - First task: Recognize the Cornell “C” Logo



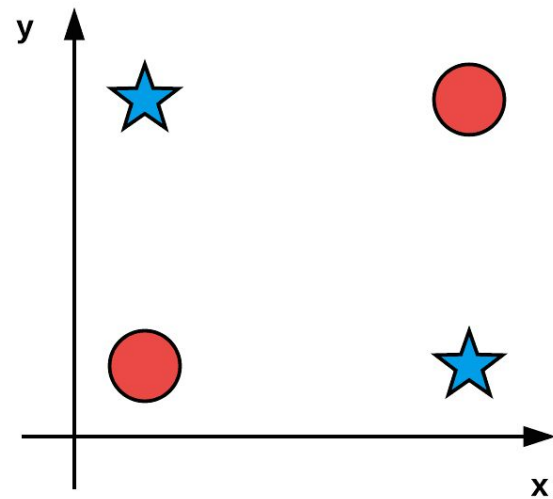
AND



OR

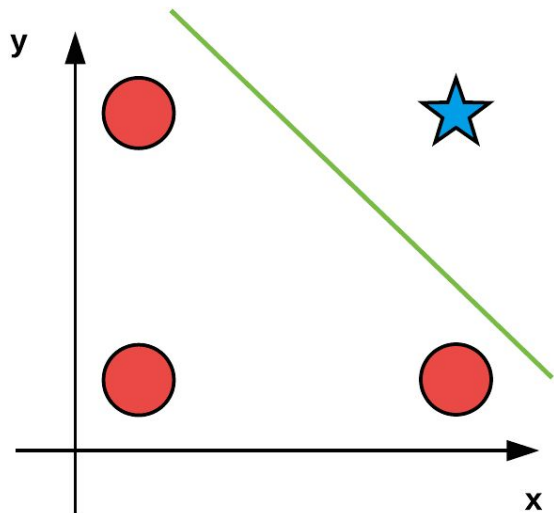


XOR

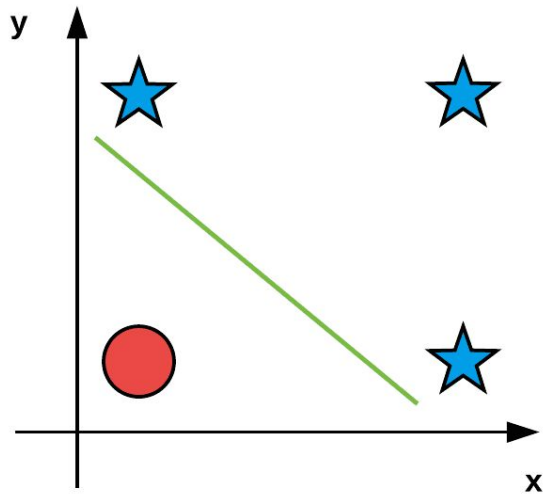


Cornell Bowers C-IS

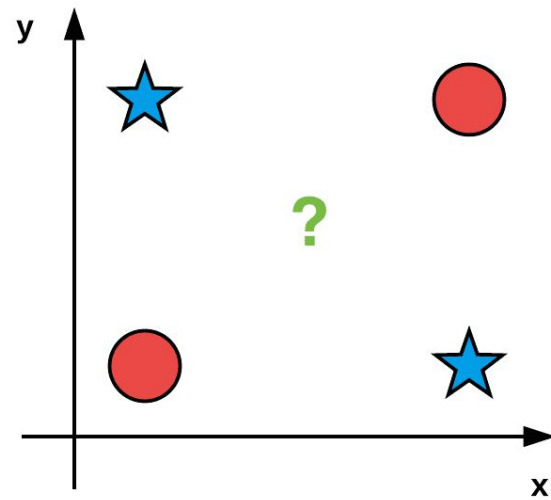
AND



OR

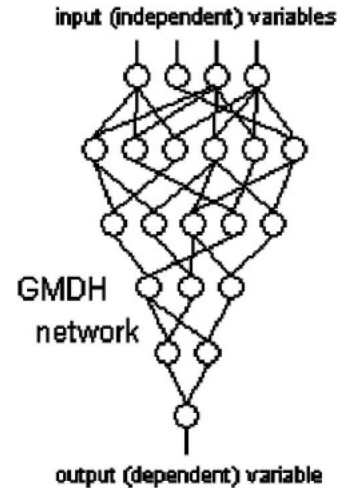


XOR



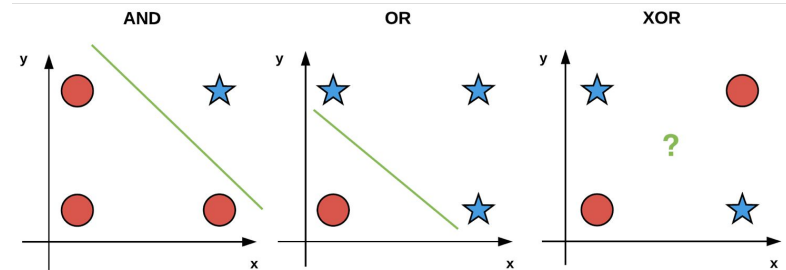
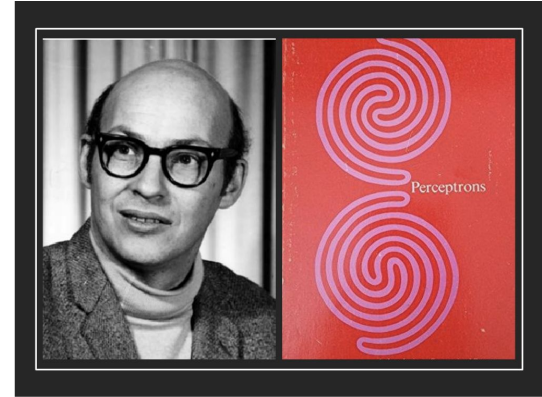
Multi-layer neural networks

- Multi-Layer Perceptron, Rosenblatt (around 1965)
- Alexey Grigoryevich Ivakhnenko 1965 Group Method of Data Handling (GMDH)
 - 1971 Eight Layer Neural Nets with skip connections!



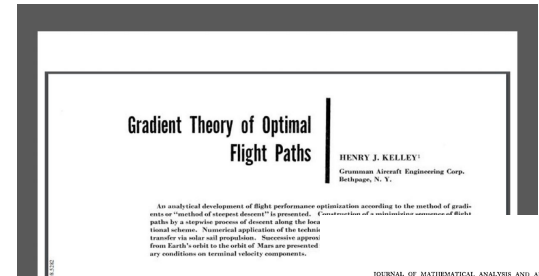
AI Winter (1974-1980)

- (1969) Minsky & Papert “killed” AI
- Burst huge expectation bubble
- Speech understanding / translation fails
- UK and US stop funding AI research



Backprop

- 1960 Henry J. Kelly Initial formulation in control theory (rocket science)
- 1962 Stuart Dreyfuss (use of chain rule)
- 1979 Seppo Linnainmaa (modern backdrop with automatic differentiation [not in context of neural nets])
- 1982 Paul Werbos proposes backprop for artificial Neural Networks in PhD thesis
- 1986 Rumelhart, Hinton, Williams (coin the term “back-propagation”) make the algorithm popular (published in Nature)



The Numerical Solution of Variational Problems
STUART DREYFUS
Computer Science Department,
The RAND Corporation, Santa Monica, California
Submitted by Richard Bellman

Seppo Linnainmaa

ALGORITMIN KUNNATIIIVIEREN PÄRISISTYVIERE
YESITTAISIT**

I. INTRODUCTION

ve derived various conditions that must
solution of a variational problem. In [1]
problem of Lagrange. In [2] we studied
and the conditions implied by the intru-
either the shape of the solution curve
he region in which it could lie (state
we shall discuss the numerical solution
conventional and inequality-constrained
at, until recently, has been the usual
gradient technique that has proved very
h is known to a few practitioners of the
Bryson [4]. Our derivation, using essen-
iques of dynamic programming that we
ew and simple. We conclude this paper
l state of a variant of the classical state
variable inequality constraint has
analytic solution is still known.
is familiar with such dynamic program-

Pro gradu-

Learning representations by back-propagating errors

David E. Rumelhart*, Geoffrey E. Hinton†
& Ronald J. Williams*

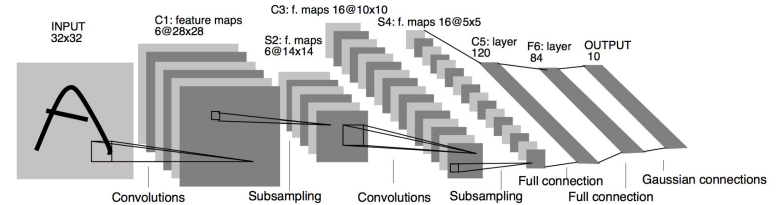
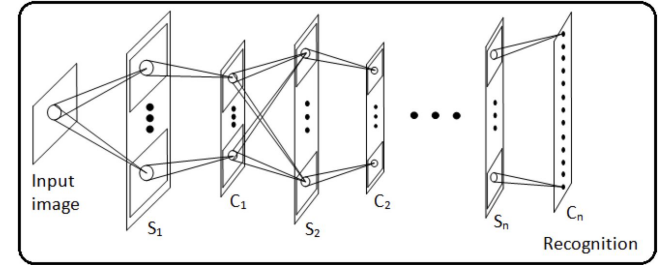
* Institute for Cognitive Science, C-015, University of California,
San Diego, La Jolla, California 92093, USA
† Department of Computer Science, Carnegie-Mellon University,
Pittsburgh, Philadelphia 15213, USA

We describe a new learning procedure, back-propagation, for networks of neurone-like units. The procedure repeatedly adjusts the weights of the connections in the network so as to minimize a measure of the difference between the actual output vector of the net and the desired output vector. As a result of the weight adjustments, internal 'hidden' units which are not part of the input or output come to represent important features of the task domain, and the regularities in the task are captured by the interactions of these units. The ability to create useful new features distinguishes back-propagation from earlier, simpler methods such as the perceptron-convergence procedure¹.

There have been many attempts to design self-organizing neural networks. The aim is to find a powerful synaptic

ConvNets

- 1979 Kunihiro Fukushima invents Neocognitron
 - Heavily inspired by human Visual Cortex
 - Alternates between Simple Cells / Complex Cells
 - Unsupervised
- 1986 Yann LeCun introduces BackProp to ConvNets for Handwritten Digits (creates MNIST)



Recurrent Neural Nets

- 1982 John Hopfield “Hopfield Networks”
- 1991 Sepp Hochreiter formulates Vanishing Gradient Problem
- 1997 S. Hochreiter and Jürgen Schmidhuber publish “Long Short-Term Memory” (LSTM)
 - <https://web.archive.org/web/20231216143334/https://people.idsia.ch/~juergen/ai-priority-disputes.html>

Proc. Natl. Acad. Sci. USA
Vol. 79, pp. 2554–2558, April 1982
Biophysics

Neural networks and physical systems with emergent collective computational abilities

(associative memory/parallel processing/categorization/content-addressable memory/fail-soft devices)

J. J. HOPFIELD

Division of Chemistry and Biology, California Institute of Technology, Pasadena, California 91125; and Bell Laboratories, Murray Hill, New Jersey 07974

Contributed by John J. Hopfield, January 15, 1982

ABSTRACT Computational properties of use to biological organisms or to the construction of computers can emerge as collective properties of systems having a large number of simple equivalent components (or neurons). The physical meaning of content-addressable memory is described by an appropriate phase space flow of the state of a system. A model of such a system is given, based on aspects of neurobiology but readily adapted to integrated circuits. The collective properties of this model produce a content-addressable memory which correctly yields an entire memory from any subpart of sufficient size. The algorithm for the time evolution of the state of the system is based on asynchronous parallel processing. Additional emergent collective properties include some capacity for generalization, familiarity recognition, categorization, error correction, and time sequence retention. The collective properties are only weakly sensitive to details of the modeling or the failure of individual devices.

Given the dynamical electrochemical properties of neurons and their interconnections (synapses), we readily understand schemes that use a few neurons to obtain elementary useful biological behavior (1–3). Our understanding of such simple circuits in electronics allows us to plan larger and more complex circuits which are essential to large computers. Because evolution has no such plan, it becomes relevant to ask whether the ability of large collections of neurons to perform “computational” tasks may in part be a spontaneous collective consequence of having a large number of interacting simple neurons.

In physical systems made from a large number of simple elements, interactions among large numbers of elementary components yield collective phenomena such as the stable magnetic orientations and domains in a magnetic system or the vortex patterns in fluid flow. Do analogous collective phenomena in

calized content-addressable memory or categorizer using extensive asynchronous parallel processing.

The general content-addressable memory of a physical system

Suppose that an item stored in memory is “H. A. Kramers & C. H. Wannier *Phys. Rev.* 60, 252 (1941).” A general content-addressable memory would be capable of retrieving this entire memory item on the basis of sufficient partial information. The input “& Wannier, (1941)” might suffice. An ideal memory could deal with errors and retrieve this reference even from the input “Wannier, (1941)”. In computers, only relatively simple forms of content-addressable memory have been made in hardware (10, 11). Sophisticated ideas like error correction in accessing information are usually introduced as software (10).

There are classes of physical systems whose spontaneous behavior can be used as a form of general (and error-correcting) content-addressable memory. Consider the time evolution of a physical system that can be described by a set of general coordinates. A point in state space then represents the instantaneous condition of the system. This state space may be either continuous or discrete (as in the case of N Ising spins).

The equations of motion of the system describe a flow in state space. Various classes of flow patterns are possible, but the systems of use for memory particularly include those that flow toward locally stable points from any where within regions around those points. A particle with frictional damping moving in a potential well with two minima exemplifies such a dynamics.

If the flow is not completely deterministic, the description is more complicated. In the two-well problems above, if the frictional force is characterized by a temperature, it must also produce a random driving force. The limit points become small

Universal Approximation

- 1989 George Cybenko proves universal approximation of single hidden-layer neural networks
- Also yields wide-spread believe that more than one layer is unnecessary

Math. Control Signals Systems (1989) 2: 303–314

Mathematics of Control,
Signals, and Systems
© 1989 Springer-Verlag New York Inc.

Approximation by Superpositions of a Sigmoidal Function*

G. Cybenko†

Abstract. In this paper we demonstrate that finite linear combinations of compositions of a fixed, univariate function and a set of affine functionals can uniformly approximate any continuous function of n real variables with support in the unit hypercube; only mild conditions are imposed on the univariate function. Our results settle an open question about representability in the class of single hidden layer neural networks. In particular, we show that arbitrary decision regions can be arbitrarily well approximated by continuous feedforward neural networks with only a single internal, hidden layer and any continuous sigmoidal nonlinearity. The paper discusses approximation properties of other possible types of nonlinearities that might be implemented by artificial neural networks.

Key words. Neural networks, Approximation, Completeness.

1. Introduction

A number of diverse application areas are concerned with the representation of general functions of an n -dimensional real variable, $x \in \mathbb{R}^n$, by finite linear combinations of the form

$$\sum_{j=1}^N \alpha_j \sigma(y_j^T x + \theta_j), \quad (1)$$

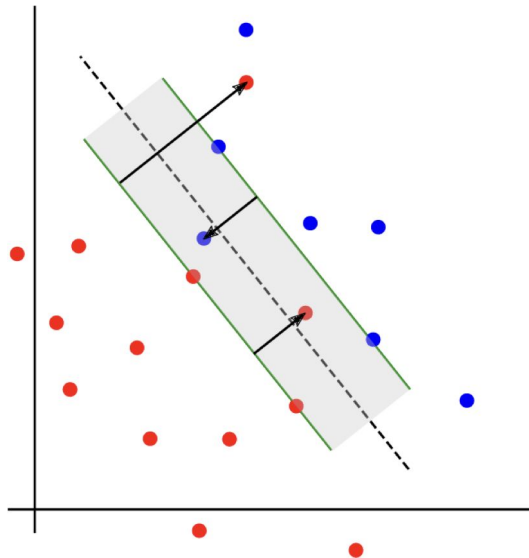
where $y_j \in \mathbb{R}^n$ and $\alpha_j, \theta_j \in \mathbb{R}$ are fixed. (y_j^T is the transpose of y_j so that $y_j^T x$ is the inner product of y_j and x .) Here the univariate function σ depends heavily on the context of the application. Our major concern is with so-called sigmoidal σ 's:

$$\sigma(t) \rightarrow \begin{cases} 1 & \text{as } t \rightarrow +\infty, \\ 0 & \text{as } t \rightarrow -\infty. \end{cases}$$

Such functions arise naturally in neural network theory as the activation function

Summer of SVMs 1995-2008

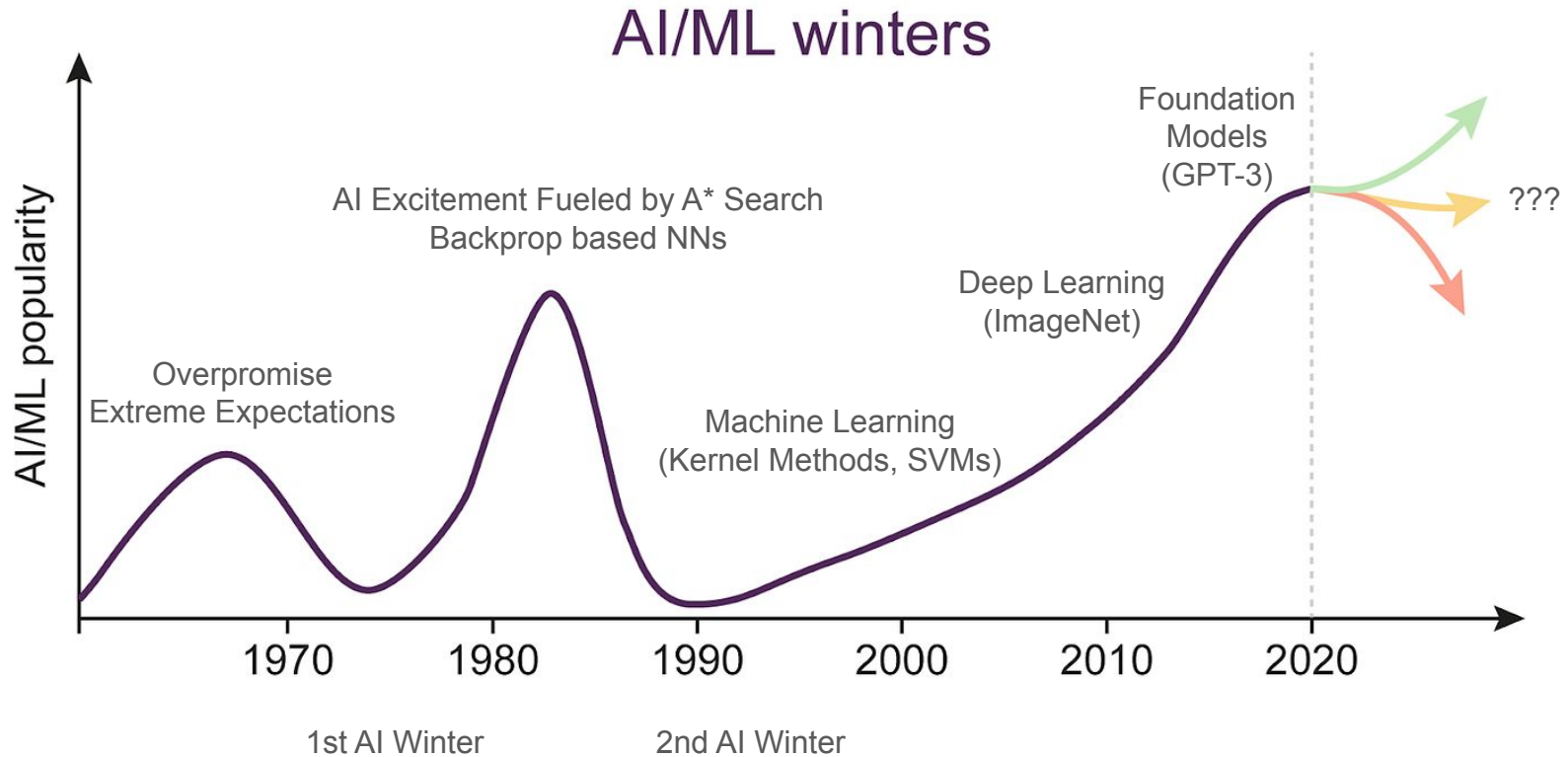
- 1993-1995 Corinna Cortes, Isabella Guyon, Vladimir Vapnik invent Support Vector Machines
- Mid 2000s ICML and NeuRIPS (NIPS) exclusively papers on non-neural network approaches
 - Mostly SVM, Graphical Models, Boosting
 - These algorithms are more efficient, easier to train / modify, have strong theoretical guarantees / frameworks



Neural Network Resurgence (2010s)

- Relentless effort by Hinton, Bengio, LeCun: Kept pushing Neural Nets when they were not cool - but did not join other communities (e.g. ICANN)
- Invent Deep Belief Nets in effort to attract experts in Graphical Models (mimics Graphical Models)
- Rename Neural Nets as “Deep Learning” (in effort to brand SVMs as “shallow”)
- Create ICLR as a venue to accept research on Neural Nets
- 2007 NeuRIPS Workshop on Deep Learning (rejected, changed to Hinton’s 60th birthday party)
- 2009 Fei-Fei Li creates ImageNet (after Caltech 4, 101, 256)
- 2012 Hinton’s deep network research creates AlexNet

Public Perception of AI/ML



Task: Predict whether it is winter from an image.

Thanks!

- If you have received a permission number
 - Enroll today if you'd like to take the course
- We will start sending out permission numbers to people on the waitlist later this week
- If you have not received a permission number and want to enroll
 - Come talk to us after class