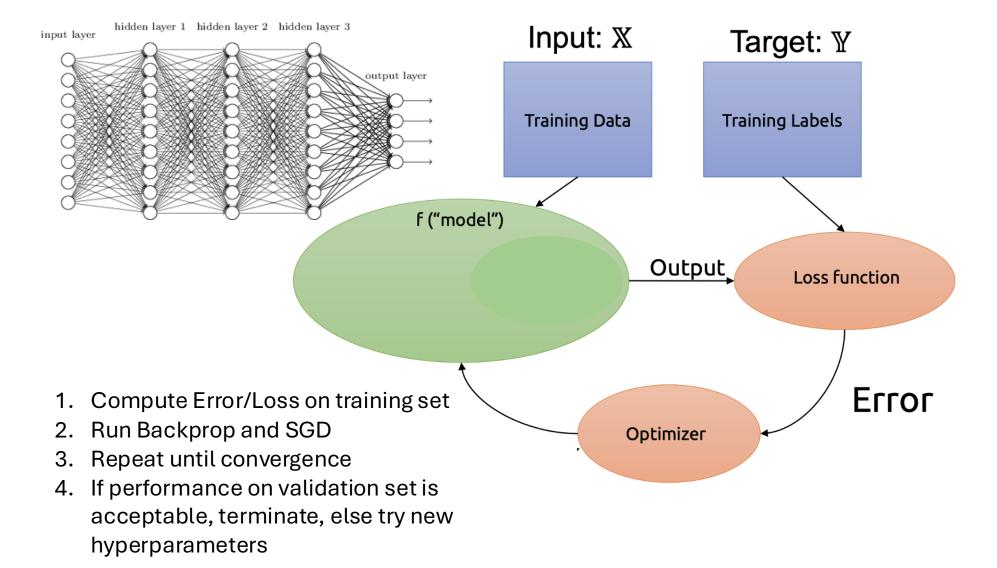
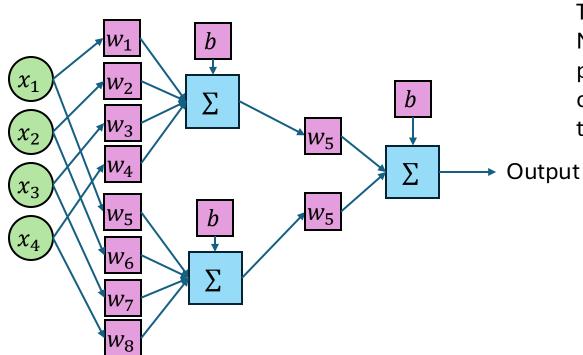


Recap: MLPs



Hyperparameters



The **parameters** of a Neural Network are what is trained (e.g., weights and biases).

The **hyperparameters** of a Neural Network are the parameters that **you** have control of that control that training.

Normalization: Shift each feature to be of a similar scale

Normalization: Shift each feature to be of a similar scale

Apply: $x' = \frac{x - min(x)}{min(x) - max(x)}$ to every feature in dataset, shifting each feature to lie within [0, 1].

Normalization: Shift each feature to be of a similar scale

Apply: $x' = \frac{x - min(x)}{min(x) - max(x)}$ to every feature in dataset, shifting each feature to lie within [0, 1].

Before normalization: different Distance A: 10 inches Distance A: 1000 inches features can be different scales, Distance B: 25 inches Distance B: 2500 inches Distance C: 4000 inches

After normalization: Each feature Normalized Features: Normalized Features:

has same impact on model [0, 0.5, 1] [0, 0.5, 1]

Normalization: Shift each feature to be of a similar scale

Apply: $x' = \frac{x - min(x)}{min(x) - max(x)}$ to every feature in dataset, shifting each feature to lie within [0, 1].

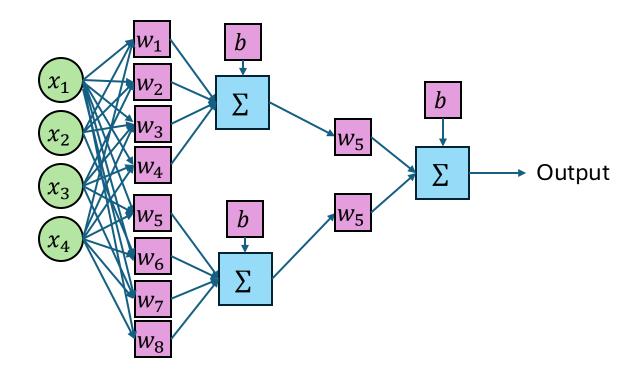
Before normalization: different features can be different scales, different units, etc.

After normalization: Each feature has same impact on model

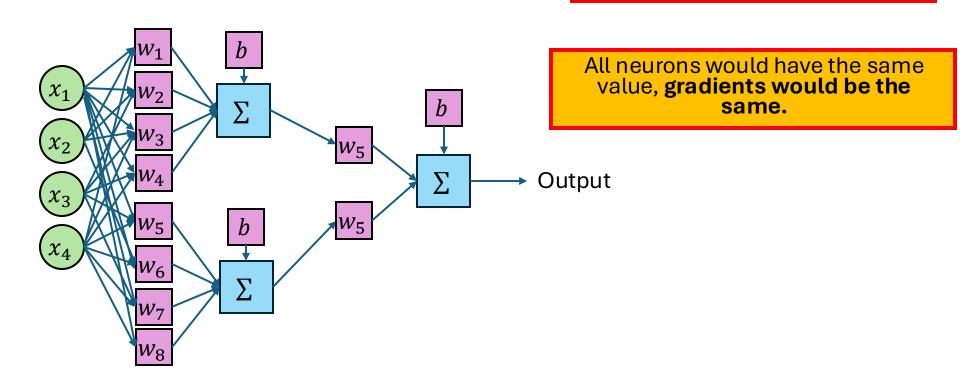
Gradient Descent converges faster when working with normalized data!

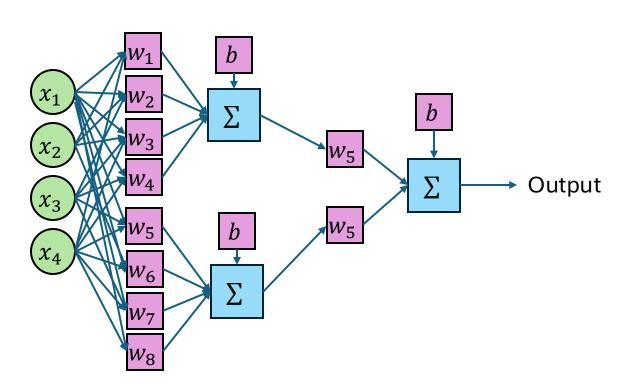
Ioffe, Sergey; Christian Szegedy (2015). "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift".

What if we begin with all parameters set to 0?

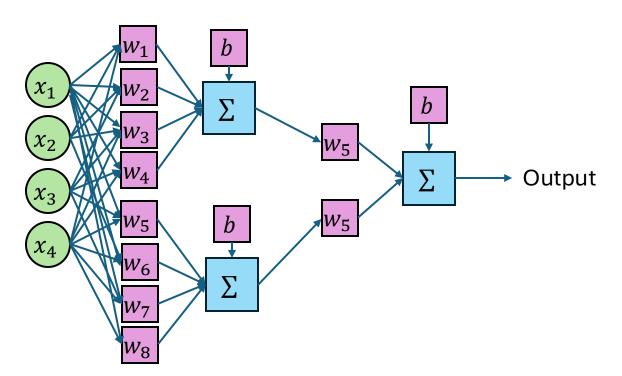


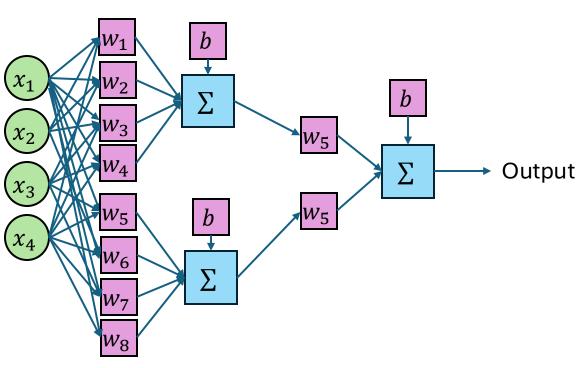
What if we begin with all parameters set to 0?





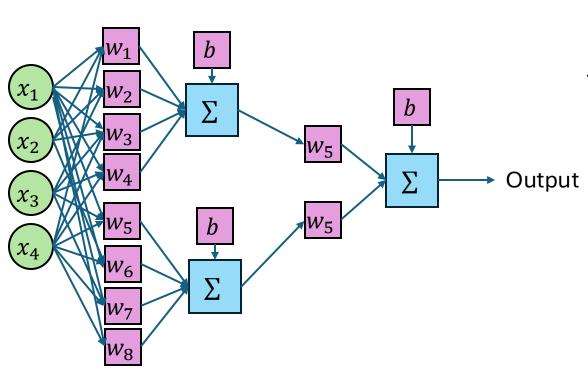
Idea #1: Uniform random weights between -1 and 1 (works fine)





Idea #1: Uniform random weights between -1 and 1 (works fine)

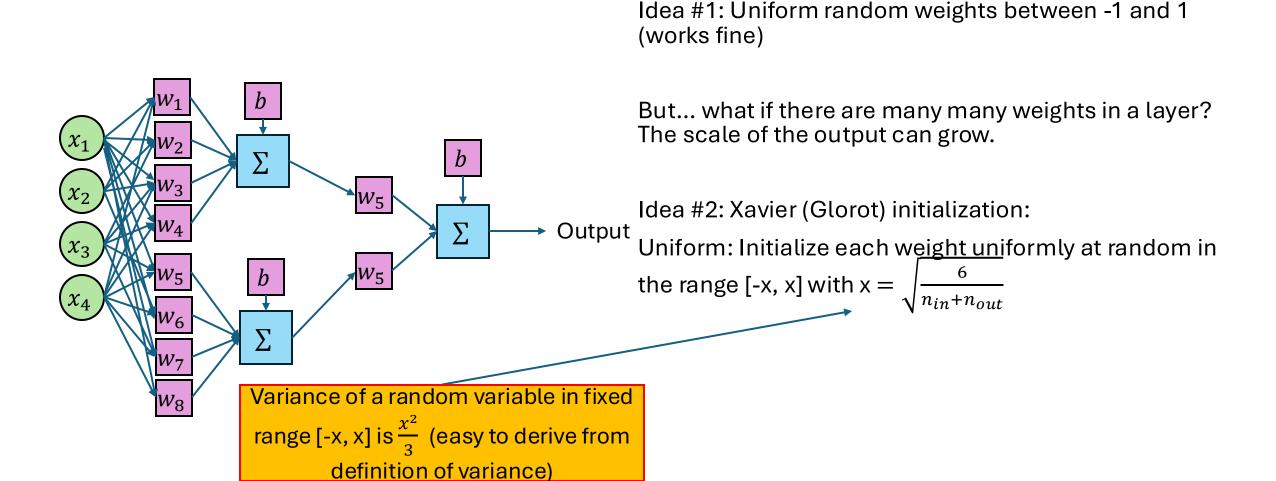
But... what if there are many many weights in a layer? The scale of the output can grow, variance of output increases

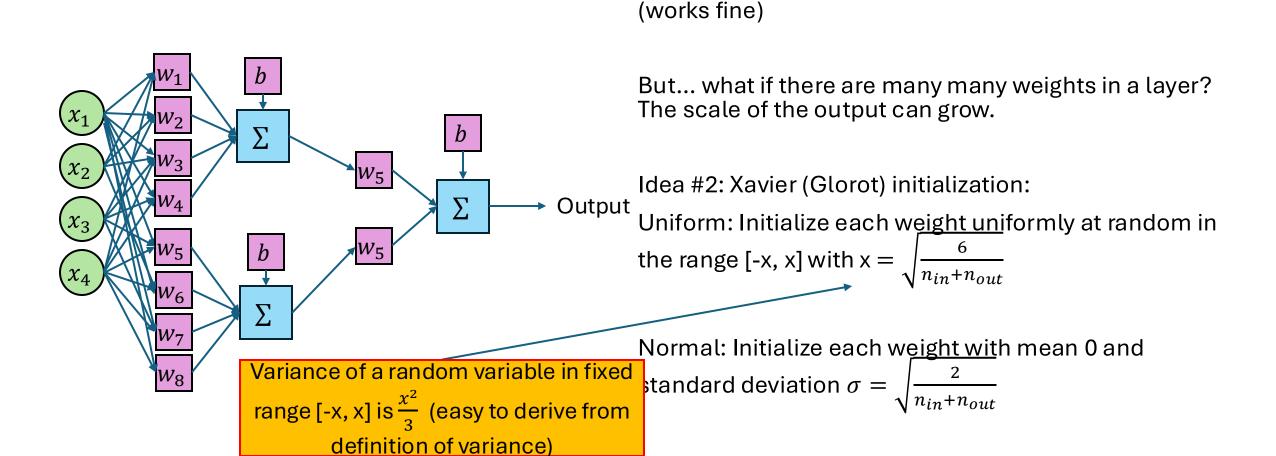


Idea #1: Uniform random weights between -1 and 1 (works fine)

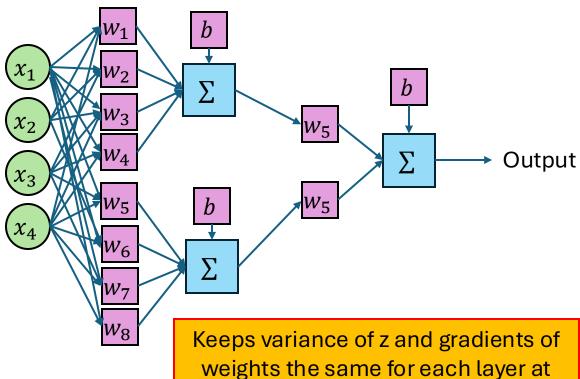
But... what if there are many many weights in a layer? The scale of the output can grow.

Idea #2: Xavier (Glorot) initialization:





Idea #1: Uniform random weights between -1 and 1



initialization.

Idea #1: Uniform random weights between -1 and 1 (works fine)

But... what if there are many many weights in a layer? The scale of the output can grow.

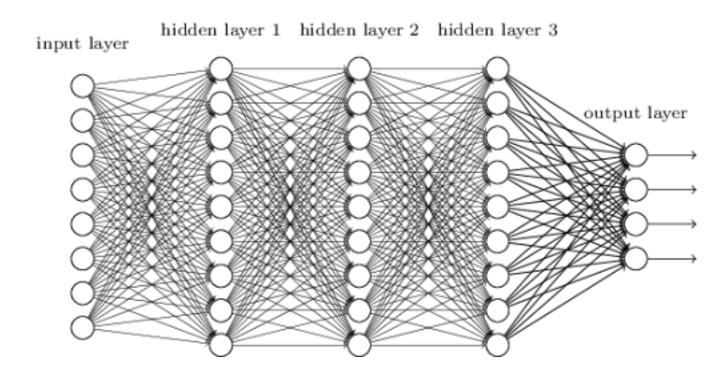
Idea #2: Xavier (Glorot) initialization:

Uniform: Initialize each weight uniformly at random in the range [-x, x] with $x = \sqrt{\frac{6}{n_{in} + n_{out}}}$

Normal: Initialize each weight with mean 0 and standard deviation $\sigma = \sqrt{\frac{2}{n_{in} + n_{out}}}$

Hidden Layers

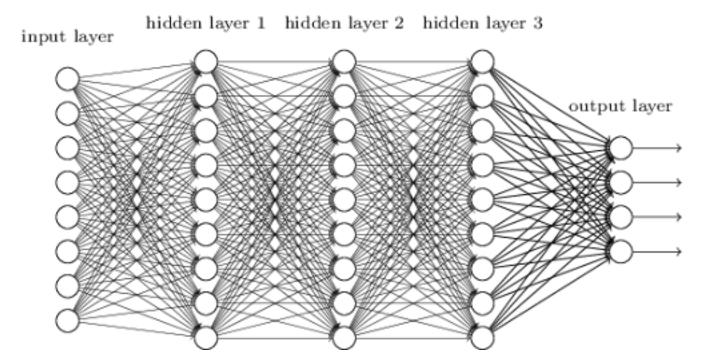
- How deep (# hidden layers) should your network be?
- How wide (# neurons in a layer) should your network be?



Hidden Layers

- How deep (# hidden layers) should your network be?
- How wide (# neurons in a layer) should your network be?

How complex is the problem you are trying to solve?



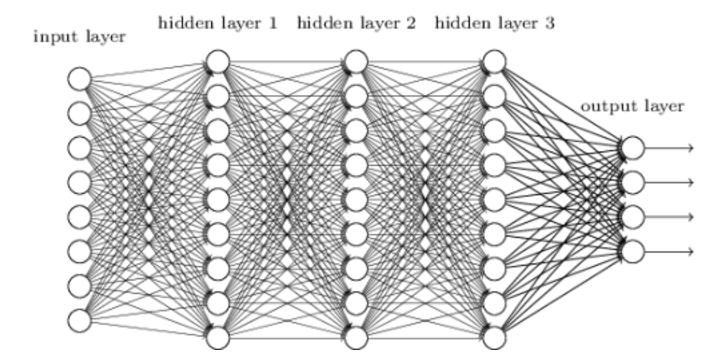
Hidden Layers

- How deep (# hidden layers) should your network be?
- How wide (# neurons in a layer) should your network be?

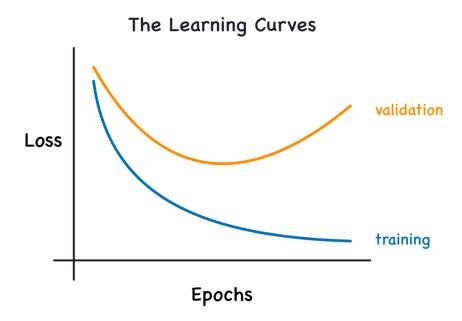
How complex is the problem you are trying to solve?

Process of (informed) trial and error.

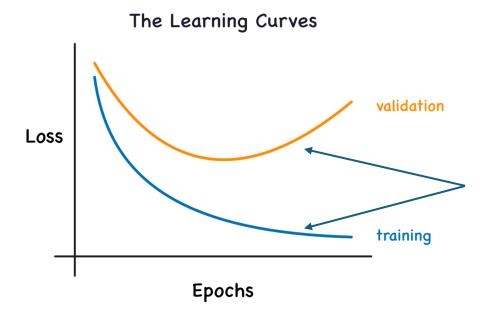
How do you know if one hyperparameter setting is better than
another?



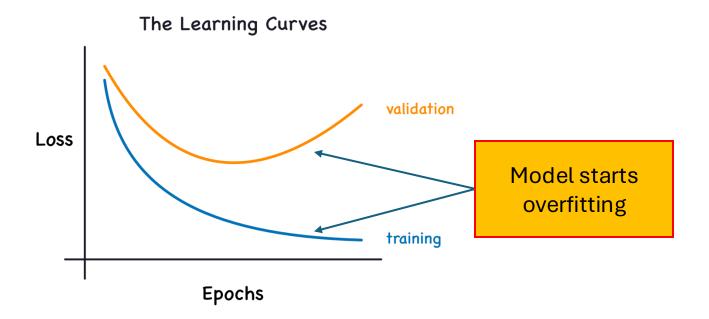
(In theory)



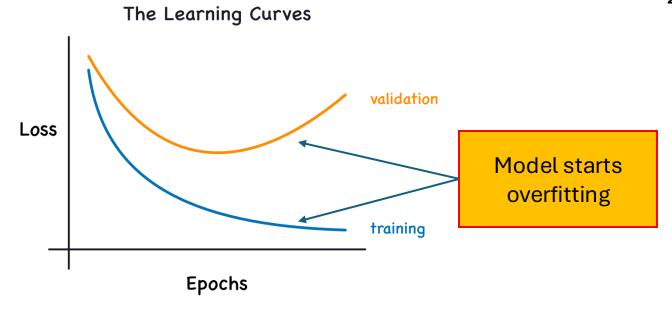
(In theory)



(In theory)

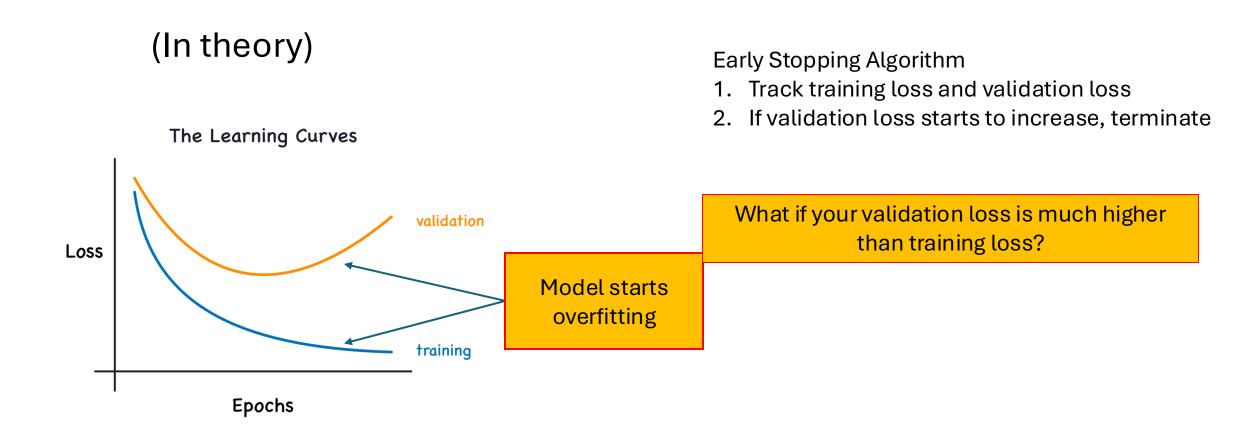


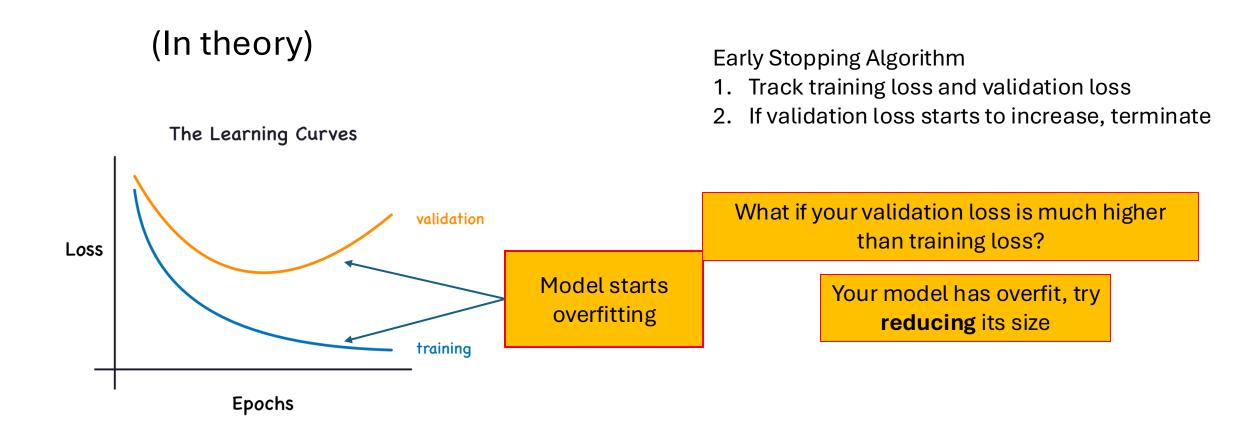
(In theory)

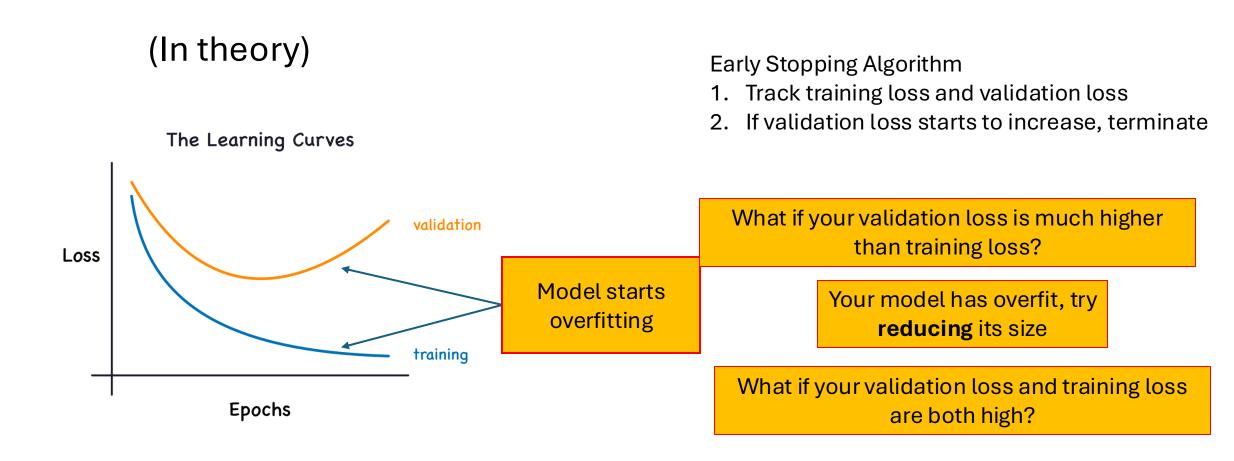


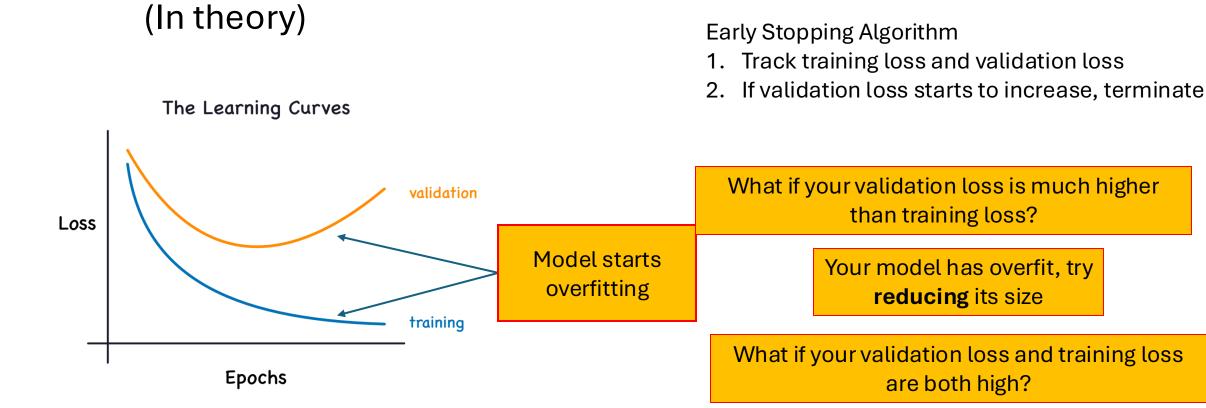
Early Stopping Algorithm

- 1. Track training loss and validation loss
- 2. If validation loss starts to increase, terminate









Your model has underfit,

try increasing its size

Is adding more width or depth better?

Is adding more width or depth better?



CSCI 1470

Deep Learning

Section S01, CRN 26629 Spring 2025

Theoretical Approaches to Understanding Depth

Proofs:

- Are there functions that deep networks can represent better than shallow networks (with similar numbers of neurons)?

Conceptual Understanding:

- Neural Networks and Manifolds for representation learning

Benefits of depth in neural networks

"For any positive integer k, there exist neural networks with $\Theta(k^3)$ layers, $\Theta(1)$ nodes per layer, and $\Theta(1)$ distinct parameters which can not be approximated by networks with O(k) layers unless they are exponentially large — they must possess $\Omega(2^k)$ nodes."

Benefits of depth in neural networks

"For any positive integer k, there exist neural networks with $\Theta(k^3)$ layers, $\Theta(1)$ nodes per layer, and $\Theta(1)$ distinct parameters which can not be approximated by networks with O(k) layers unless they are exponentially large — they must possess $\Omega(2^k)$ nodes."

There exist functions that shallow networks cannot represent as efficiently as deep networks

Matus Telgarsky "Benefits of depth in neural networks", JMLR

Benefits of depth in neural networks

"For any positive integer k, there exist neural networks with $\Theta(k^3)$ layers, $\Theta(1)$ nodes per layer, and $\Theta(1)$ distinct parameters which can not be approximated by networks with O(k) layers unless they are exponentially large — they must possess $\Omega(2^k)$ nodes."

There exist functions that shallow networks cannot represent as efficiently as deep networks

How well does theory match real world applications? Are these functions pathological?

Matus Telgarsky "Benefits of depth in neural networks", JMLR

Depth-Width Tradeoffs in Approximating Natural Functions with Neural Networks

Itay Safran Weizmann Institute of Science

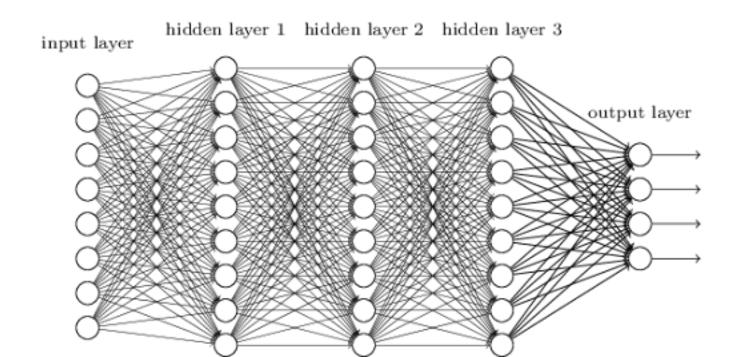
Ohad Shamir Weizmann Institute of Science itay.safran@weizmann.ac.il ohad.shamir@weizmann.ac.il

With the same number of total parameters, deep networks can learn more complex functions.

Recall that NNs are compositions of functions for which we are learning parameters:

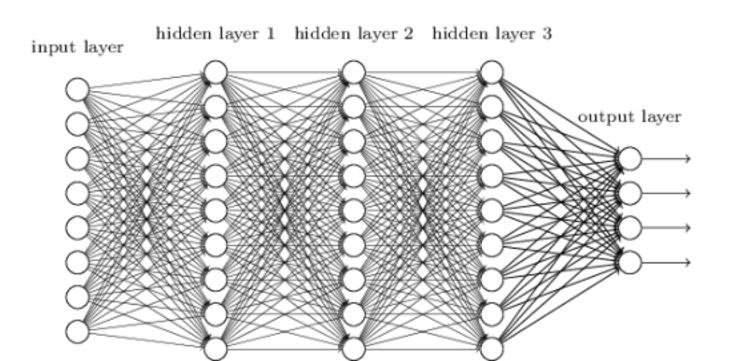
It's better (in general) to have more functions composed than it is to have more complex functions

• If there are 10 inputs, 3 layers of 10 neurons, and 4 outputs, how many weights are there total?

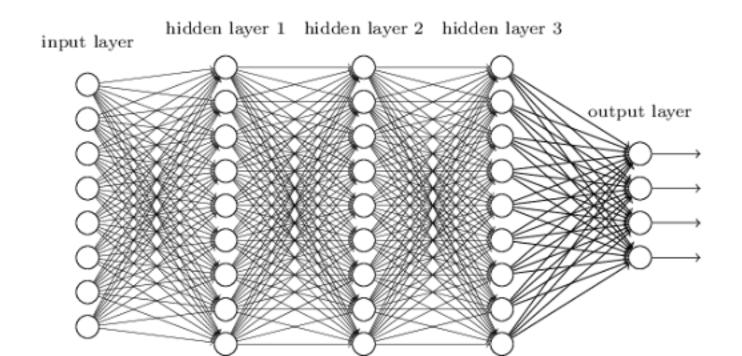


• If there are 10 inputs, 3 layers of 10 neurons, and 4 outputs, how many weights are there total? $W_1 \in \mathbb{R}^1$

 $W_1 \in \mathbb{R}^{10 \times 10}$ $W_2 \in \mathbb{R}^{10 \times 10}$ $W_3 \in \mathbb{R}^{10 \times 10}$ $W_4 \in \mathbb{R}^{10 \times 4}$ Total = 340

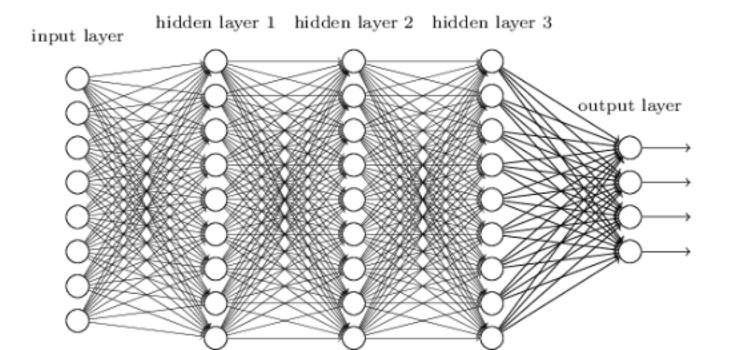


- If there are 10 inputs, 3 layers of 10 neurons, and 4 outputs, how many weights are there total?
- What if we double the width of each hidden layer?



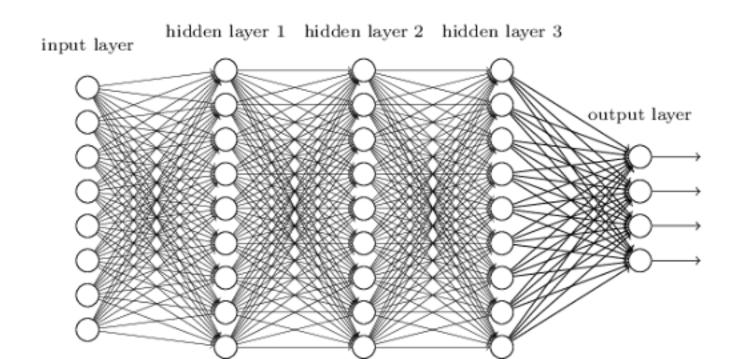
• If there are 10 inputs, 3 layers of 10 neurons, and 4 outputs, how many weights are there total? $W_1 \in \mathbb{R}^1$

• What if we double the width of each hidden layer?



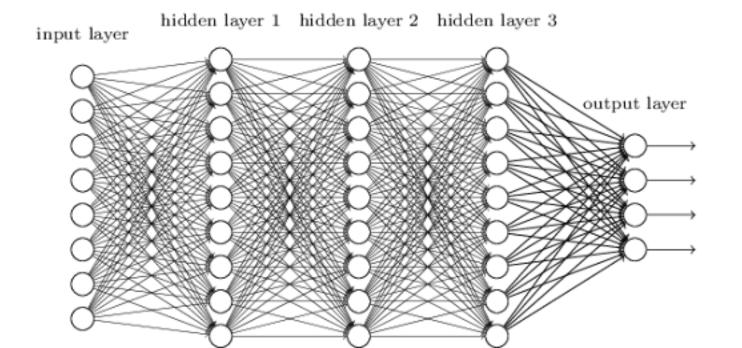
Total = 1080

- If there are 10 inputs, 3 layers of 10 neurons, and 4 outputs, how many weights are there total?
- What if we double the depth?



• If there are 10 inputs, 3 layers of 10 neurons, and 4 outputs, how many weights are there total?

What if we double the depth?



 $W_{1} \in \mathbb{R}^{10 \times 10}$ $W_{2} \in \mathbb{R}^{10 \times 10}$ $W_{3} \in \mathbb{R}^{10 \times 10}$ $W_{4} \in \mathbb{R}^{10 \times 10}$ $W_{5} \in \mathbb{R}^{10 \times 10}$ $W_{6} \in \mathbb{R}^{10 \times 10}$ $W_{7} \in \mathbb{R}^{10 \times 4}$ Total = 640

Manifold: A space that appears locally like Euclidean space

Manifold: A space that appears locally like Euclidean space

Locally, the surface of the earth appears like a flat plane in \mathbb{R}^2 , while the earth itself is a sphere(-ish) in \mathbb{R}^3



Manifold: A space that appears locally like Euclidean space

Locally, the surface of the earth appears like a flat plane in \mathbb{R}^2 , while the earth itself is a sphere(-ish) in \mathbb{R}^3



Hypothesis: real-world high-dimensional data lies on low-dimensional manifolds embedded within the high-dimensional space.

Manifold: A space that appears locally like Euclidean space

Locally, the surface of the earth appears like a flat plane in \mathbb{R}^2 , while the earth itself is a sphere(-ish) in \mathbb{R}^3



Hypothesis: real-world high-dimensional data lies on low-dimensional manifolds embedded within the high-dimensional space.

Even though we may have d features in your data, it may require many fewer features to fully represent.

Hypothesis: real-world high-dimensional data lies on low-dimensional manifolds embedded within the high-dimensional space.

Hypothesis: real-world high-dimensional data lies on low-dimensional manifolds embedded within the high-dimensional space.

MNIST images are 28x28 or 784 pixels total.

Hypothesis: real-world high-dimensional data lies on low-dimensional manifolds embedded within the high-dimensional space.

MNIST images are 28x28 or 784 pixels total.

If we restrict our pixels to only being black or white (0 or 1), then there are 2^{784} possible images we can create.

Hypothesis: real-world high-dimensional data lies on low-dimensional manifolds embedded within the high-dimensional space.

MNIST images are 28x28 or 784 pixels total.

If we restrict our pixels to only being black or white (0 or 1), then there are 2^{784} possible images we can create.

How many of these images are digits?

Hypothesis: real-world high-dimensional data lies on low-dimensional manifolds embedded within the high-dimensional space.

MNIST images are 28x28 or 784 pixels total.

If we restrict our pixels to only being black or white (0 or 1), then there are 2^{784} possible images we can create.

How many of these images are digits?

Our high-dimensional data is very sparse in high dimensions, perhaps there is a lower dimensional space where it can be better represented.

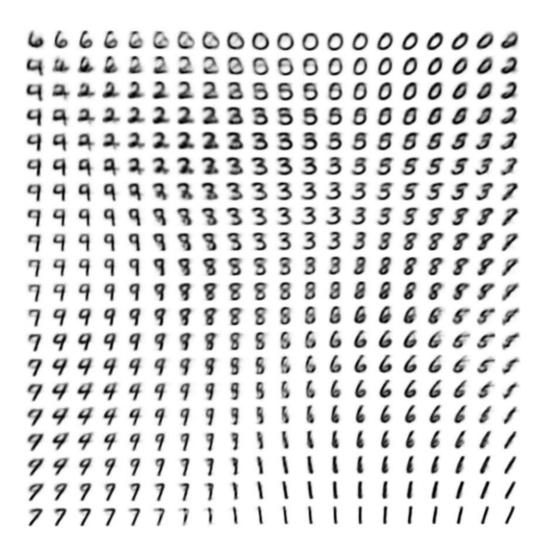
Hypothesis: real-world high-dimensional data lies on low-dimensional manifolds embedded within the high-dimensional space.

MNIST images are 28x28 or 784 pixels total.

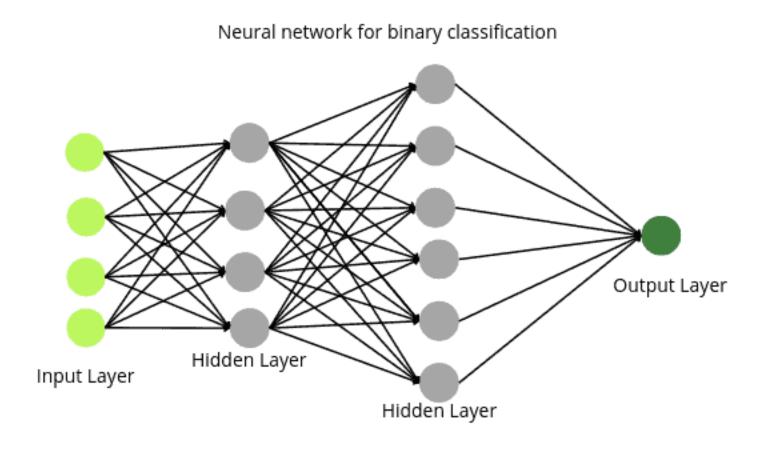
If we restrict our pixels to only being black or white (0 or 1), then there are 2^{784} possible images we can create.

How many of these images are digits?

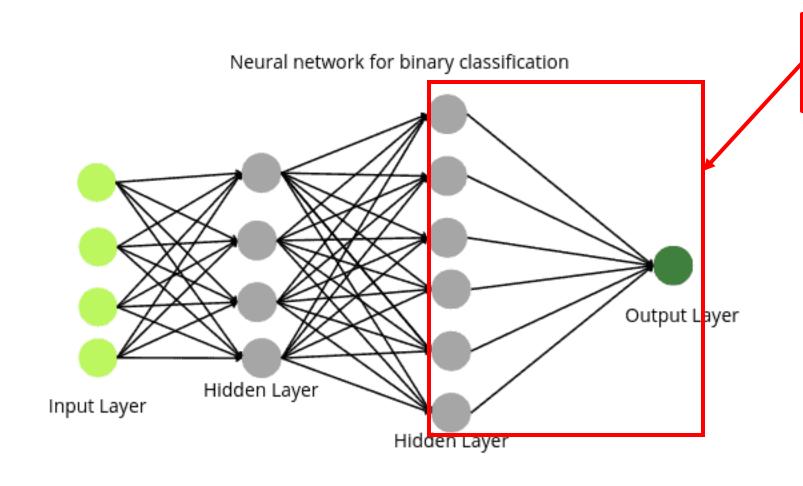
Our high-dimensional data is very sparse in high dimensions, perhaps there is a lower dimensional space where it can be better represented.



Deep Networks and Representation Learning

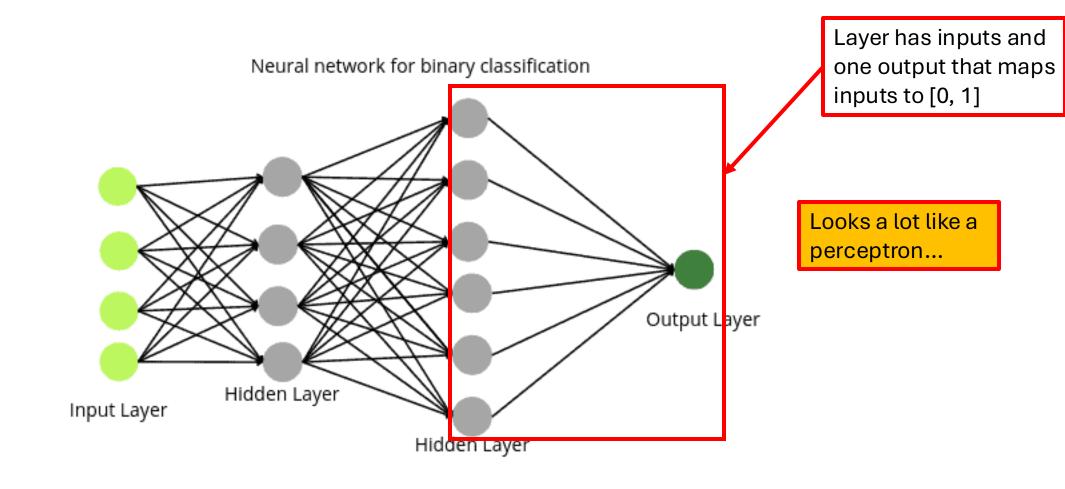


Deep Networks and Representation Learning

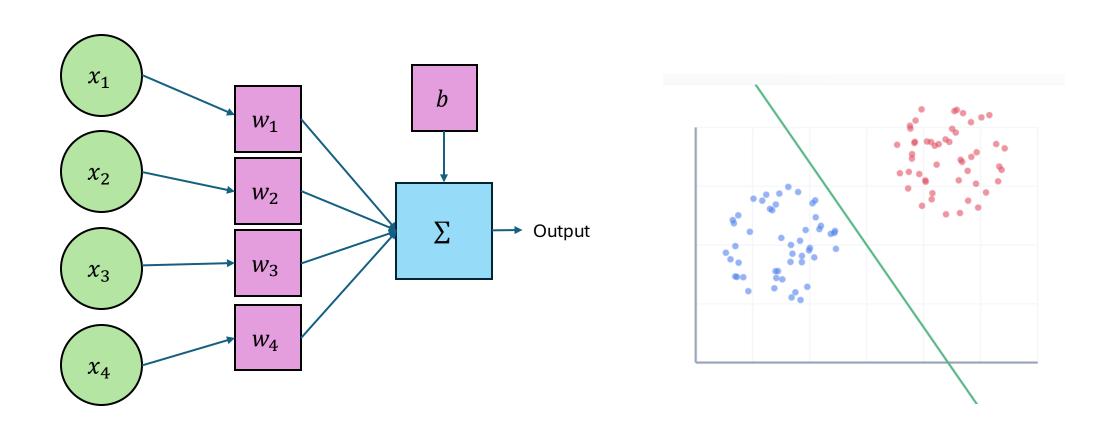


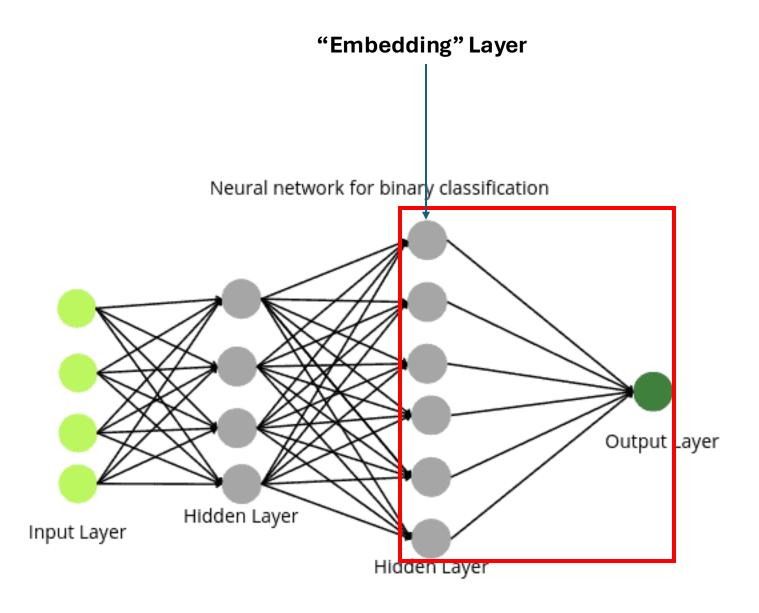
Layer has inputs and one output that maps inputs to [0, 1]

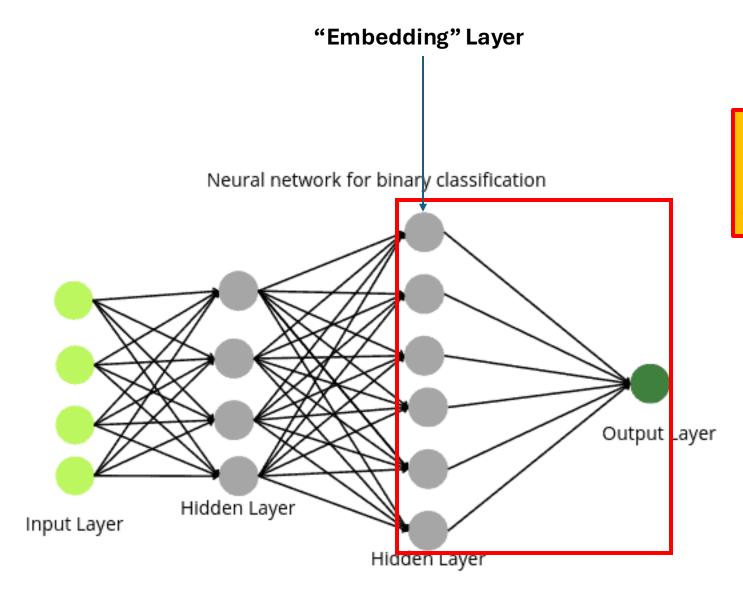
Deep Networks and Representation Learning



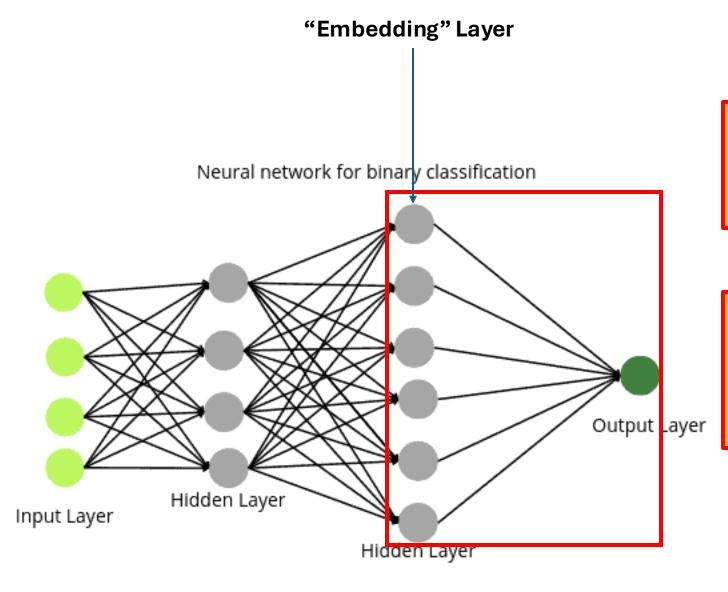
Perceptrons are Linear Separators







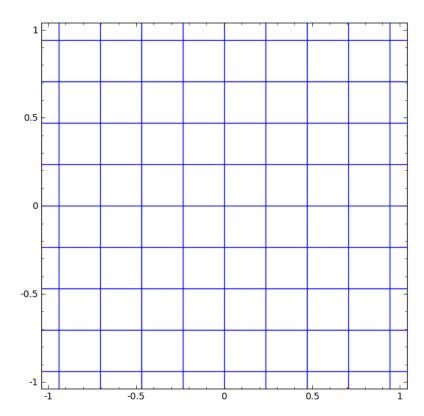
If the network can achieve 100% accuracy and the final layer is a linear separator (ala a perceptron), what does that imply about the embedding layer?



If the network can achieve 100% accuracy and the final layer is a linear separator (ala a perceptron), what does that imply about the embedding layer?

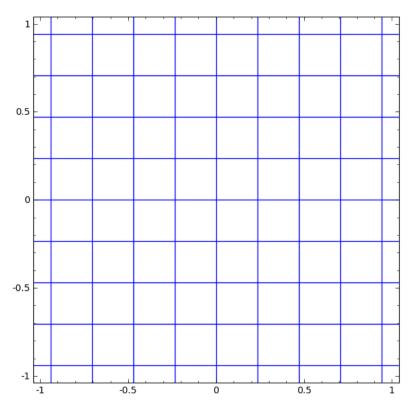
Neural Networks are learning to transform data into new learned "features" in the embedding layer. In the case of classification, the NN tries to learn linearly separable features.

A Linear Transformation applied to (x, y) coordinates

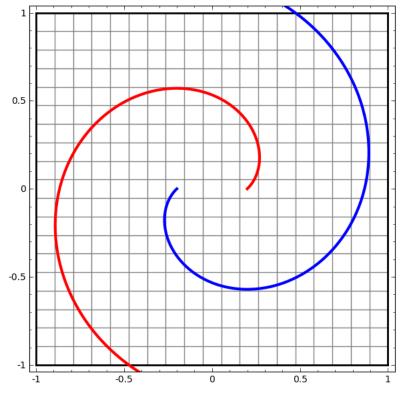


https://colah.github.io/posts/2014-03-NN-Manifolds-Topology/

A Linear Transformation applied to (x, y) coordinates



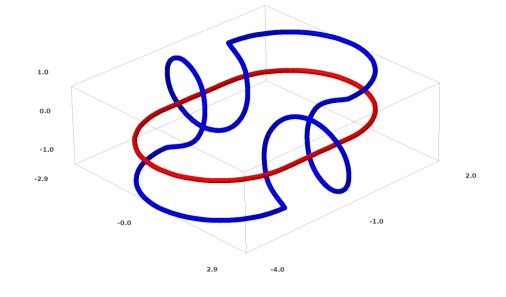
A series of linear transformations (4) applied to (x, y) coordinates to separate a spiral

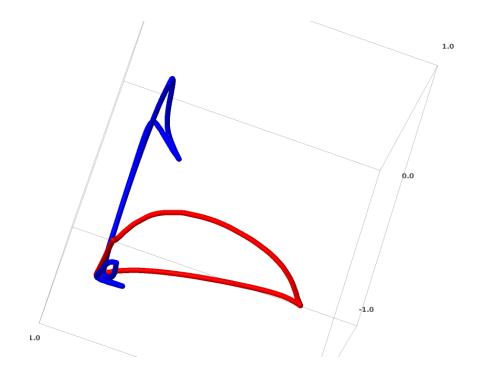


https://colah.github.io/posts/2014-03-NN-Manifolds-Topology/

Data may be hard to classify in its original form, but a series of transformations can transform it to a representation where classification is easy.

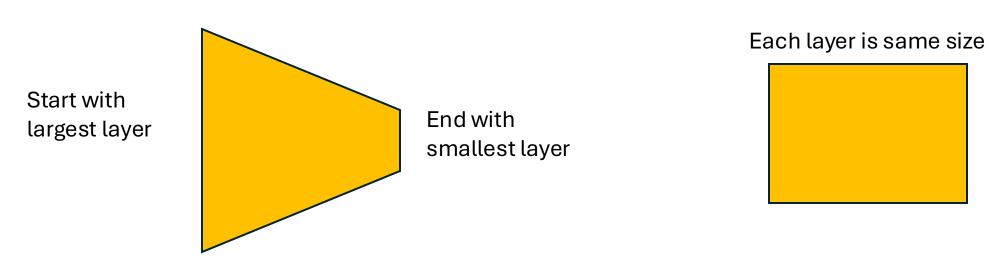
Neural Networks may be knot "untanglers"

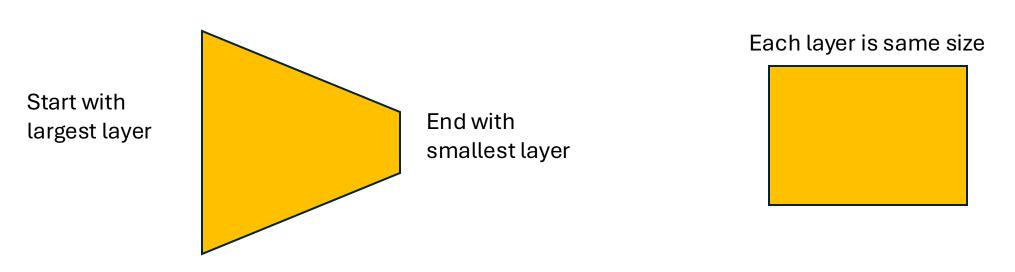


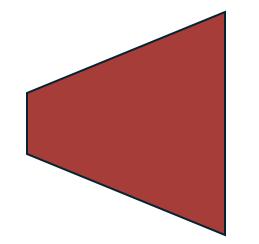


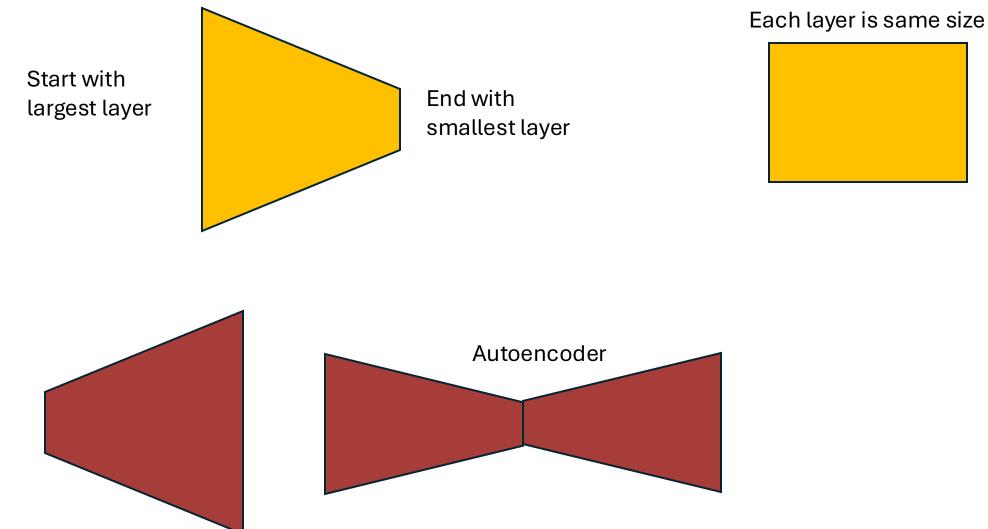
Each layer is same size

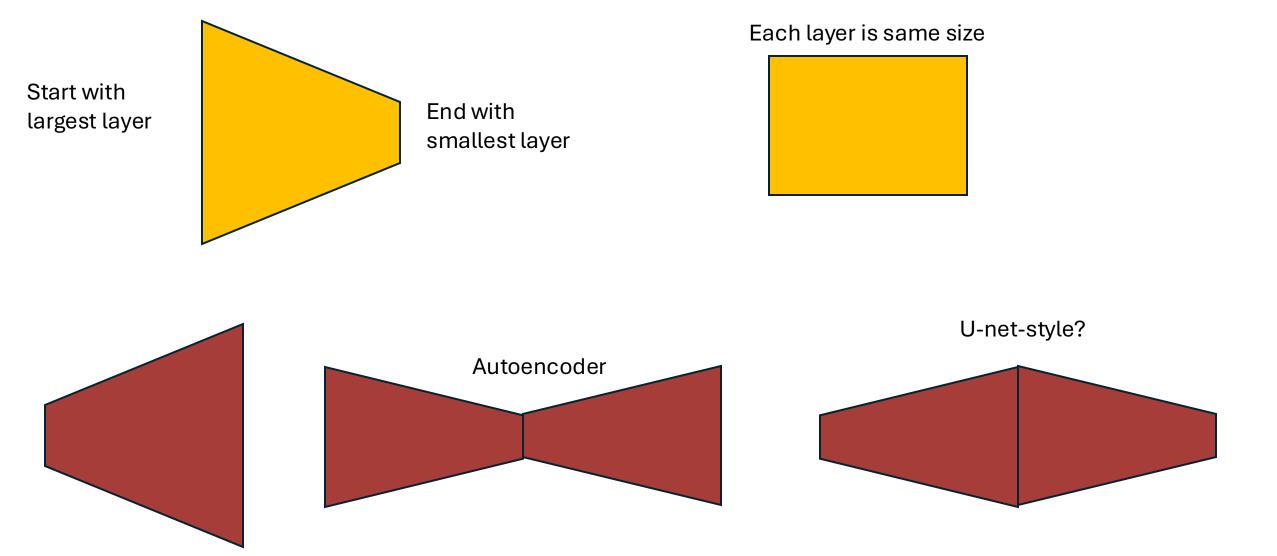












So How Many Layers/How Large Should the be?

- Final embedding needs to be expressive enough to represent your data in meaningful learned features
- Layer(s) before your embedding layer should be complex enough to transform your data into the embedding features.
- You are unlikely to need more than three sequential hidden linear layers for most common tasks

Overparameterization

Overparametization: Using more parameters than necessary for a ML problem.

Playing Atari with Deep Reinforcement Learning

Volodymyr Mnih Koray Kavukcuoglu David Silver Alex Graves Ioannis Antonoglou

Daan Wierstra Martin Riedmiller

DeepMind Technologies

{vlad, koray, david, alex.graves, ioannis, daan, martin.riedmiller} @ deepmind.com

~10,000 parameters in network

Overparameterization

Overparametization: Using more parameters than necessary for a ML problem.

Most of the time, networks use many more parameters than *necessary*.

In general, it's impossible to know the fewest amount of parameters that could solve a problem.

Playing Atari with Deep Reinforcement Learning

Volodymyr Mnih Koray Kavukcuoglu David Silver Alex Graves Ioannis Antonoglou

Daan Wierstra Martin Riedmiller

DeepMind Technologies

{vlad, koray, david, alex.graves, ioannis, daan, martin.riedmiller} @ deepmind.com

~10,000 parameters in network

PLAYING ATARI WITH SIX NEURONS

Giuseppe Cuccu

eXascale Infolab
Department of Computer Science
University of Fribourg, Switzerland
name.surname@unifr.ch

Julian Togelius

Game Innovation Lab
Tandon School of Engineering
New York University, NY, USA
julian@togelius.com

Philippe Cudré-Mauroux

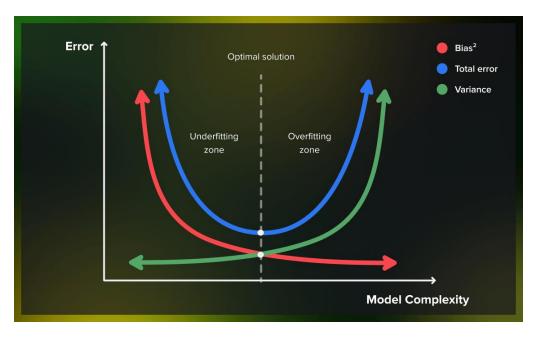
eXascale Infolab
Department of Computer Science
University of Fribourg, Switzerland
name.surname@unifr.ch

(This paper doesn't use SGD or backprop, but another optimization method)

ABSTRACT

Deep reinforcement learning, applied to vision-based problems like Atari games, maps pixels directly to actions; internally, the deep neural network bears the responsibility of both extracting useful information and making decisions based on it. By separating the image processing from decision-making, one could better understand the complexity of each task, as well as potentially find smaller policy representations that are easier for humans to understand and may generalize better. To this end, we propose a new method for learning policies and compact state representations separately but simultaneously for policy approximation in reinforcement learning. State representations are generated by an encoder based on two novel algorithms: Increasing Dictionary Vector Quantization makes the encoder capable of growing its dictionary size over time, to address new observations as they appear in an open-ended online-learning context; Direct Residuals Sparse Coding encodes observations by disregarding reconstruction error minimization, and aiming instead for highest information inclusion. The encoder autonomously selects observations online to train on, in order to maximize code sparsity. As the dictionary size increases, the encoder produces increasingly larger inputs for the neural network: this is addressed by a variation of the Exponential Natural Evolution Strategies algorithm which adapts its probability distribution dimensionality along the run. We test our system on a selection of Atari games using tiny neural networks of only 6 to 18 neurons (depending on the game's controls). These are still capable of achieving results comparable—and occasionally superior—to state-of-the-art techniques which use two orders of magnitude more neurons.

Bias-Variance Tradeoff (Traditional Understanding)



A Farewell to the Bias-Variance Tradeoff? An Overview of the Theory of Overparameterized Machine Learning

Yehuda Dar* Vidya Muthukumar[†] Richard G. Baraniuk[‡]

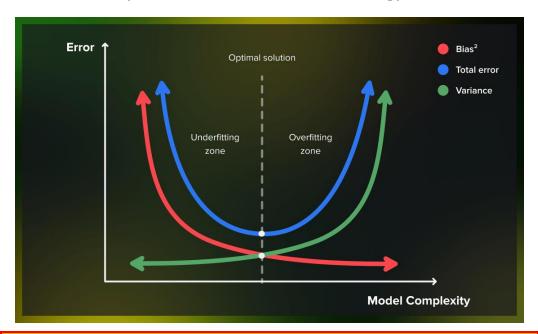
Abstract

The rapid recent progress in machine learning (ML) has raised a number of scientific questions that challenge the longstanding dogma of the field. One of the most important riddles is the good empirical generalization of *overparameterized* models. Overparameterized models are excessively complex with respect to the size of the training dataset, which results in them perfectly fitting (i.e., *interpolating*) the training data, which is usually noisy. Such interpolation of noisy data is traditionally associated with detrimental overfitting, and yet a wide range of interpolating models – from simple linear models to deep neural networks – have recently been observed to generalize extremely well on fresh test data. Indeed, the recently discovered *double descent* phenomenon has revealed that highly overparameterized models often improve over the best underparameterized model in test performance.

Understanding learning in this overparameterized regime requires new theory and foundational empirical studies, even for the simplest case of the linear model. The underpinnings of this understanding have been laid in very recent analyses of overparameterized linear regression and related statistical learning tasks, which resulted in precise analytic characterizations of double descent. This paper provides a succinct overview of this emerging theory of overparameterized ML (henceforth abbreviated as TOPML) that explains these recent findings through a statistical signal processing perspective. We emphasize the unique aspects that define the TOPML research area as a subfield of modern ML theory and outline interesting open questions that remain.

https://serokell.io/blog/bias-variance-tradeoff

Bias-Variance Tradeoff (Traditional Understanding)



If you are overfitting, reduce model complexity (smaller width/fewer layers). If underfitting, add more model complexity.

https://serokell.io/blog/bias-variance-tradeoff

A Farewell to the Bias-Variance Tradeoff? An Overview of the Theory of Overparameterized Machine Learning

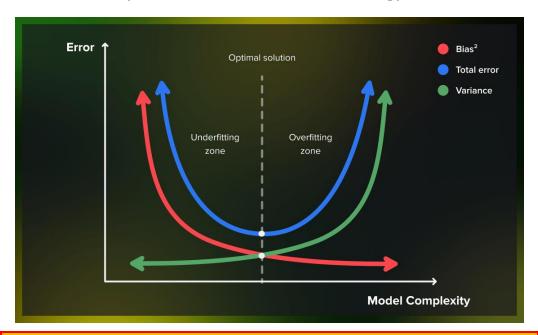
Yehuda Dar* Vidya Muthukumar[†] Richard G. Baraniuk[‡]

Abstract

The rapid recent progress in machine learning (ML) has raised a number of scientific questions that challenge the longstanding dogma of the field. One of the most important riddles is the good empirical generalization of *overparameterized* models. Overparameterized models are excessively complex with respect to the size of the training dataset, which results in them perfectly fitting (i.e., *interpolating*) the training data, which is usually noisy. Such interpolation of noisy data is traditionally associated with detrimental overfitting, and yet a wide range of interpolating models – from simple linear models to deep neural networks – have recently been observed to generalize extremely well on fresh test data. Indeed, the recently discovered *double descent* phenomenon has revealed that highly overparameterized models often improve over the best underparameterized model in test performance.

Understanding learning in this overparameterized regime requires new theory and foundational empirical studies, even for the simplest case of the linear model. The underpinnings of this understanding have been laid in very recent analyses of overparameterized linear regression and related statistical learning tasks, which resulted in precise analytic characterizations of double descent. This paper provides a succinct overview of this emerging theory of overparameterized ML (henceforth abbreviated as TOPML) that explains these recent findings through a statistical signal processing perspective. We emphasize the unique aspects that define the TOPML research area as a subfield of modern ML theory and outline interesting open questions that remain.

Bias-Variance Tradeoff (Traditional Understanding)



If you are overfitting, reduce model complexity (smaller width/fewer layers). If underfitting, add more model complexity.

https://serokell.io/blog/bias-variance-tradeoff

A Farewell to the Bias-Variance Tradeoff? An Overview of the Theory of Overparameterized Machine Learning

Yehuda Dar* Vidya Muthukumar[†] Richard G. Baraniuk[‡]

Abstract

The rapid recent progress in machine learning (ML) has raised a number of scientific questions that challenge the longstanding dogma of the field. One of the most important riddles is the good empirical generalization of *overparameterized* models. Overparameterized models are excessively complex with respect to the size of the training dataset, which results in them perfectly fitting (i.e., *interpolating*) the training data, which is usually noisy. Such interpolation of noisy data is traditionally associated with detrimental overfitting, and yet a wide range of interpolating models – from simple linear models to deep neural networks – have recently been observed to generalize extremely well on fresh test data. Indeed, the recently discovered *double descent* phenomenon has revealed that highly overparameterized models often improve over the best underparameterized model in test performance.

Understanding learning in this overparameterized regime requires new theory and foundational empirical studies, even for the simplest case of the linear model. The underpinnings of this understanding have been laid in very recent analyses of overparameterized linear regression and related statistical learning tasks, which resulted in precise analytic characterizations of double descent. This paper provides a succinct overview of this emerging theory of overparameterized ML (henceforth abbreviated as TOPML) that explains these recent findings through a statistical signal processing perspective. We emphasize the unique aspects that define the TOPML research area as a subfield of modern ML theory and outline interesting open questions that remain.

• SGD, SGD + Momentum, SGD + Adaptive Momentum (Adam), RMSProp,... the list is ever growing

• SGD, SGD + Momentum, SGD + Adaptive Momentum (Adam), RMSProp,... the list is ever growing

How do you choose between them?

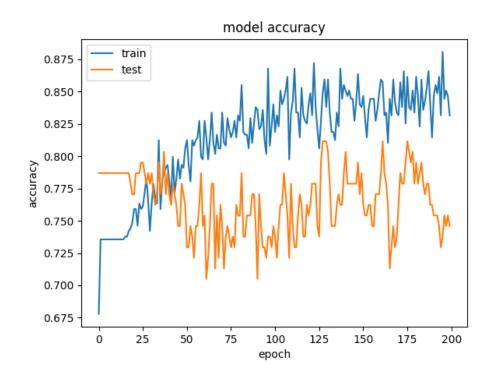
• SGD, SGD + Momentum, SGD + Adaptive Momentum (Adam), RMSProp,... the list is ever growing

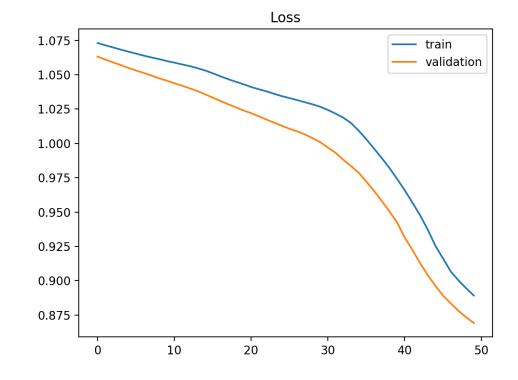
How do you choose between them?

- Just use Adam.
 - The only downside is that it might work so well that you end up overfitting.
 - Suggested initial learning rate of 3e-4

Batch Size and Learning Rate

Having too small a batch or too high a learning rate can cause variance in training/validation loss – symptoms often look similar





- Don't change too much at once.

- Don't change too much at once.
- Keep track of parameters you've tested and track their performance

- Don't change too much at once.
- Keep track of parameters you've tested and track their performance
- Don't just randomly guess parameters, apply critical thinking, come up with a hypothesis and test your hypothesis.

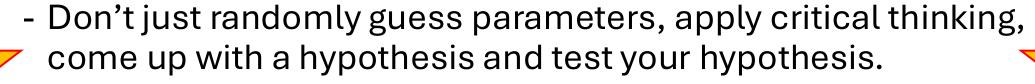
- Don't change too much at once.
- Keep track of parameters you've tested and track their performance
- Don't just randomly guess parameters, apply critical thinking, come up with a hypothesis and test your hypothesis.

(Use the scientific method)

- Don't change too much at once.
- Keep track of parameters you've tested and track their performance
- Don't just randomly guess parameters, apply critical thinking,
 come up with a hypothesis and test your hypothesis.

(Use the scientific method)

- Don't change too much at once.
- Keep track of parameters you've tested and track their performance

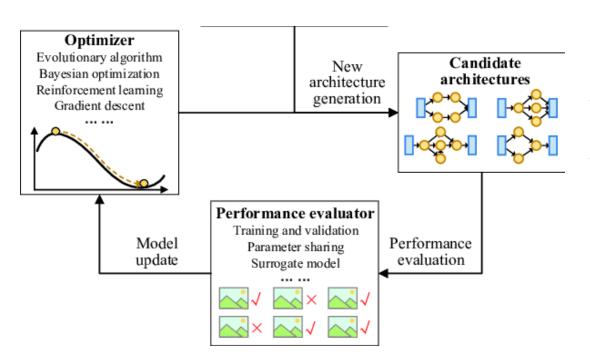


(Use the scientific method)

Andrej Karpathy: A recipe for training neural networks https://karpathy.github.io/2019/04/25/recipe/

Neural Architecture Search (NAS)

Changing hyperparameters results in different performance, can we run an optimization algorithm on our hyperparameters?



Pros:

- No longer need human input
- May find better
 hyperparameters
 than humans

Cons:

- Takes a very long time...
- Hyperparameters are discrete and highly dependent (e.g., width/depth), it's a really hard optimization problem...

Option #1: Grid search

Define a set of possible parameters (i.e., learning rates, width, depth, etc)

Try every combination of hyperparameters possible, pick setting with best validation set performance.

Option #1: Grid search

Define a set of possible parameters (i.e., learning rates, width, depth, etc)

Try every combination of hyperparameters possible, pick setting with best validation set performance.

What are some downsides of grid search?

Option #2: Bayesian Optimization

Option #2: Bayesian Optimization

We believe the performance of hyperparameters that are close together, should have similar results.

Option #2: Bayesian Optimization

We believe the performance of hyperparameters that are close together, should have similar results.

Define a set of possible parameters (i.e., learning rates, width, depth, etc)

Option #2: Bayesian Optimization

We believe the performance of hyperparameters that are close together, should have similar results.

Define a set of possible parameters (i.e., learning rates, width, depth, etc)

Try sets of possible hyperparameters, each with some probability.

Option #2: Bayesian Optimization

We believe the performance of hyperparameters that are close together, should have similar results.

Define a set of possible parameters (i.e., learning rates, width, depth, etc)

Try sets of possible hyperparameters, each with some probability.

- The probability that you try a specific hyperparameter setting depends on the performance of nearby hyperparameter settings.

Option #2: Bayesian Optimization

We believe the performance of hyperparameters that are close together, should have similar results.

Define a set of possible parameters (i.e., learning rates, width, depth, etc)

Try sets of possible hyperparameters, each with some probability.

- The probability that you try a specific hyperparameter setting depends on the performance of nearby hyperparameter settings.
- Also track uncertainty of hyperparameters (i.e., settings you have not tried something close to before)

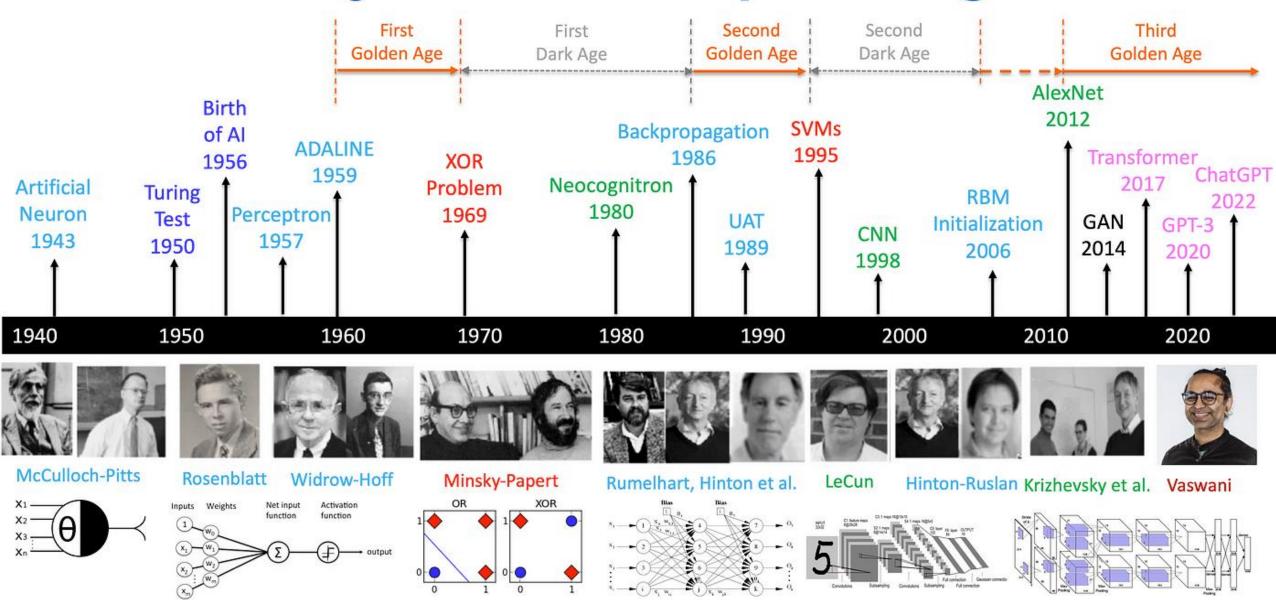
Keras tuner is compatible with Tensorflow, Pytorch, and Jax and has various automatic hyperparameter tuning methods

KerasTuner



KerasTuner is an easy-to-use, scalable hyperparameter optimization framework that solves the pain points of hyperparameter search. Easily configure your search space with a define-by-run syntax, then leverage one of the available search algorithms to find the best hyperparameter values for your models. KerasTuner comes with Bayesian Optimization, Hyperband, and Random Search algorithms built-in, and is also designed to be easy for researchers to extend in order to experiment with new search algorithms.

A Brief History of Al with Deep Learning



What has happened in the last 15 years?

What has changed?

- 1. Power and efficiency of compute (GPUs)
- 2. Availability of data (the internet)
- 3. New Architectures (e.g., CNNs, Transformers)



Issues with MLPs

- 1. Resource Intensive
- 2. Difficult to incorporate certain types of information
- 3. (and more)

Issues with MLPs

- 1. Resource Intensive
- 2. Difficult to incorporate certain types of information
- 3. (and more)

GPUs to the rescue!



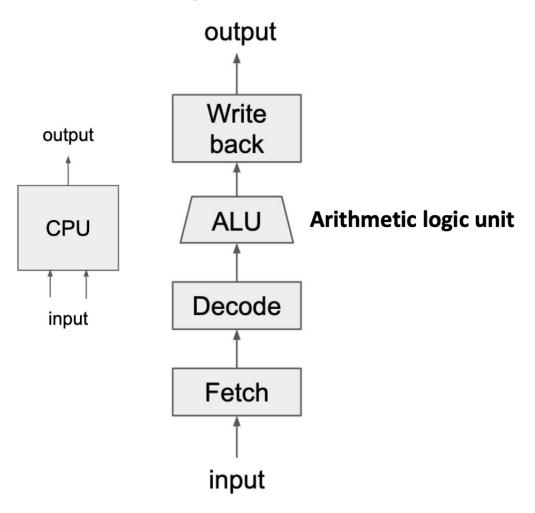


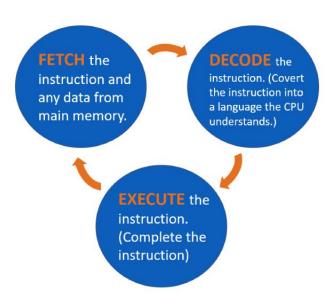
- GPUs are really good at computing mathematical operations in parallel!
- Matrix multiplication == many independent multiply and add operations

Easily parallelizable

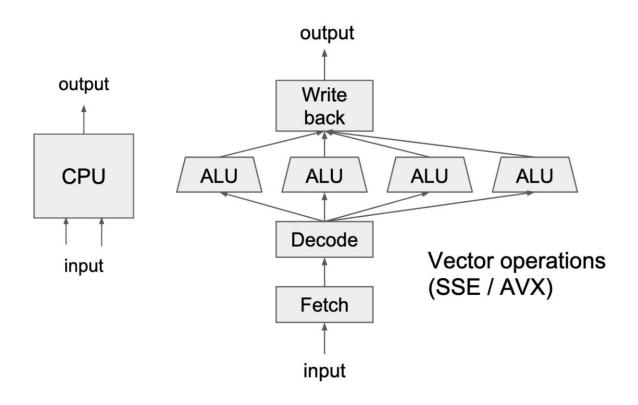
GPUs are great for this!

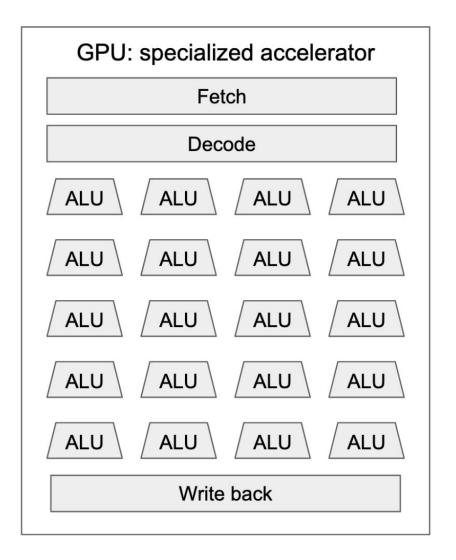
CPU v/s GPU





CPU v/s GPU





GPU-Parallel Acceleration

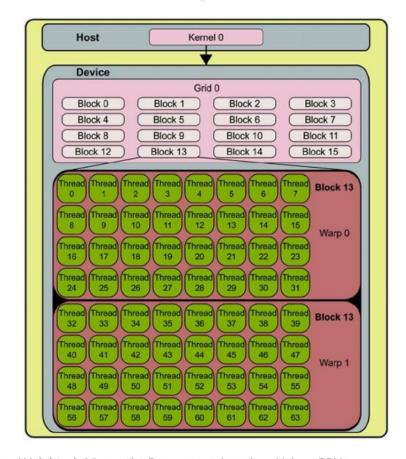
- User code (*kernels*) is compiled on the *host* (the CPU) and then transferred to the *device* (the GPU)
- Kernel is executed as a grid
- Each grid has multiple thread blocks
- Each thread block has multiple warps

A warp is the basic schedule unit in kernel execution

A warp consists of 32 threads

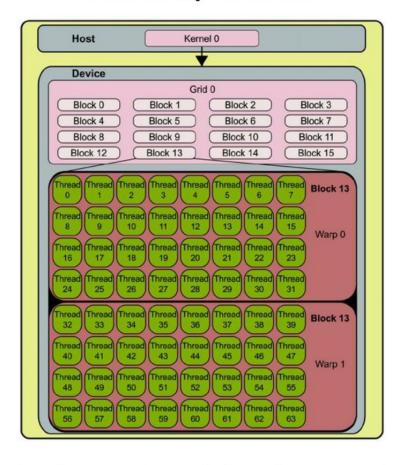
Compute Unified Device Architecture is a parallel computing platform and application programming interface (API)

CUDA compute model



GPU-Parallel Acceleration

CUDA compute model



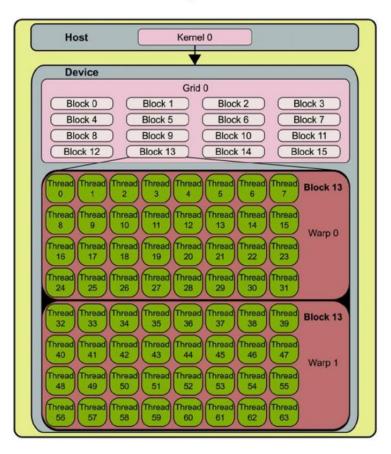
- Programmer decides how they want to parallelize the computation across grids and blocks
 - Modern deep learning frameworks take care of this for you
- CUDA compiler figures out how to schedule these units of computation on to the physical hardware

Any questions?

???

GPU-Parallel Acceleration

CUDA compute model



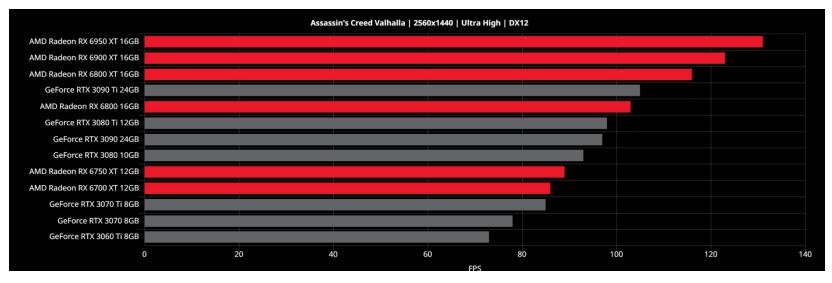
- Upshot: order of magnitude speedups!
- Example: training CNN on CIFAR-10 dataset

Device	Speed of training, examples/sec
2 x AMD Opteron 6168	440
i7-7500U	415
GeForce 940MX	1190
GeForce 1070	6500

From:

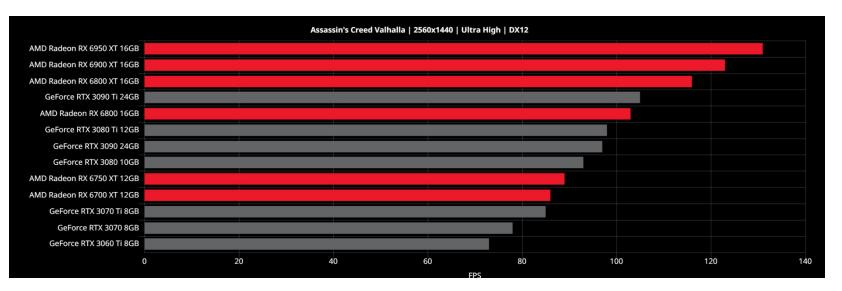
https://medium.com/@andriylazorenko/tensorflow-performance-test-cpu-vs-gpu-79fcd39170c

AMD GPUs are competitive for gaming and graphics, why not for AI?



(With a benchmarking tool made by AMD)

AMD GPUs are competitive for gaming and graphics, why not for AI?

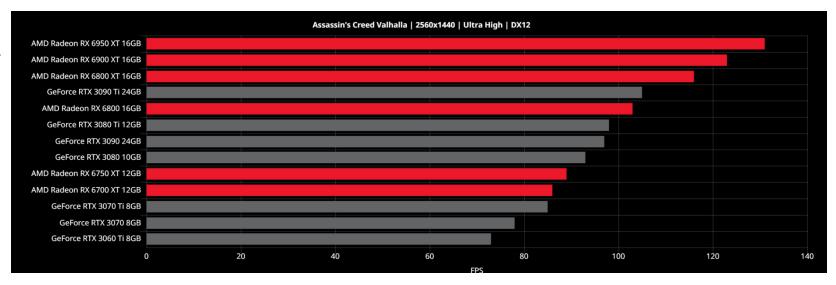


CUDA is far better than competitors (AMD)

(With a benchmarking tool made by AMD)

- Easier to use
- Better optimization
- AMD makes GPUs for graphics, NVIDIA makes GPUs for AI

AMD GPUs are competitive for gaming and graphics, why not for Al?



CUDA is far better than competitors (AMD)

(With a benchmarking tool made by AMD)

- Easier to use
- Better optimization
- AMD makes GPUs for graphics, NVIDIA makes GPUs for Al

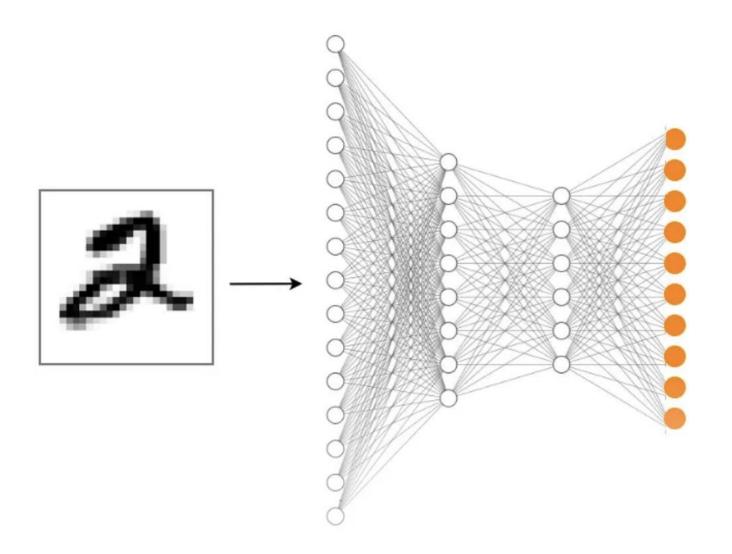
CUDA is Still a Giant Moat for NVIDIA

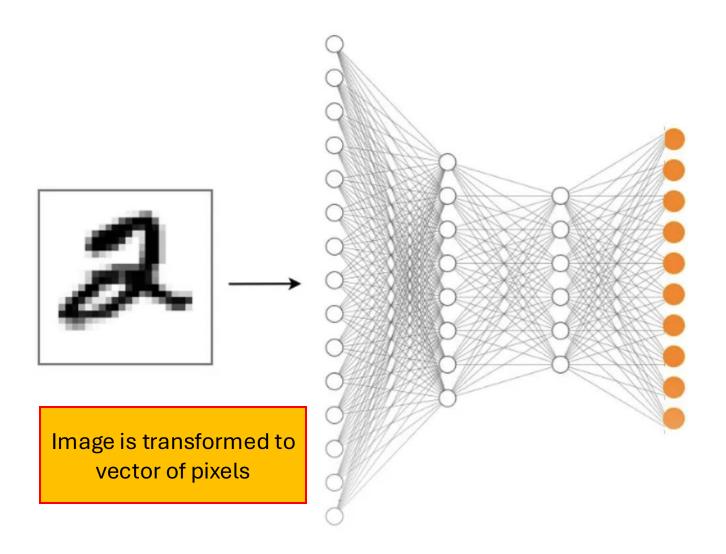
Despite everyone's focus on hardware, the software of AI is what protects NVIDIA

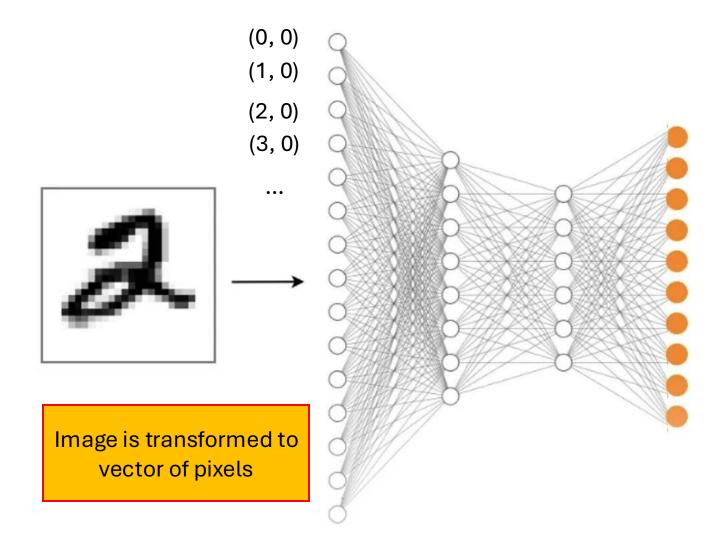


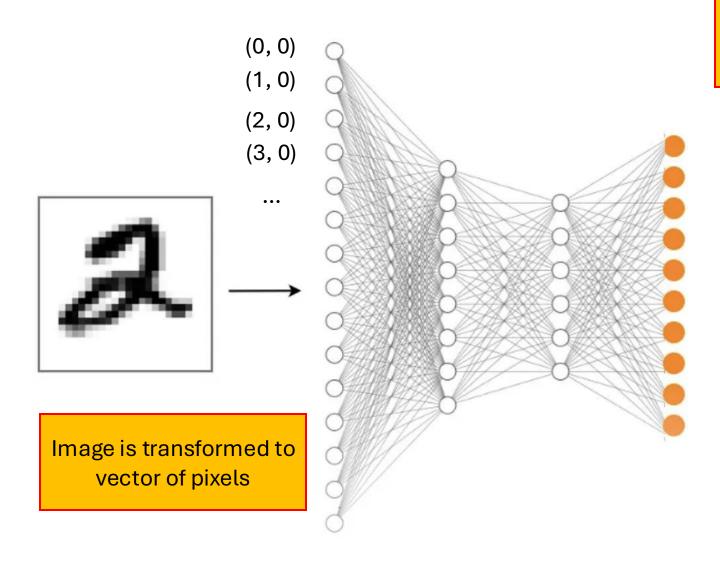
Issues with MLPs

- 1. Resource Intensive
- 2. Difficult to incorporate certain types of information

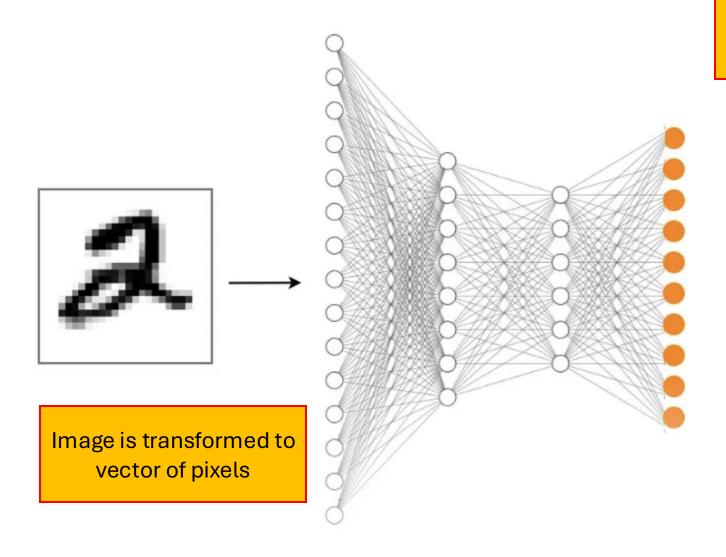




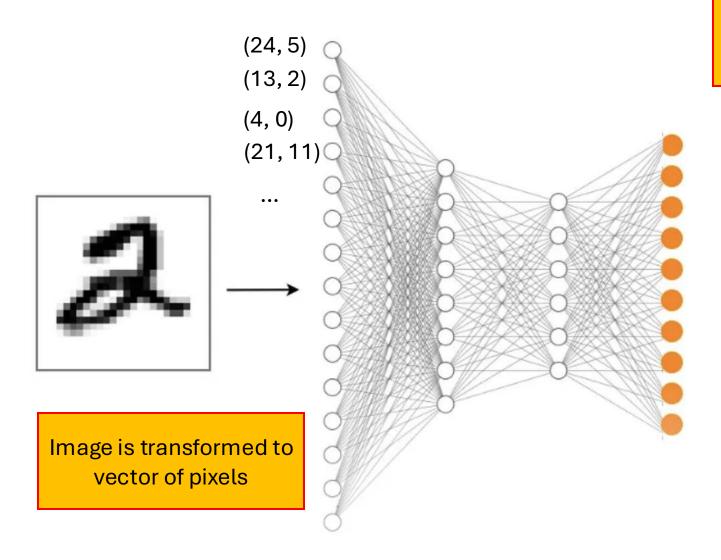




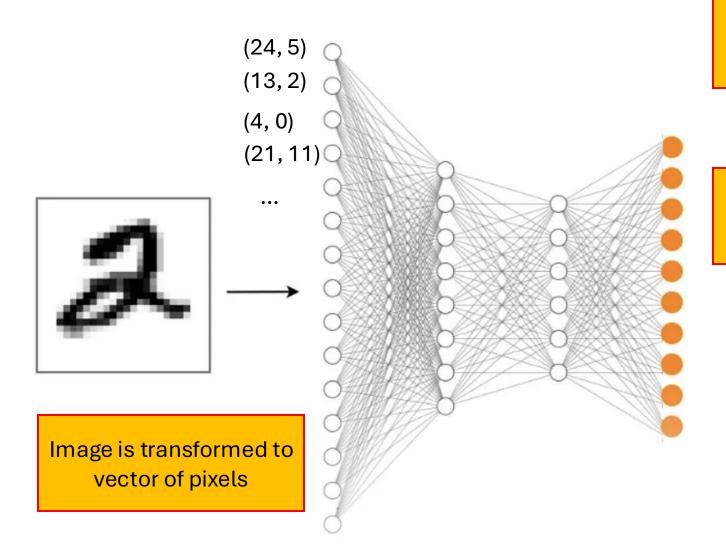
What would happen if we permuted the ordering of the pixels?



What would happen if we permuted the ordering of the pixels?

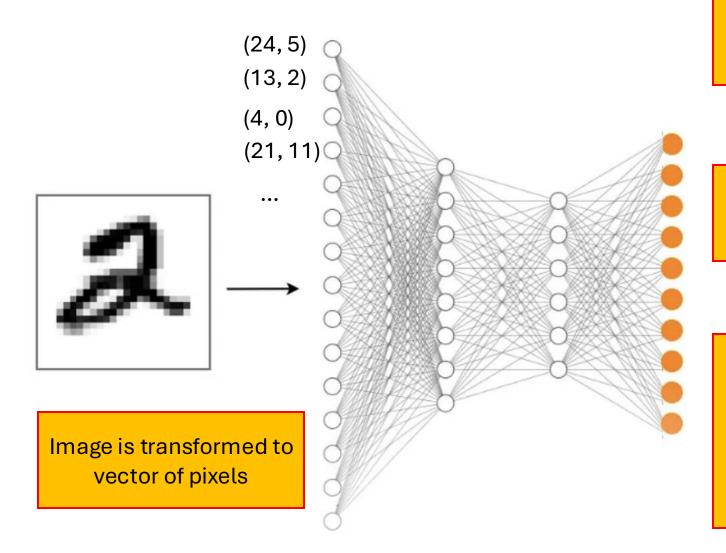


What would happen if we permuted the ordering of the pixels?



What would happen if we permuted the ordering of the pixels?

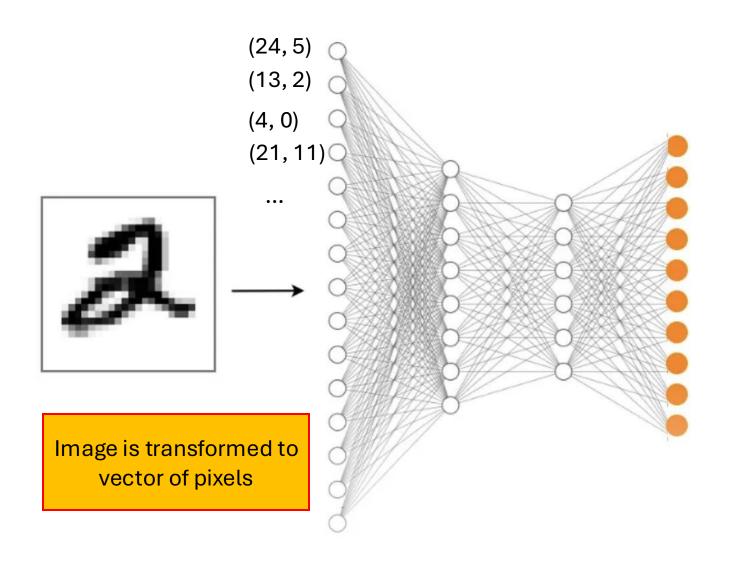
Will the training of the neural network differ?



What would happen if we permuted the ordering of the pixels?

Will the training of the neural network differ?

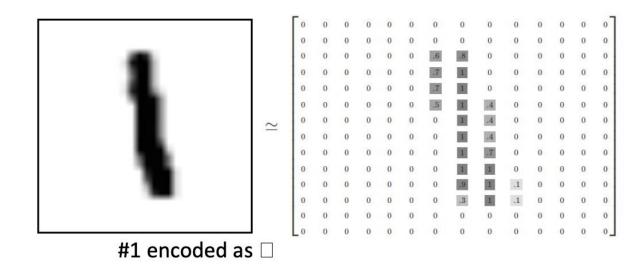
No! MLPs do not use spatial information, it does not matter which order the pixels are fed in so long as it is the same ordering for every input



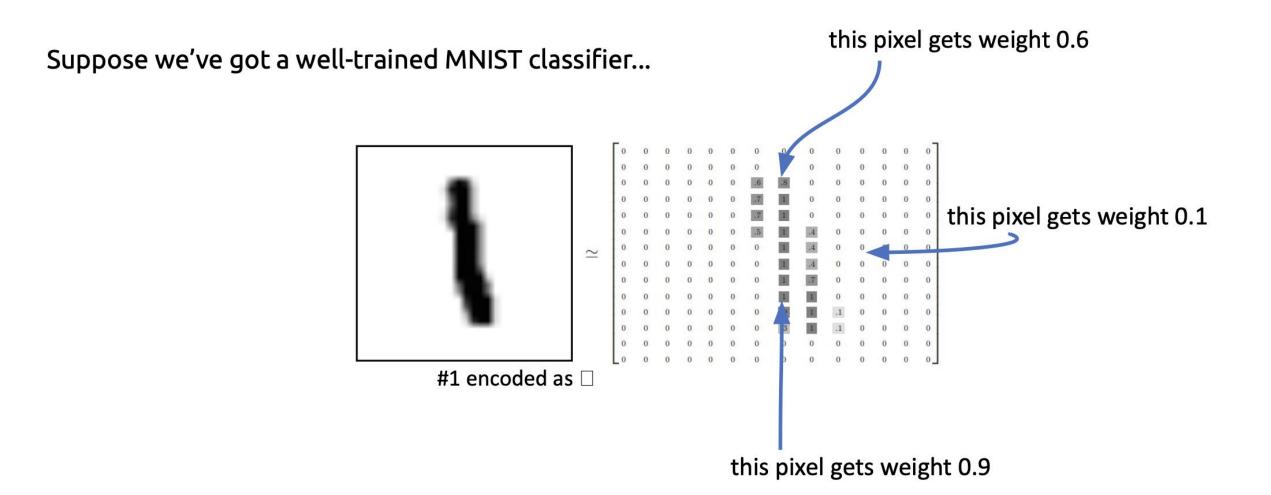
Isn't this actually a hard problem that we are trying to learn?

Limitations of Full Connections for MNIST

Suppose we've got a well-trained MNIST classifier...



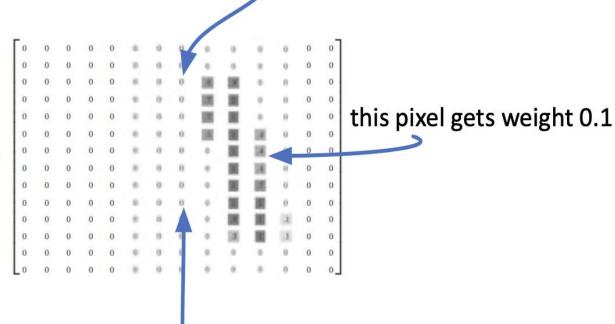
Limitations of Full Connections for MNIST



Limitations of Full Connections for MNIST

#1 encoded as \square

If we shift the digit to the right, then a different set of weights becomes relevant $\Box \rightarrow$ etwork might have trouble classifying this as a 1...

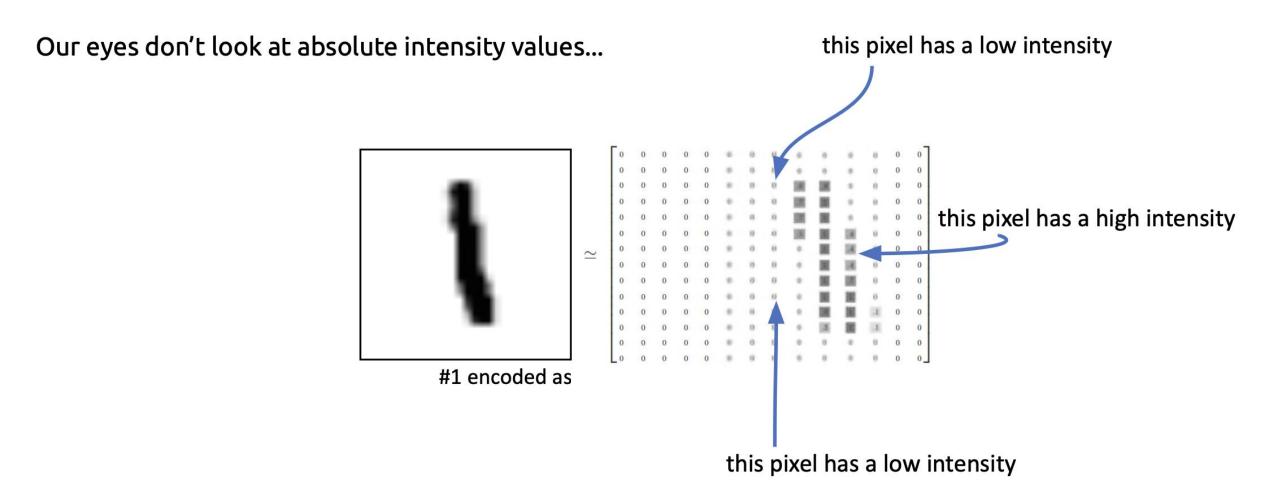


this pixel gets weight 0.6

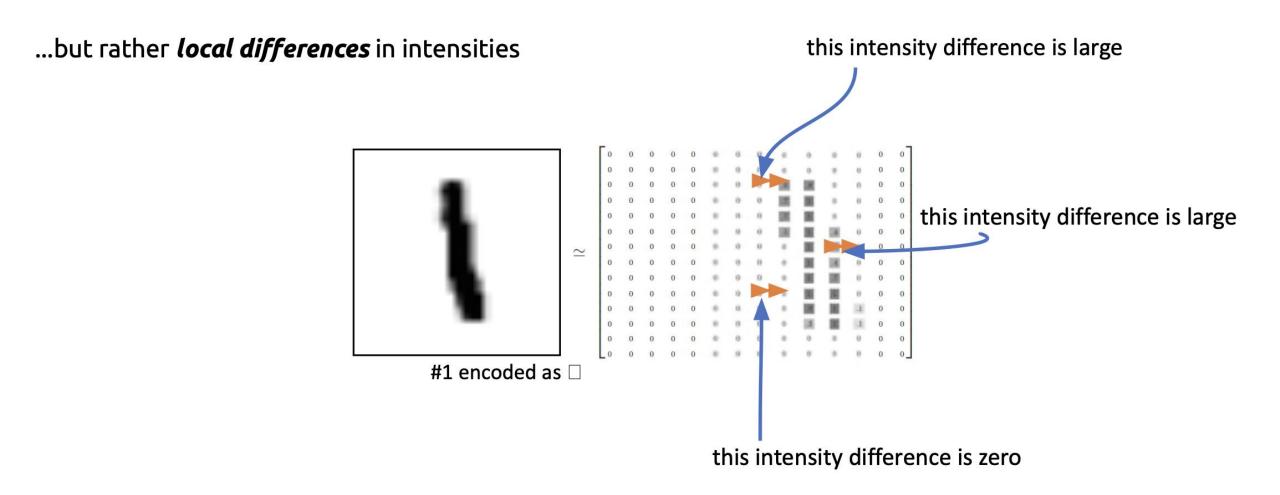
Can you tell this is a 1?

this pixel gets weight 0.9

This would **not** be a problem for the human visual system

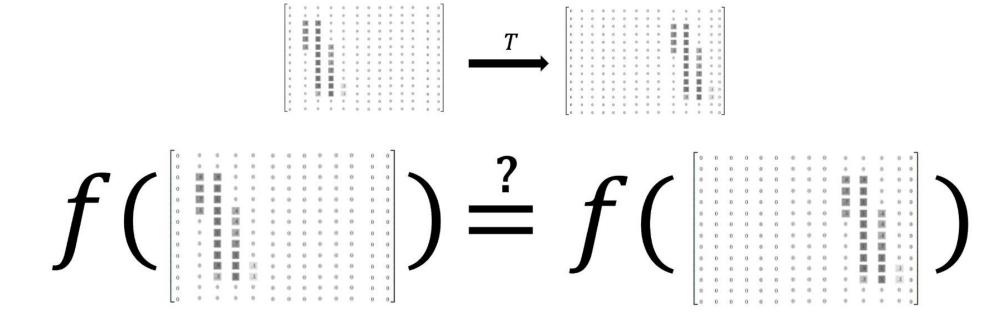


This would **not** be a problem for the human visual system

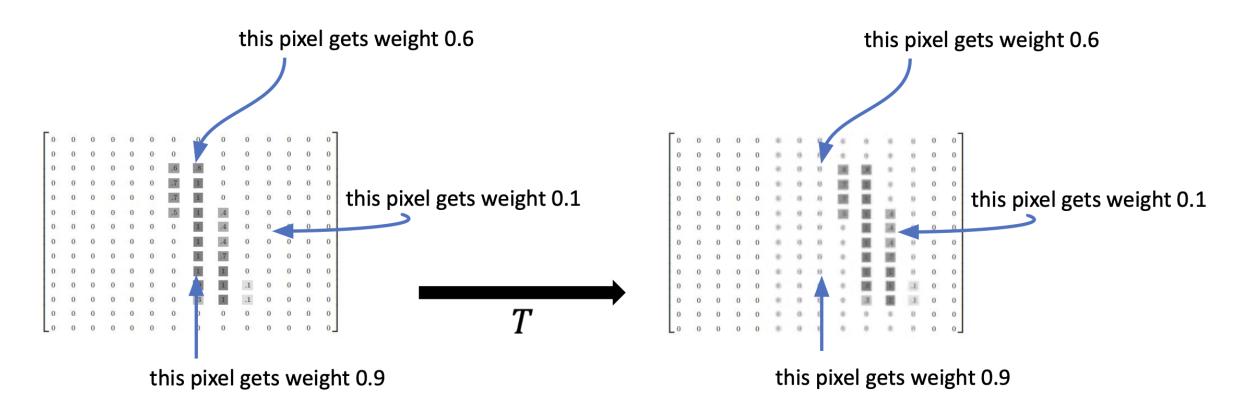


Translational Invariance

- To make a neural net f robust in this same way, it should ideally satisfy **translational invariance**: f(T(x)) = f(x), where
 - x is the input image
 - T is a translation (i.e. a horizonal and/or vertical shift)



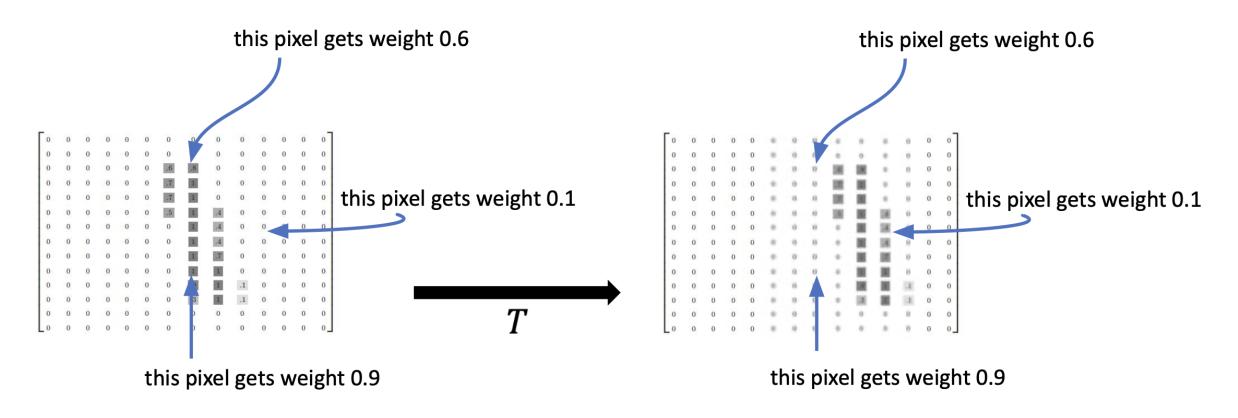
Fully Connected Nets are *not*Translationally Invariant



Sum of these three: $0.6 \cdot 0.8 + 0.1 \cdot 0 + 0.9 \cdot 1 = 1.38$

Sum of these three: $0.6 \cdot 0 + 0.1 \cdot 0.4 + 0.9 \cdot 0 = 0.4$

Fully Connected Nets are *not*Translationally Invariant



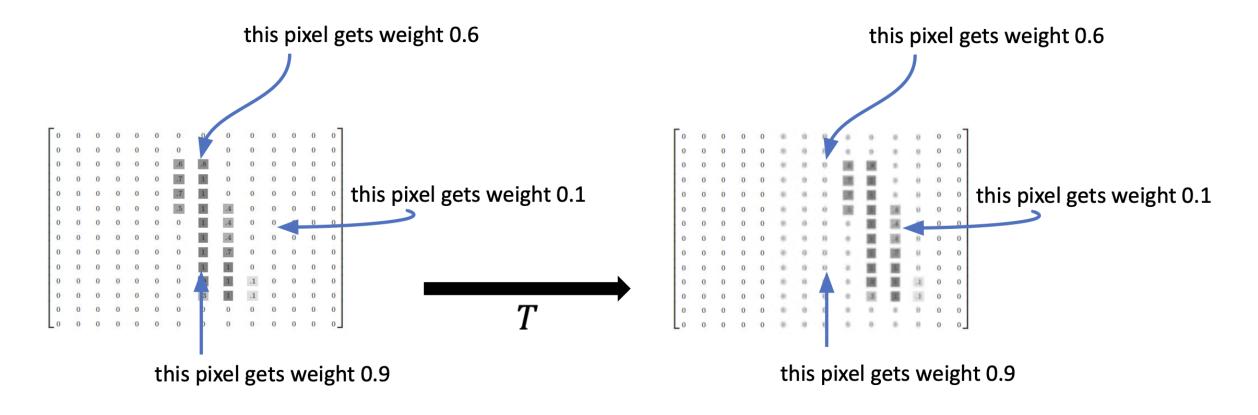
Sum of these three: $0.6 \cdot 0.8 + 0.1 \cdot 0 + 0.9 \cdot 1 = 1.38$

Sum of these three: $0.6 \cdot 0 + 0.1 \cdot 0.4 + 0.9 \cdot 0 = 0.4$

Fully Connected Nets are *not*Translationally Invariant

How to make the network translationally invariant?

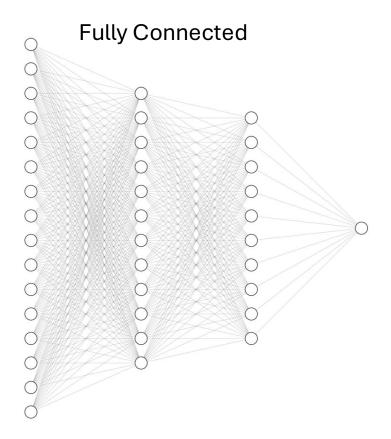
Focus on local differences/patterns



Sum of these three: $0.6 \cdot 0.8 + 0.1 \cdot 0 + 0.9 \cdot 1 = 1.38$

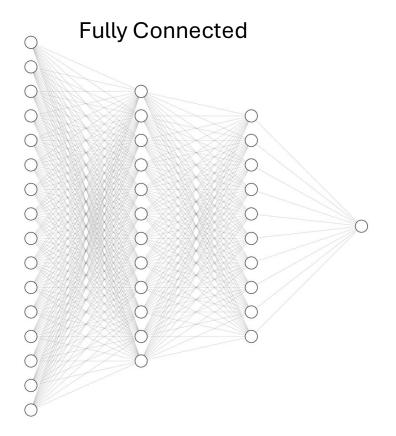
Sum of these three: $0.6 \cdot 0 + 0.1 \cdot 0.4 + 0.9 \cdot 0 = 0.4$

MLPs (also called fully-connected networks) have weights from every pixel to every neuron



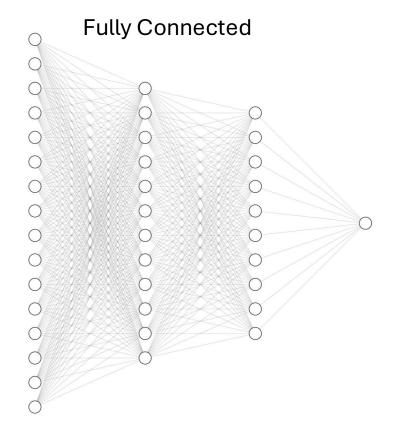
How can we change a fullyconnected network to account for spatial information?

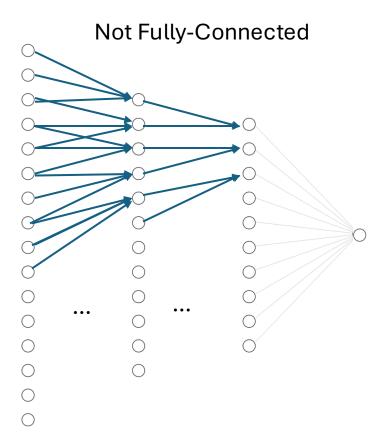
MLPs (also called fully-connected networks) have weights from every pixel to every neuron



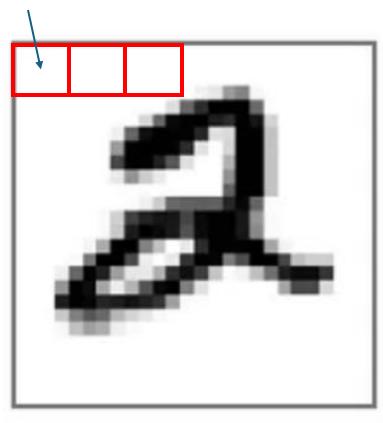
How can we change a fullyconnected network to account for spatial information?

MLPs (also called fully-connected networks) have weights from every pixel to every neuron



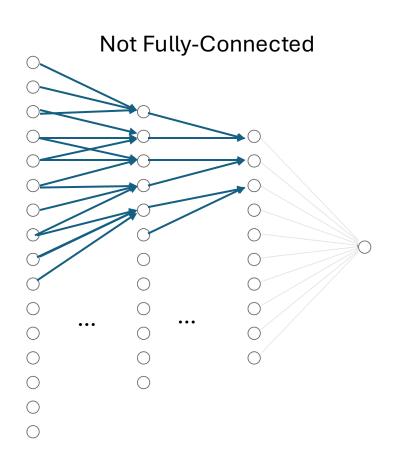


Patches: Pixels close to each other



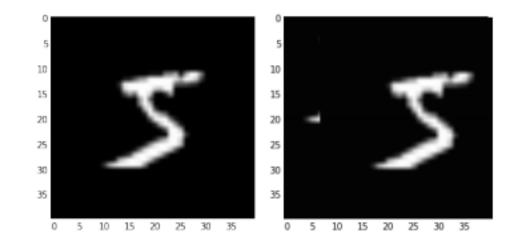
Advantages of Not Fully Connected Layers

- Fewer weights → Faster?
- The outputs of neurons are "features" for local "patches"
- Incorporates spatial information (pixels that are close together matter)



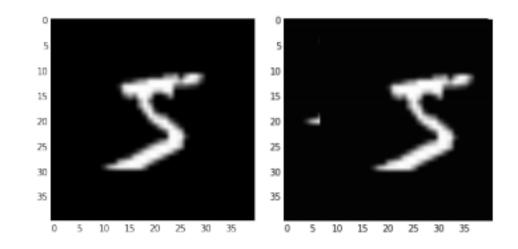
Disadvantages of Not Fully Connected Layers

- What happens if the image is Translated?
- The patches on the right side were never trained with 5's in that side.



Disadvantages of Not Fully Connected Layers

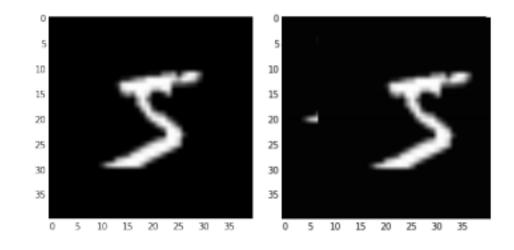
- What happens if the image is Translated?
- The patches on the right side were never trained with 5's in that side.



Even though we include spatial *information*, we still don't have spatial *reasoning*. (Can't recognize a shifted 5 is still a 5)

Disadvantages of Not Fully Connected Layers

- What happens if the image is Translated?
- The patches on the right side were never trained with 5's in that side.



Even though we include spatial *information*, we still don't have spatial *reasoning*. (Can't recognize a shifted 5 is still a 5)

What if we used the same weights for each patch? (Weight Sharing)

The Main Building Block: Convolution

Convolution is an operation that takes two inputs:

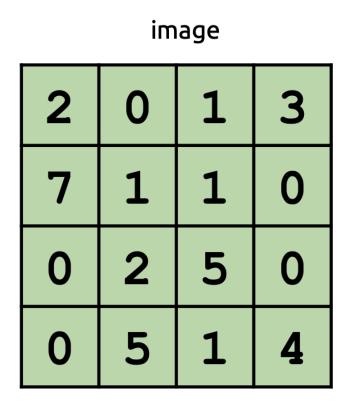
(1) An image (2D - B/W)

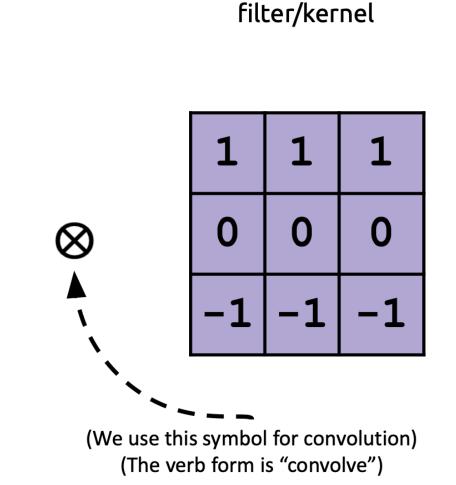
(2) A filter (also called a kernel)



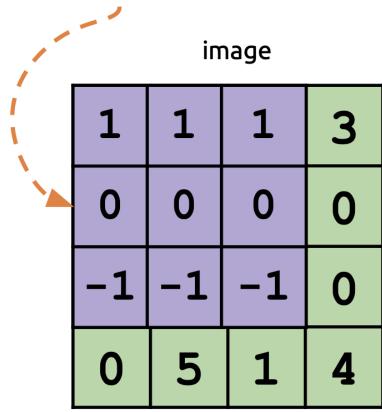
1	1	1
0	0	0
-1	-1	-1

2D array of numbers; could be any values

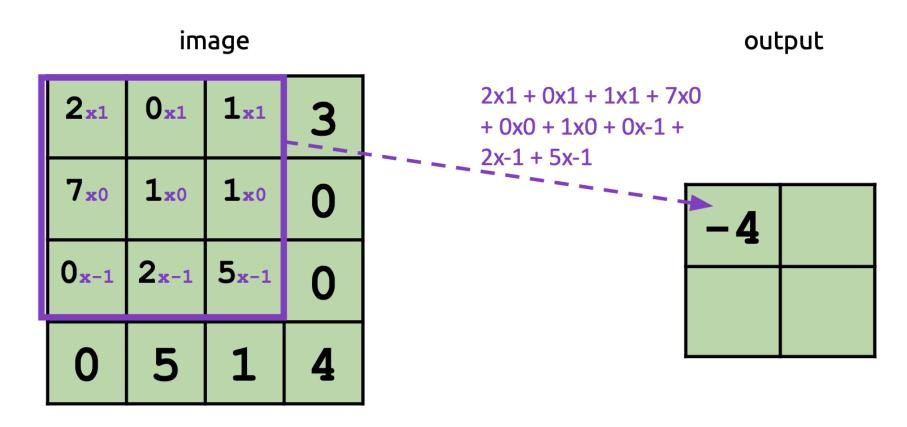




Overlay the filter on the image



Sum up multiplied values to produce output value

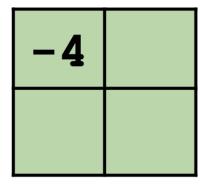


Move the filter over by one pixel

image

1	1	1	3
0	0	0	0
-1	-1	-1	0
0	5	1	4

output

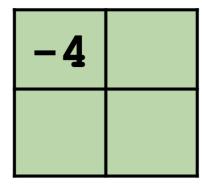


Move the filter over by one pixel

image

2	1	1	1
7	0	0	0
0	-1	-1	-1
0	5	1	4

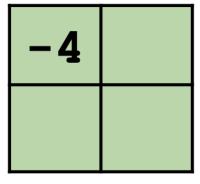
output



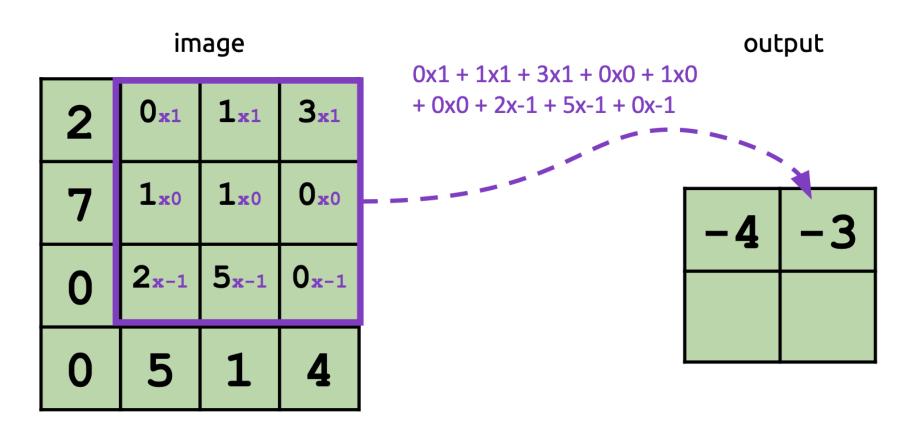
Repeat (multiply, sum up)

image

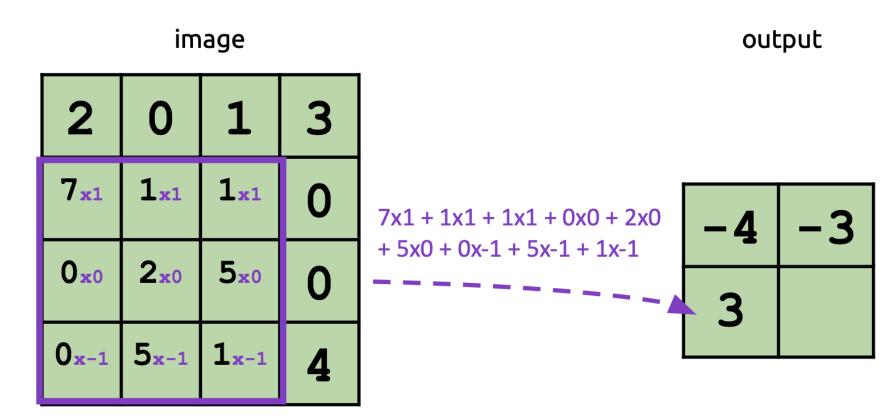
 output



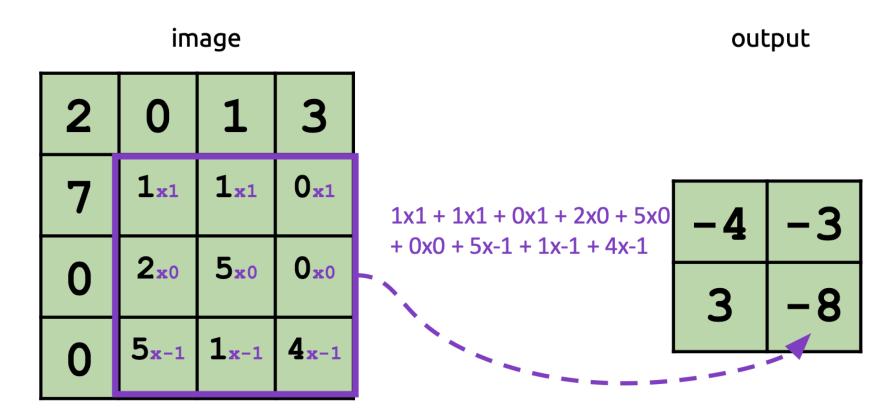
Repeat (multiply, sum up)



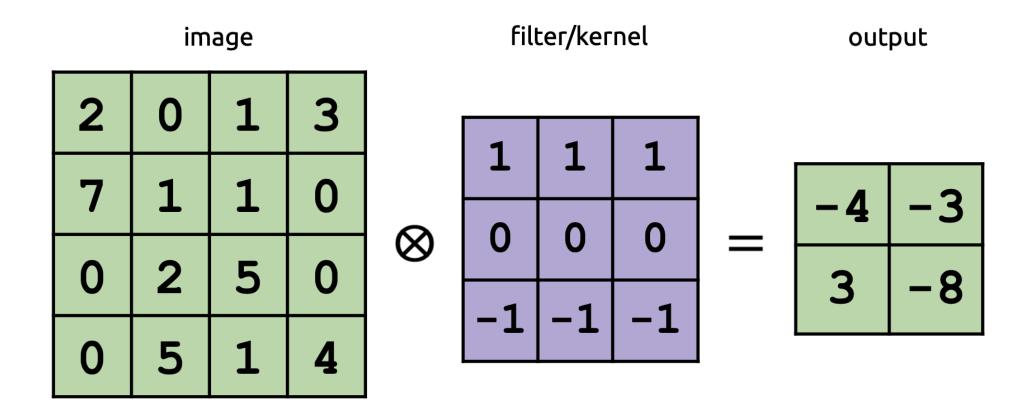
Repeat...



Repeat...



In summary:



Handmade Kernels and Filters

$$\begin{bmatrix}
 0 & 0 & 0 \\
 0 & 1 & 0 \\
 0 & 0 & 0
 \end{bmatrix}$$

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

$$\left[egin{array}{ccc} 0 & -1 & 0 \ -1 & 5 & -1 \ 0 & -1 & 0 \ \end{array}
ight]$$

Identity kernel

Edge detection

Sharpen kernel

$$\frac{1}{9} \left[\begin{array}{rrr} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{array} \right]$$

Box blur

$$\frac{1}{256} \begin{bmatrix} 1 & 4 & 6 & 4 & 1 \\ 4 & 16 & 24 & 16 & 4 \\ 6 & 24 & 36 & 24 & 6 \\ 4 & 16 & 24 & 16 & 4 \\ 1 & 4 & 6 & 4 & 1 \end{bmatrix}$$

Gaussian blurr kernel

Operation Image result g(x,y) Kernel ω Identity Ridge or edge detection Sharpen **Box blur** (normalized) Gaussian blur 3 x 3 (approximation)

Source: Wikipedia

What Comes Next?

Can we learn a filter for our images rather than "hand crafting" one?