

Machine Learning, Artificial Intelligence, and Natural Language Processing

Marcelo C. Medeiros

Department of Economics
The University of Illinois at Urbana-Champaign

Neural Networks

Introduction

Mathematical definition

Deep Neural Networks

Convolutional Neural Networks

Long Short Term Neural Networks

Introduction

What are Neural Networks (NN)?

Let's play Minecraft



Introduction

What are Neural Networks (NN)?

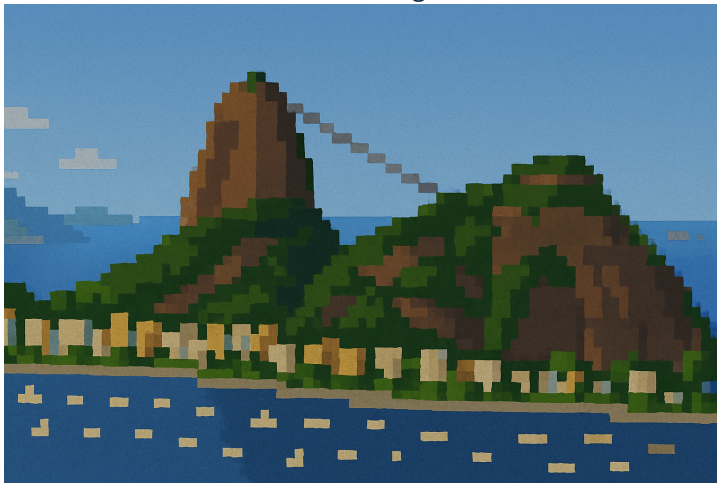
Let's reproduce this image with blocks



Introduction

What are Neural Networks (NN)?

We start with big blocks



Introduction

What are Neural Networks (NN)?

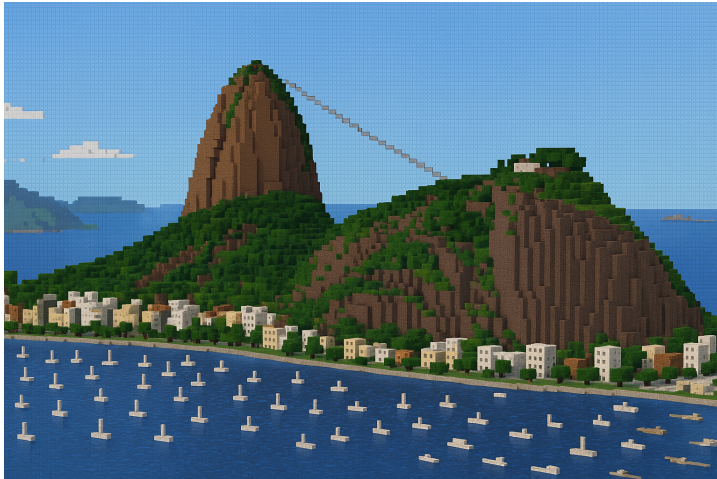
A bit smaller blocks



Introduction

What are Neural Networks (NN)?

Even smaller blocks



Introduction

What are Neural Networks (NN)?

Now let's polish the edges



Introduction

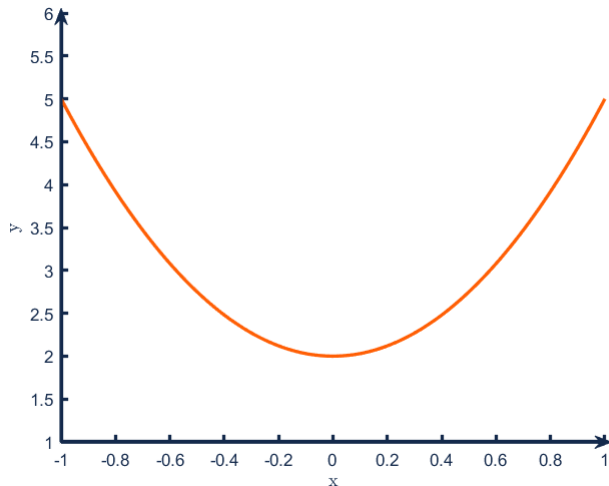
What are Neural Networks (NN)?

Our final result



Introduction

Neural Networks and function approximation

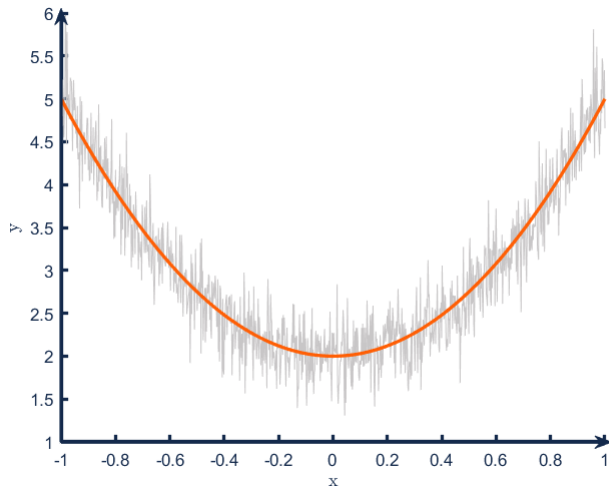


► Consider the function

$$\begin{aligned} y &= f(x) \\ &= 2 + 3x^2 \end{aligned}$$

Introduction

Neural Networks and function approximation



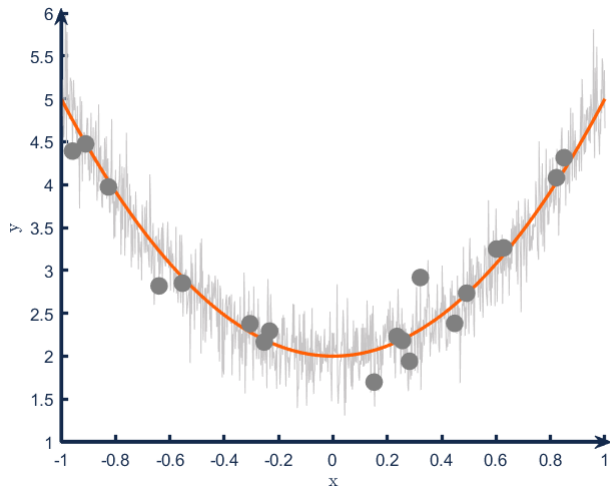
► Consider the function

$$\begin{aligned} y &= f(x) \\ &= 2 + 3x^2 \end{aligned}$$

► We observe y with error...

Introduction

Neural Networks and function approximation



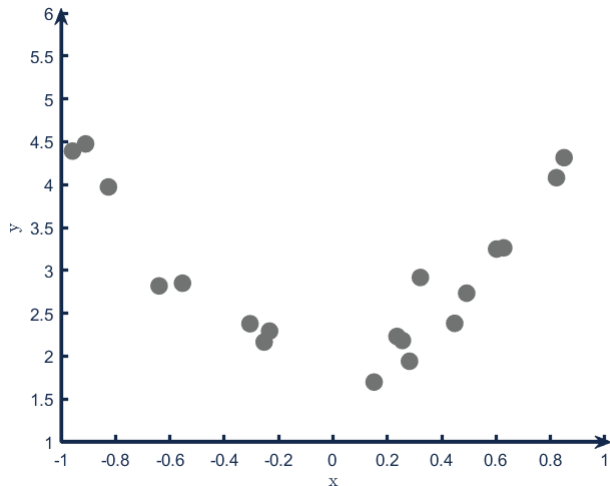
► Consider the function

$$\begin{aligned} y &= f(x) \\ &= 2 + 3x^2 \end{aligned}$$

► We observe y with error...
... and only a sample

Introduction

Neural Networks and function approximation



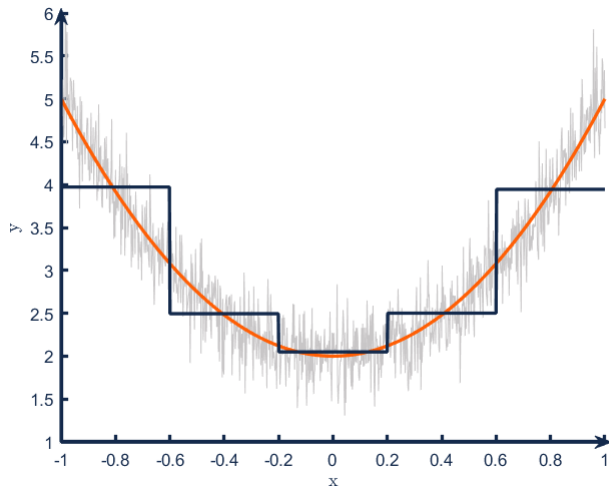
► Consider the function

$$\begin{aligned} y &= f(x) \\ &= 2 + 3x^2 \end{aligned}$$

- We observe y with error...
... and only a sample
- How can we recover the function from data?

Introduction

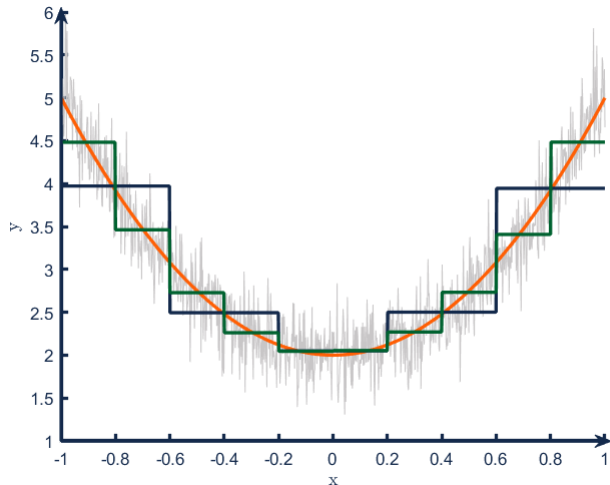
Neural Networks and function approximation



- ▶ Divide the domain of x into $K = 5$ regions
- ▶ Consider a locally constant approximation

Introduction

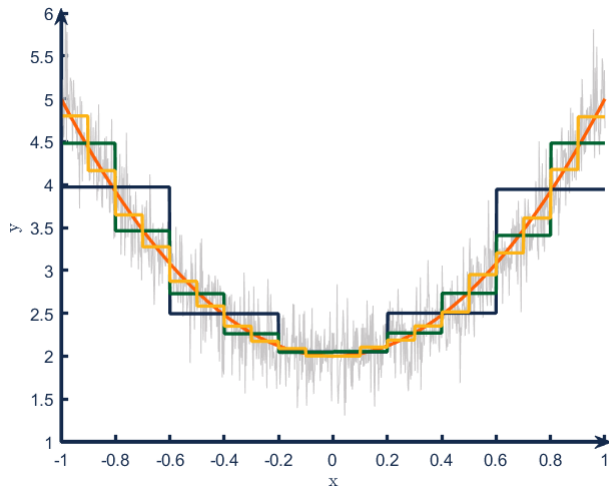
Neural Networks and function approximation



- ▶ Divide the domain of x into $K = 5$ regions
- ▶ Consider a locally constant approximation
- ▶ Approximation gets better as the number of regions grows

Introduction

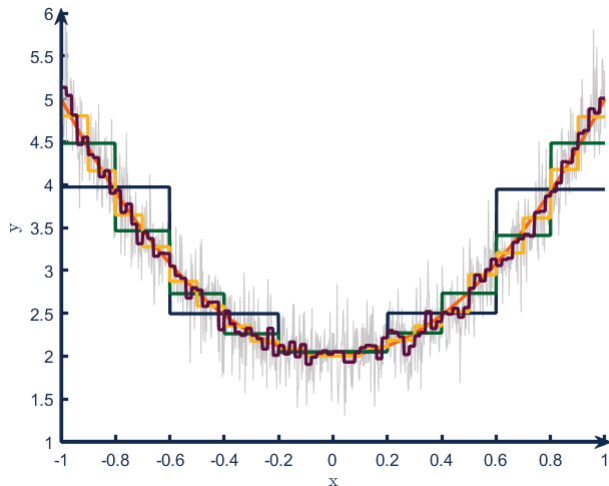
Neural Networks and function approximation



- ▶ Divide the domain of x into $K = 5$ regions
- ▶ Consider a locally constant approximation
- ▶ Approximation gets better as the number of regions grows
- ▶ $K = 10, 20, 100, \dots$

Introduction

Neural Networks and function approximation



- ▶ Divide the domain of x into $K = 5$ regions
- ▶ Consider a locally constant approximation
- ▶ Approximation gets better as the number of regions grows
- ▶ $K = 10, 20, 100, \dots$

- ▶ In the previous slides, the function $f(x) = 2 + x^3$ was being approximated by

$$h(x) = \beta_0 + \sum_{k=1}^{K-1} \beta_k \mathbb{I}(x \geq c_k),$$

where

$$\mathbb{I}(x \geq c_k) = \begin{cases} 1 & \text{if } x \geq c_k \\ 0 & \text{otherwise} \end{cases}$$

- ▶ c_1, \dots, c_K are split points and $\beta_0, \dots, \beta_{K-1}$ represent the local approximation

Introduction

Neural Networks and function approximation

► Can we do better than this?

Introduction

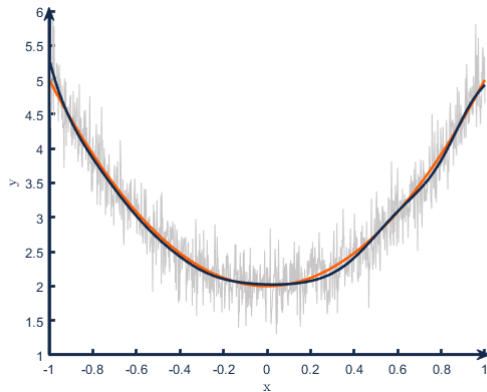
Neural Networks and function approximation

- ▶ Can we do better than this?
- ▶ Yes, we can smooth the edges

Introduction

Neural Networks and function approximation

- ▶ Can we do better than this?
- ▶ Yes, we can smooth the edges
- ▶ Five regions with smooth edges



Introduction

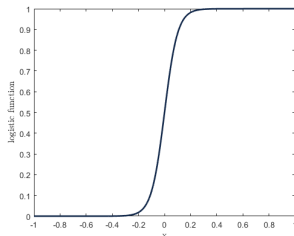
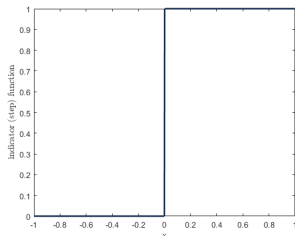
Neural Networks and function approximation

- ▶ The edges can be smoothed by replacing

$$I(x \geq c_k) = \begin{cases} 1 & \text{if } x \geq c_k \\ 0 & \text{otherwise} \end{cases}$$

by a smooth function

- ▶ One possible choice is the logistic function



This is the idea behind Neural Networks

Neural Networks

Introduction

Mathematical definition

Deep Neural Networks

Convolutional Neural Networks

Long Short Term Neural Networks

The Problem

- ▶ Observe the target variable Y and the inputs $\mathbf{X} = (X_1, \dots, X_p)'$

The Problem

- ▶ Observe the target variable Y and the inputs $\mathbf{X} = (X_1, \dots, X_p)'$
- ▶ Unknown mapping (relation) between Y and \mathbf{X} :

$$Y = f(\mathbf{X}) + U,$$

where U is a random error \rightsquigarrow the relation between Y and \mathbf{X} is not perfect

The Problem

- ▶ Observe the target variable Y and the inputs $\mathbf{X} = (X_1, \dots, X_p)'$
- ▶ Unknown mapping (relation) between Y and \mathbf{X} :

$$Y = f(\mathbf{X}) + U,$$

where U is a random error \rightsquigarrow the relation between Y and \mathbf{X} is not perfect

- ▶ From a random sample $\{Y_i, \mathbf{X}_i\}_{i=1}^n$, we would like to learn (estimate) f to predict Y^* from a new \mathbf{X}^* :

$$Y^* := f(\mathbf{X}^*)$$

The Neural Network Approach

- **NN idea:** Approximates the unknown $f(\cdot)$ by

$$H(\mathbf{X}; \boldsymbol{\theta}) = \beta_0 + \sum_{j=1}^J \beta_j S(\boldsymbol{\gamma}'_j \mathbf{X} + \gamma_{0,j}),$$

where:

- * $\mathbf{X} \mapsto \boldsymbol{\gamma}'_j \mathbf{X} + \gamma_{0,j}$ is an **affine transformation** (linear combination plus a shift) of the input vector \mathbf{X}

The Neural Network Approach

- **NN idea:** Approximates the unknown $f(\cdot)$ by

$$H(\mathbf{X}; \boldsymbol{\theta}) = \beta_0 + \sum_{j=1}^J \beta_j S(\boldsymbol{\gamma}'_j \mathbf{X} + \gamma_{0,j}),$$

where:

- * $\mathbf{X} \mapsto \boldsymbol{\gamma}'_j \mathbf{X} + \gamma_{0,j}$ is an **affine transformation** (linear combination plus a shift) of the input vector \mathbf{X}
- * $S : \mathbb{R} \rightarrow \mathbb{R}$ is a **basis function**

The Neural Network Approach

- **NN idea:** Approximates the unknown $f(\cdot)$ by

$$H(\mathbf{X}; \boldsymbol{\theta}) = \beta_0 + \sum_{j=1}^J \beta_j S(\boldsymbol{\gamma}'_j \mathbf{X} + \gamma_{0,j}),$$

where:

- * $\mathbf{X} \mapsto \boldsymbol{\gamma}'_j \mathbf{X} + \gamma_{0,j}$ is an **affine transformation** (linear combination plus a shift) of the input vector \mathbf{X}
- * $S : \mathbb{R} \rightarrow \mathbb{R}$ is a **basis function**
- * $\boldsymbol{\theta} := (\beta_0, \dots, \beta_{J_T}, \boldsymbol{\gamma}'_1, \dots, \boldsymbol{\gamma}'_J, \gamma_{0,1}, \dots, \gamma_{0,J})'$ is the vector of parameters that must be estimated

$$H(\mathbf{X}; \boldsymbol{\theta}) = \beta_0 + \sum_{j=1}^{J_T} \beta_j S(\boldsymbol{\gamma}'_j \mathbf{X} + \gamma_{0,j})$$

- ▶ The basis functions S are called **activation functions**

Some Neural Network Nomenclature

$$H(\mathbf{X}; \boldsymbol{\theta}) = \beta_0 + \sum_{j=1}^{J_T} \beta_j S(\boldsymbol{\gamma}'_j \mathbf{X} + \gamma_{0,j})$$

- ▶ The basis functions S are called **activation functions**
- ▶ The parameters θ are called **weights**

Some Neural Network Nomenclature

$$H(\mathbf{X}; \boldsymbol{\theta}) = \beta_0 + \sum_{j=1}^{J_T} \beta_j S(\boldsymbol{\gamma}'_j \mathbf{X} + \gamma_{0,j})$$

- ▶ The basis functions S are called **activation functions**
- ▶ The parameters θ are called **weights**
- ▶ In particular, β_0 and $\gamma_{0,j}$ are called **bias** \rightsquigarrow Unrelated to the statistical concept of bias

Some Neural Network Nomenclature

$$H(\mathbf{X}; \boldsymbol{\theta}) = \beta_0 + \sum_{j=1}^{J_T} \beta_j S(\boldsymbol{\gamma}'_j \mathbf{X} + \gamma_{0,j})$$

- ▶ The basis functions S are called **activation functions**
- ▶ The parameters θ are called **weights**
- ▶ In particular, β_0 and $\gamma_{0,j}$ are called **bias** \rightsquigarrow Unrelated to the statistical concept of bias
- ▶ Note that $\gamma_{0,j}$ shifts the whole S curve to the left and right

Some Neural Network Nomenclature

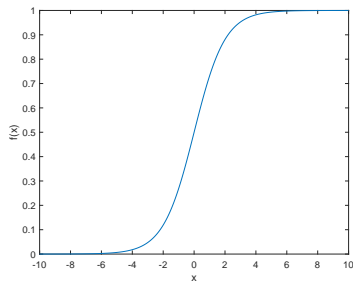
$$H(\mathbf{X}; \boldsymbol{\theta}) = \beta_0 + \sum_{j=1}^{J_T} \beta_j S(\boldsymbol{\gamma}'_j \mathbf{X} + \gamma_{0,j})$$

- ▶ The basis functions S are called **activation functions**
- ▶ The parameters θ are called **weights**
- ▶ In particular, β_0 and $\gamma_{0,j}$ are called **bias** \rightsquigarrow Unrelated to the statistical concept of bias
- ▶ Note that $\gamma_{0,j}$ shifts the whole S curve to the left and right
- ▶ While γ_j controls for the “slope” of S

Popular Choices for S

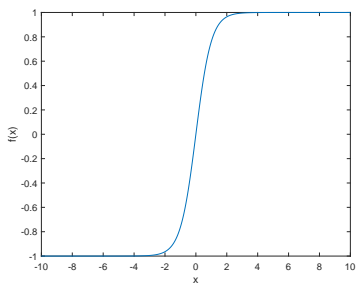
Logistic

$$S(X) = \frac{1}{1 + \exp(-X)}$$



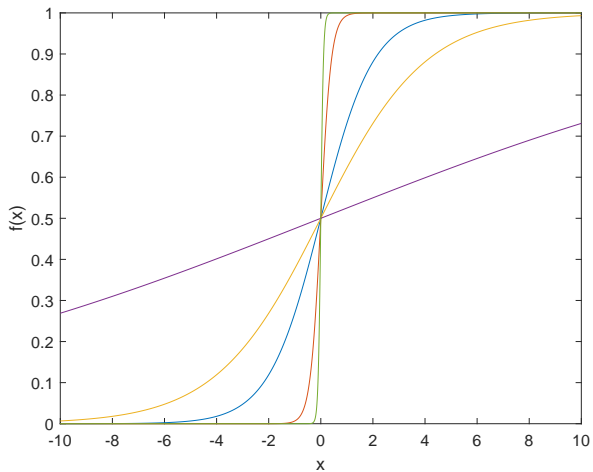
Hyperbolic tangent

$$S(X) = \frac{\exp(X) - \exp(-X)}{\exp(X) + \exp(-X)}$$



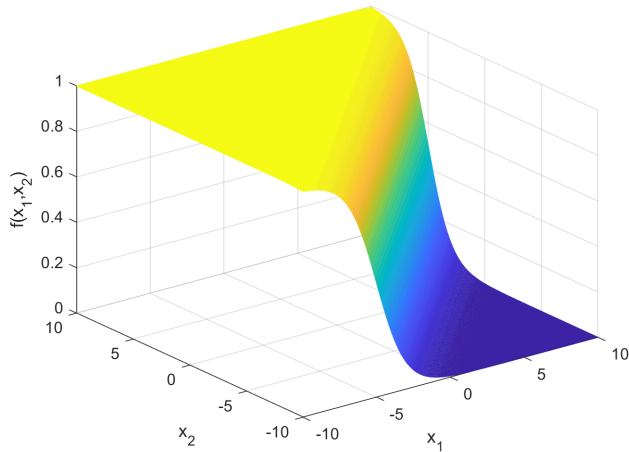
Logistic Activation Function

Varying slope parameter γ



Logistic NN with 2 Inputs

Geometric Interpretation



- ▶ In general, but not always, $S(\cdot)$ is a *squashing* function.

Squashing function

A function $S : \mathbb{R} \rightarrow [a, b]$, $a < b$, is a squashing (sigmoid) function if it is non-decreasing, $\lim_{X \rightarrow \infty} S(X) = b$ and $\lim_{X \rightarrow -\infty} S(X) = a$

- ▶ In general, but not always, $S(\cdot)$ is a *squashing* function.

Squashing function

A function $S : \mathbb{R} \rightarrow [a, b]$, $a < b$, is a squashing (sigmoid) function if it is non-decreasing, $\lim_{X \rightarrow \infty} S(X) = b$ and $\lim_{X \rightarrow -\infty} S(X) = a$

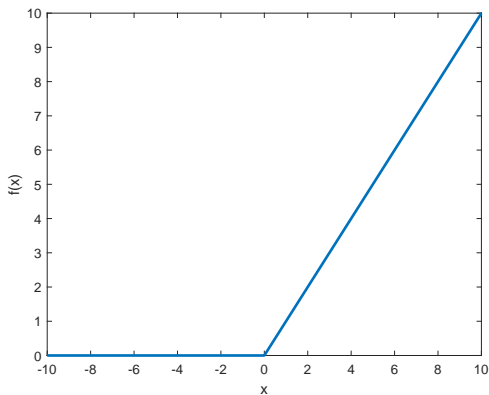
- ▶ Some popular (old) squashing functions:

- * Heaviside: $S(X) = I(X \geq 0)$
- * Logistic: $S(X) = 1/[1 + \exp(-X)]$
- * Hyperbolic tangent:
$$S(X) = [\exp(X) - \exp(-X)]/[\exp(X) + \exp(-X)]$$
- * Gaussian sigmoid: $S(X) = (2\pi)^{-1/2} \int_{-\infty}^X \exp(-u^2/2) du$
- * Cosine squasher: $S(X) = \frac{1 + \cos(X + 3\pi/2)}{2} I(|X| \leq \pi/2) + I(X > \pi/2)$

Activation Function: Rectified Linear Units

Example of a non-squashing Activation Function

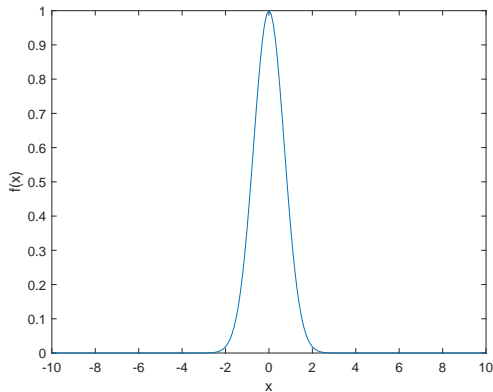
$$S(X) := \text{ReLU}(X) = \max(0, X)$$



Activation Function: Radial Basis

Example of a non-squashing Activation Function

$$S(X) := \text{RBF}(X) = \exp(-X^2)$$



Why NNs?

Universal Approximation Theorem

Feed-forward NN with a single hidden layer with “arbitrary” **squashing** functions can approximate any Borel-measurable function from one finite dimensional space to another to any desired degree of accuracy, provided sufficiently many (finite) hidden units are available

Universal Approximation Theorem

Feed-forward NN with a single hidden layer with “arbitrary” **squashing** functions can approximate any Borel-measurable function from one finite dimensional space to another to any desired degree of accuracy, provided sufficiently many (finite) hidden units are available

- ▶ How big is the class of **Borel Measurable functions**? It contains all functions with at most countable discontinuities

Universal Approximation Theorem

Feed-forward NN with a single hidden layer with “arbitrary” **squashing** functions can approximate any Borel-measurable function from one finite dimensional space to another to any desired degree of accuracy, provided sufficiently many (finite) hidden units are available

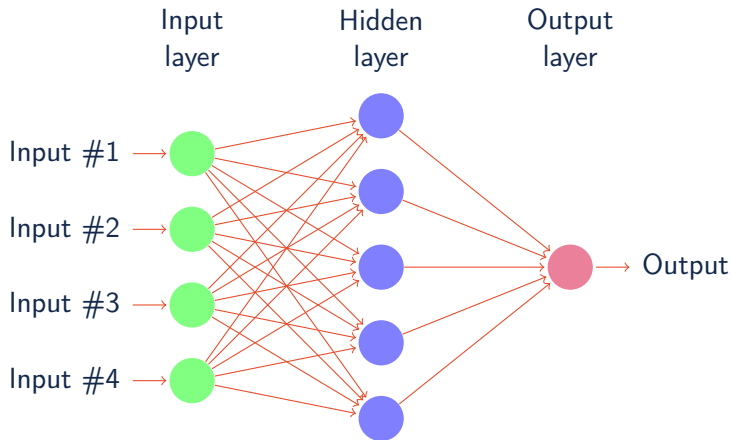
- ▶ How big is the class of **Borel Measurable functions**? It contains all functions with at most countable discontinuities
- ▶ References:
 - * Cybenko (1989); Hornik, Stimchombe, and White (NN, 1989); Gallant and White (1988); Gallant and White (NN, 1992); Hornik (NN, 1991); and many others

Universal Approximation Theorem

Feed-forward NN with a single hidden layer with “arbitrary” **squashing** functions can approximate any Borel-measurable function from one finite dimensional space to another to any desired degree of accuracy, provided sufficiently many (finite) hidden units are available

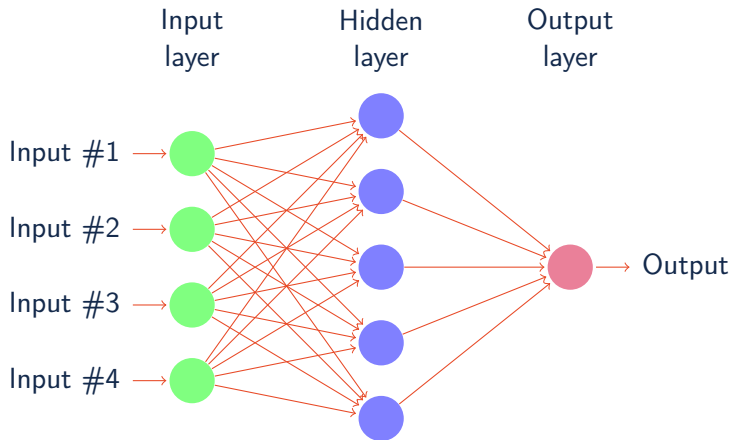
- ▶ How big is the class of **Borel Measurable functions**? It contains all functions with at most countable discontinuities
- ▶ References:
 - * Cybenko (1989); Hornik, Stimchombe, and White (NN, 1989); Gallant and White (1988); Gallant and White (NN, 1992); Hornik (NN, 1991); and many others
- ▶ The same NN can approximate the derivatives of the function

NN as a Graph



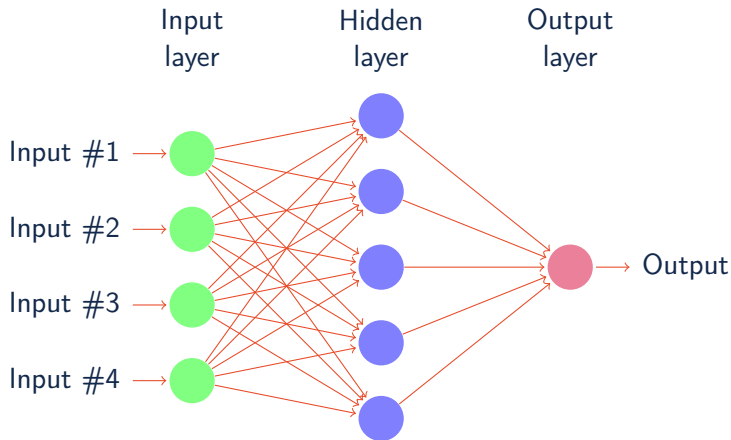
- The input layer is just the vector of explanatory variables.

NN as a Graph



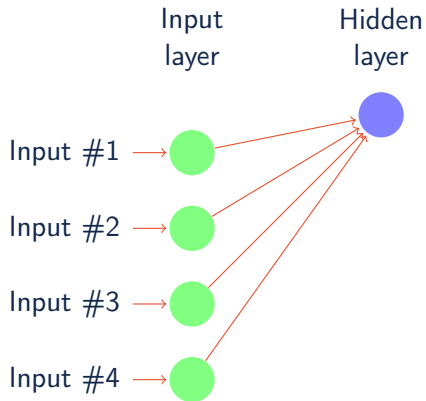
- ▶ The input layer is just the vector of explanatory variables.
- ▶ The hidden layer consists of a set of hidden units (*neurons*)

NN as a Graph



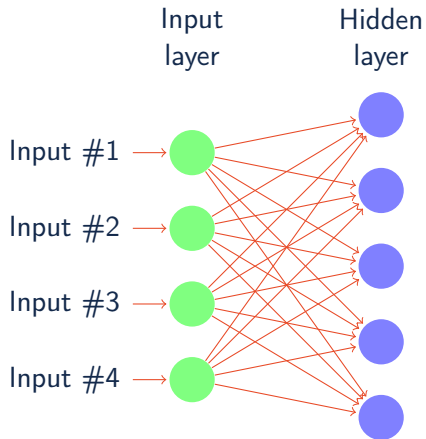
- ▶ The input layer is just the vector of explanatory variables.
- ▶ The hidden layer consists of a set of hidden units (*neurons*)
- ▶ The output layer is the predicted value for the dependent variables

NN as a Graph



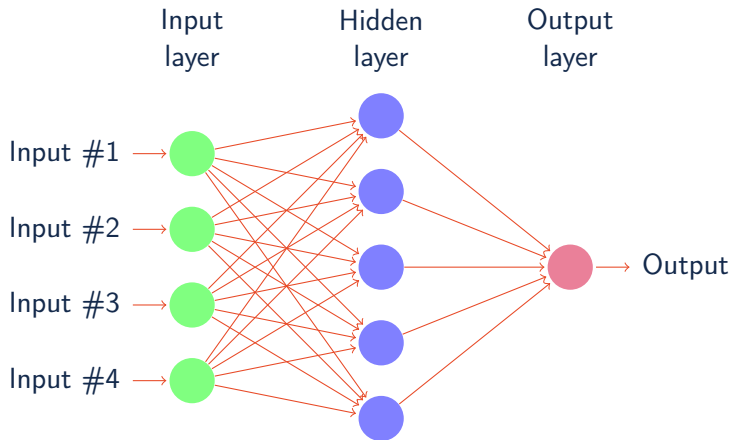
- At each unit inputs are linearly combined: $\gamma'_j \mathbf{X} + \gamma_{0,j}$

NN as a Graph



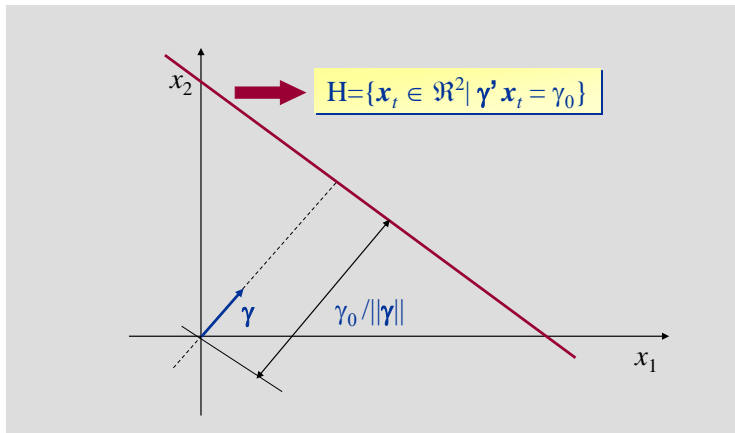
- ▶ At each unit inputs are linearly combined: $\gamma'_j \mathbf{X} + \gamma_{0,j}$
- ▶ Nonlinear transformation in the hidden layer: $S(\gamma'_j \mathbf{X} + \gamma_{0,j})$

NN as a Graph

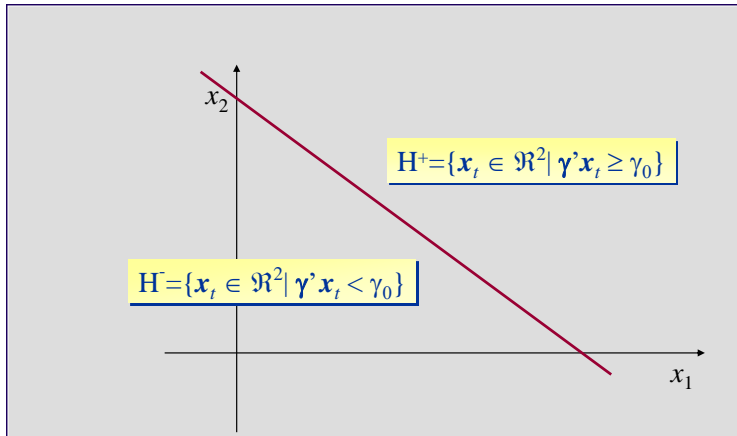


- ▶ At each unit inputs are linearly combined: $\gamma_j' \mathbf{X} + \gamma_{0,j}$
- ▶ Nonlinear transformation in the hidden layer: $S(\gamma_j' \mathbf{X} + \gamma_{0,j})$
- ▶ Outputs of the hidden layer are linearly combined $\beta_0 + \sum_{j=1}^5 \beta_j S(\gamma_j' \mathbf{X} + \gamma_{0,j})$

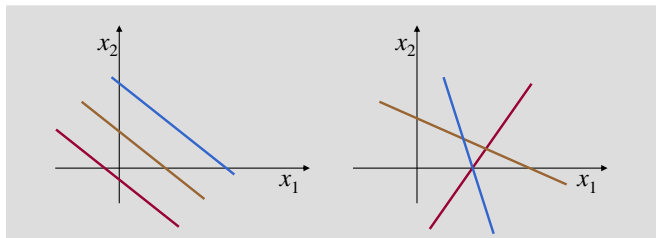
Geometric Interpretation



Geometric Interpretation



Geometric Interpretation



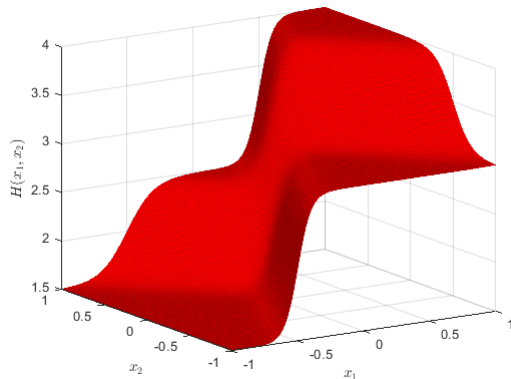
- ▶ J hyperplanes divide the space of covariates into several polyhedral regions. The maximum number of regions is given by

$$M(J, p) = M(J - 1, p) + M(J - 1, p - 1),$$

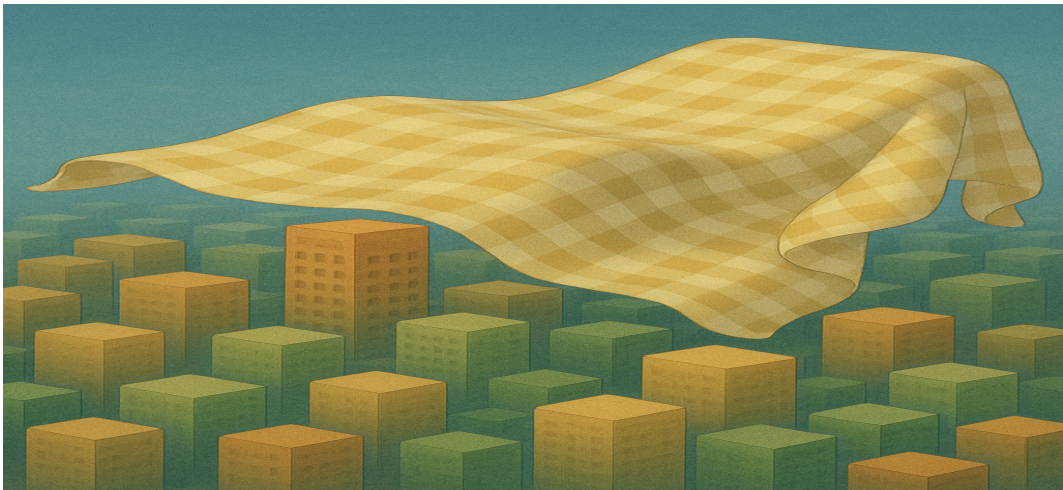
where $M(1, p) = 2$ and $M(J, 1) = J + 1$.

Geometric Interpretation

- ▶ In each region, the local model is a constant
- ▶ Smooth transition between regions
- ▶ The number of regions and the degree of smoothness determine the quality of the approximation



Geometric Interpretation



Agenda for Today

Neural Networks

Introduction

Mathematical definition

Deep Neural Networks

Convolutional Neural Networks

Long Short Term Neural Networks

From Shallow to Deep: What is a Deep Neural Network?

- ▶ It is “just” a neural network with more than one hidden layer

From Shallow to Deep: What is a Deep Neural Network?

- ▶ It is “just” a neural network with more than one hidden layer
- ▶ The layer might be fully connected or not

From Shallow to Deep: What is a Deep Neural Network?

- ▶ It is “just” a neural network with more than one hidden layer
- ▶ The layer might be fully connected or not
- ▶ Different number of units in each hidden layer

From Shallow to Deep: What is a Deep Neural Network?

- ▶ It is “just” a neural network with more than one hidden layer
- ▶ The layer might be fully connected or not
- ▶ Different number of units in each hidden layer
- ▶ Different activation function: rectified linear units (ReLU)

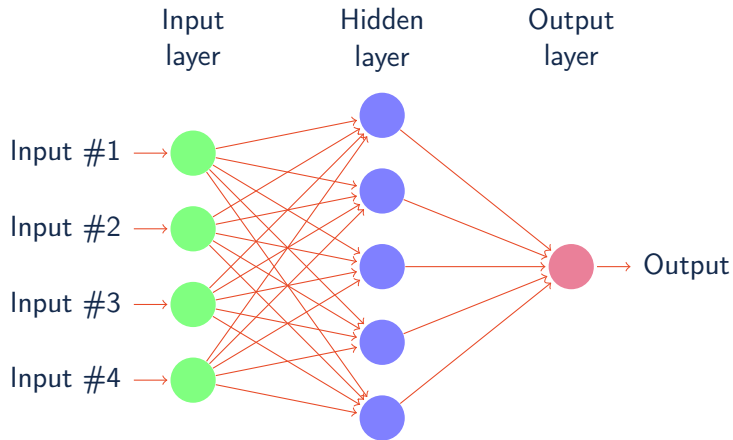
What are the Potential Advantages Over Shallow NNs?

- ▶ It has been very successful in many complex applications:
 - * Google Neural Machine Translation
 - * Lip reading
 - * Google Maps and Street View

What are the Potential Advantages Over Shallow NNs?

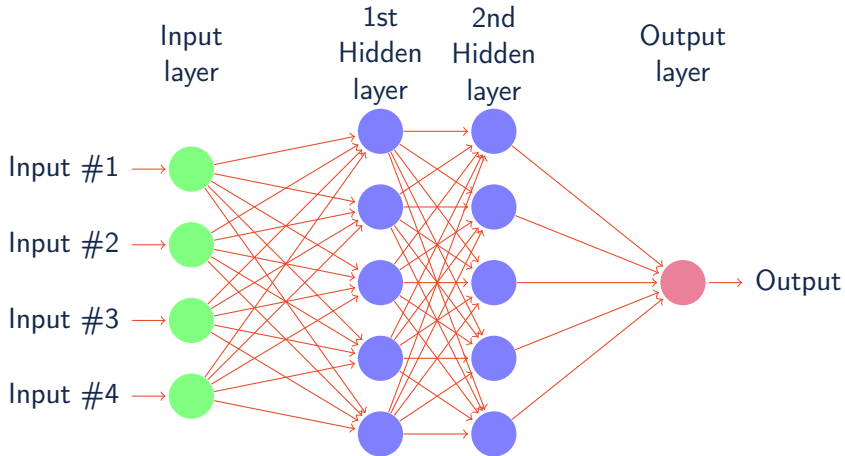
- ▶ It has been very successful in many complex applications:
 - * Google Neural Machine Translation
 - * Lip reading
 - * Google Maps and Street View
- ▶ Less hidden units per layer. “While the universal approximation property holds both for hierarchical and shallow networks, deep networks can approximate the class of compositional functions as well as shallow networks but with exponentially lower number of training parameters and sample complexity.”

Graphical Representation of a (Shallow) NN



Graphical Representation of a Deep Neural Networks

2 fully connected layers



Agenda for Today

Neural Networks

Introduction

Mathematical definition

Deep Neural Networks

Convolutional Neural Networks

Long Short Term Neural Networks

Convolutional Neural Networks

Introduction

- ▶ Convolutional Neural Networks (CNNs) are a class of Neural Network models that have proven successful in image recognition and classification.
- ▶ Multi-layer network consisting of different key elements:
 - * Convolutional layer (one or more)
 - * Nonlinear transformation
 - * Pooling (dimension reduction)
 - * Fully-connected (deep) feed-forward neural network

Convolutional Neural Networks

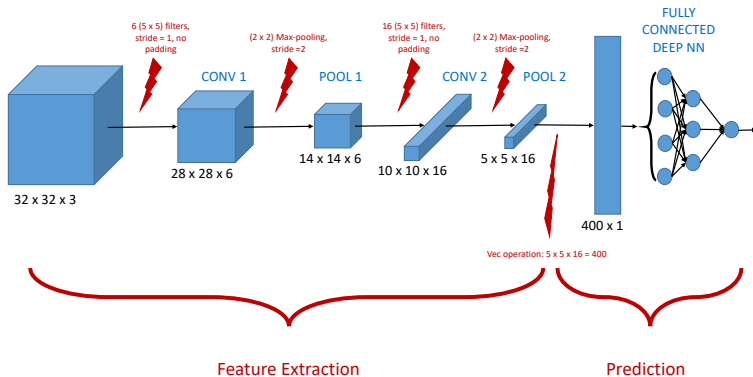
Introduction

► Sequence of layers:

convolution + nonlinear transformation \rightarrow pooling \rightarrow
convolution + nonlinear transformation \rightarrow pooling $\rightarrow \dots \rightarrow$
convolution + nonlinear transformation \rightarrow pooling \rightarrow
Fully-connected (deep) NN

Convolutional Neural Networks

Introduction



Convolutional Neural Networks

Filters (Kernels)

- ▶ To a computer, an image is a matrix of pixels.
- ▶ Each entry of the matrix is the intensity of the pixel: 0 – 255 (grayscale)
- ▶ The dimension of the matrix is the resolution of the image.
- ▶ For colored images, there is a third dimension to represent the color channels: Red (R), Green (G) and Blue (B).
- ▶ Therefore the image is a three-dimensional matrix (tensor): $\text{Height} \times \text{Width} \times 3$.

Convolutional Neural Networks

Filters (Kernels)

- ▶ An **image kernel** is a small matrix used to apply effects, such as blurring, sharpening, outlining or embossing.
- ▶ In Machine Learning, kernels are used for “feature extraction”, a technique for determining the most important portions of an image.
- ▶ In this context, the process is referred to more generally as **convolution**.
- ▶ A nice webpage:
<https://setosa.io/ev/image-kernels/>

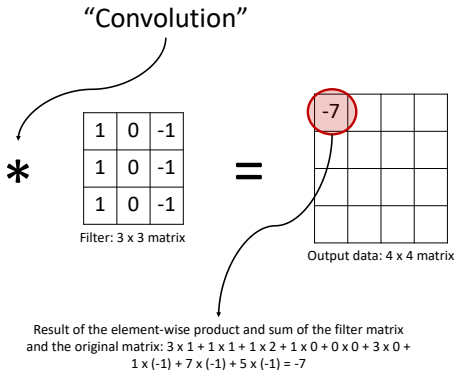
The Convolutional Layer

- ▶ Input data: $\mathbf{X} \in \mathbb{R}^{M \times N}$
- ▶ Filter (kernel): $\mathbf{W} \in \mathbb{R}^{Q \times R}$
 - * \mathbf{W} is usually unknown
- ▶ Output \mathbf{O} is of a smaller dimension than the input due to **border effects**. For $i = 1, \dots, M - Q + 1, j = 1, \dots, N - R + 1$:
 - * $O_{ij} = \sum_{q=1}^Q \sum_{r=1}^R [\mathbf{W} \odot [\mathbf{X}]_{i:i+Q-1, j:j+R-1}]_{q,r}$.

The Convolutional Layer

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

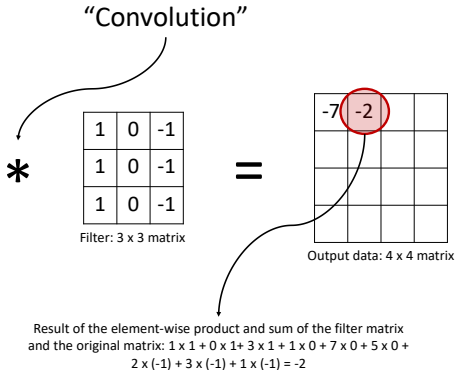
Input data: 6 x 6 matrix



The Convolutional Layer

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

Input data: 6 x 6 matrix



The Convolutional Layer

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

Input data: 6 x 6 matrix

“Convolution”

*

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

Filter: 3 x 3 matrix

=

-7	-2	2	

Output data: 4 x 4 matrix

Result of the element-wise product and sum of the filter matrix and the original matrix: $1 \times 1 + 7 \times 1 + 5 \times 1 + 2 \times 0 + 3 \times 0 + 1 \times 0 + 8 \times (-1) + 2 \times (-1) + 1 \times (-1) = 2$

The Convolutional Layer

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

Input data: 6 x 6 matrix

*

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

Filter: 3 x 3 matrix

=

-7	-2	2	-7

Output data: 4 x 4 matrix

Result of the element-wise product and sum of the filter matrix
and the original matrix: $2 \times 1 + 3 \times 1 + 1 \times 1 + 8 \times 0 + 2 \times 0 + 1 \times 0 +$
 $4 \times (-1) + 6 \times (-1) + 3 \times (-1) = -7$

The Convolutional Layer

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

Input data: 6 x 6 matrix

“Convolution”

*

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

Filter: 3 x 3 matrix

=

-7	-2	2	-7
-9			

Output data: 4 x 4 matrix

Result of the element-wise product and sum of the filter matrix and the original matrix: $1 \times 1 + 2 \times 1 + 1 \times 1 + 0 \times 0 + 3 \times 0 + 4 \times 0 + 7 \times (-1) + 5 \times (-1) + 1 \times (-1) = -9$

The Convolutional Layer

- For $i = 1, \dots, M - Q + 1, j = 1, \dots, N - R + 1$:

$$O_{ij} = \sum_{q=1}^Q \sum_{r=1}^R [\mathbf{W} \odot [\mathbf{X}]_{i:i+Q-1,j:j+R-1}]_{q,r}$$

$$O_{ij} = \boldsymbol{\iota}'_Q (\mathbf{W} \odot [\mathbf{X}]_{i:i+Q-1,j:j+R-1}) \boldsymbol{\iota}_R$$

where:

- * \odot is the element-by-element multiplication;
- * $\boldsymbol{\iota}_Q \in \mathbb{R}^Q$ and $\boldsymbol{\iota}_R^R$ are vector of ones;
- * $[\mathbf{X}]_{i:i+Q-1,j:j+R-1}$ is the block of the matrix \mathbf{X} running from row i to row $i + Q - 1$ and from column j to column $j + R - 1$; and
- * $[\mathbf{X}]_{i:i+Q-1,j:j+R-1}]_{q,r}$ is the element of $[\mathbf{X}]_{i:i+Q-1,j:j+R-1}$ in position (q, r) .

The Convolutional Layer

- ▶ O_{ij} is the discrete convolution between \mathbf{W} and $[\mathbf{X}]_{i:i+Q-1,j:j+R-1}$:

$$O_{ij} = \mathbf{W} * [\mathbf{X}]_{i:i+Q-1,j:j+R-1}$$

The Convolutional Layer

Stride

Stride (downsampling) \rightarrow reduce problem dimension

How many “pixels” we move at each step.

“Convolution”
with stride 1

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

Input data: 6 x 6 matrix

“Convolution”
with stride 2

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

Input data: 6 x 6 matrix

“Convolution”
with stride 3

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

Input data: 6 x 6 matrix



The Convolutional Layer

Stride

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

Stride = 1: slide one “pixel” each step

For a 6 x 6 original matrix, the convolution results in a 4 x 4 matrix

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

Stride = 2: slide two “pixels” each step

For a 6 x 6 original matrix, the convolution results in a 3 x 3 matrix (problems with edge effects)

3	1	1	2	8	4
1	0	7	3	2	6
2	3	5	1	1	3
1	4	1	2	6	5
3	2	1	3	7	2
9	2	6	2	5	1

Stride = 3: slide three “pixels” each step

For a 6 x 6 original matrix, the convolution results in a 2 x 2 matrix

The Convolutional Layer

Padding

Padding \rightarrow avoid border effects
Output and input with the same dimension

0	0	0	0	0	0	0	0
0	3	1	1	2	8	4	0
0	1	0	7	3	2	6	0
0	2	3	5	1	1	3	0
0	1	4	1	2	6	5	0
0	3	2	1	3	7	2	0
0	9	2	6	2	5	1	0
0	0	0	0	0	0	0	0

The Convolutional Layer

Padding

0	0	0	0	0	0	0	0
0	3	1	1	2	8	4	0
0	1	0	7	3	2	6	0
0	2	3	5	1	1	3	0
0	1	4	1	2	6	5	0
0	3	2	1	3	7	2	0
0	9	2	6	2	5	1	0
0	0	0	0	0	0	0	0

 $*$

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

 $=$

3					

The Convolutional Layer

Padding

0	0	0	0	0	0	0	0
0	3	1	1	2	8	4	0
0	1	0	7	3	2	6	0
0	2	3	5	1	1	3	0
0	1	4	1	2	6	5	0
0	3	2	1	3	7	2	0
0	9	2	6	2	5	1	0
0	0	0	0	0	0	0	0

 $*$

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

 $=$

3	-4				

The Convolutional Layer

Padding

0	0	0	0	0	0	0	0
0	3	1	1	2	8	4	0
0	1	0	7	3	2	6	0
0	2	3	5	1	1	3	0
0	1	4	1	2	6	5	0
0	3	2	1	3	7	2	0
0	9	2	6	2	5	1	0
0	0	0	0	0	0	0	0

 $*$

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

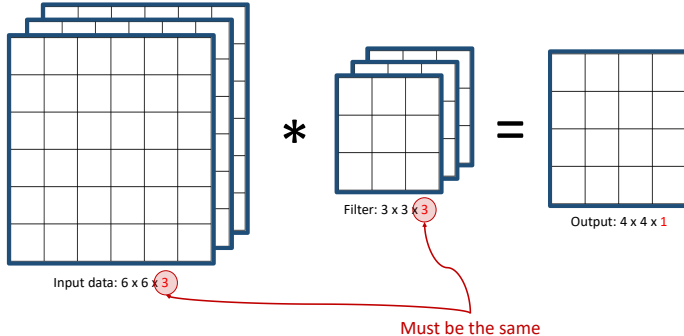
 $=$

3	-4	-2			

The Convolutional Layer

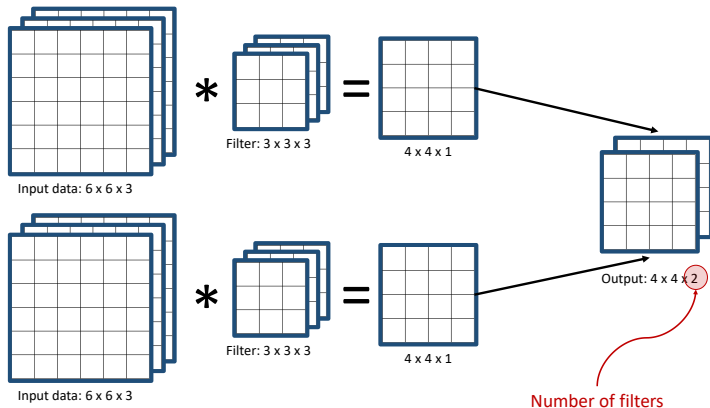
Multiple Input Channels

Three channels define colored images: R (red), B (blue), and G (green).



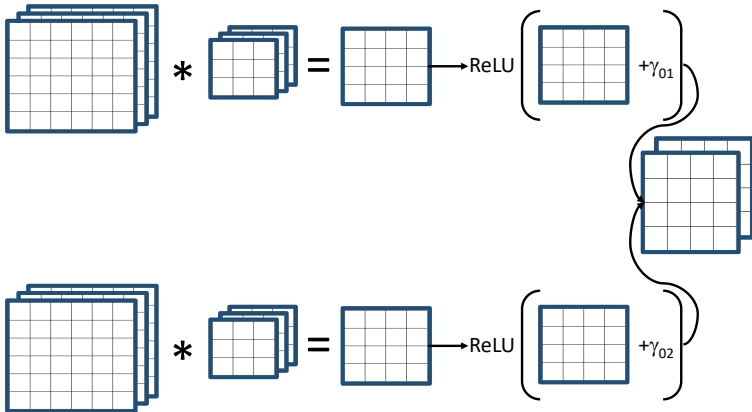
The Convolutional Layer

Multiple Filters

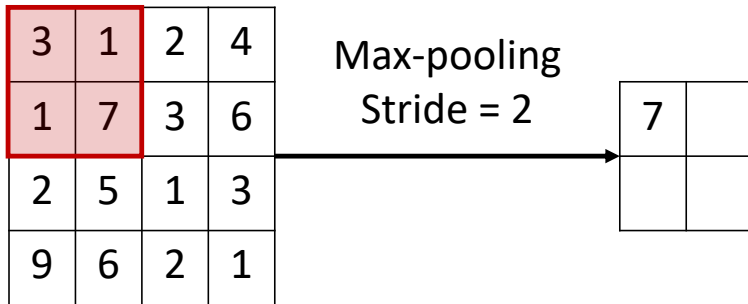


The Convolutional Layer

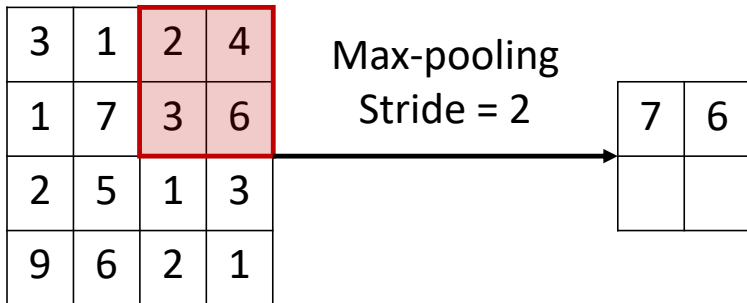
Nonlinearity



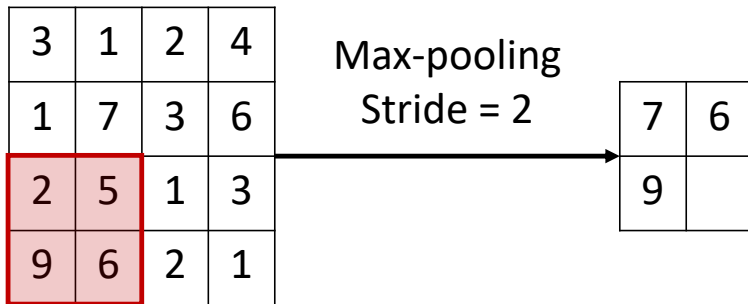
The Pooling Layer



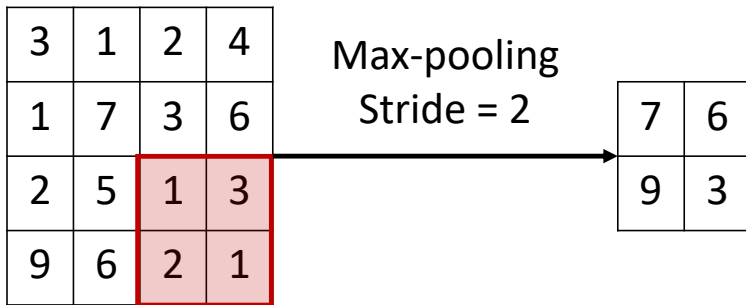
The Pooling Layer



The Pooling Layer

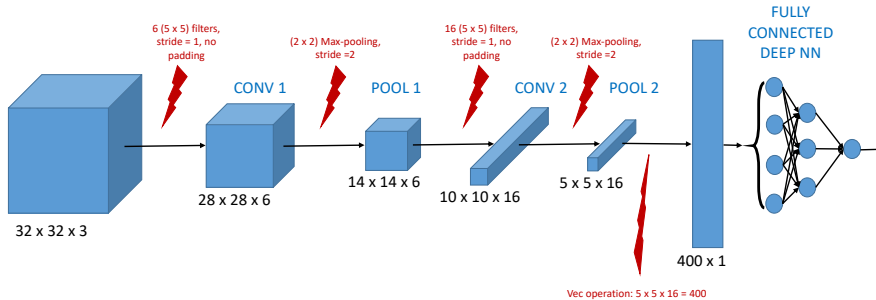


The Pooling Layer



The Convolutional Neural Network

Putting Pieces Together



The Convolutional Neural Network

Architecture Setting and Estimation

► Hyperparameters:

1. number of convolution layers (C);
2. number of pooling layers (P);
3. number (K_c) and dimensions (Q_c height, R_c width and S_c depth) of filters in each convolution layer $c = 1, \dots, C$;
4. architecture of the deep neural network.

► Parameters:

1. Filter weights: $\mathbf{W}_{ic} \in \mathbb{R}^{Q_c \times R_c \times S_c}$, $i = 1, \dots, K_c$, $c = 1, \dots, C$;
2. ReLU biases: $\gamma_c \in \mathbb{R}^{K_c}$, $c = 1, \dots, C$;
3. All the parameters of the fully connected deep NN: ψ .

Agenda for Today

Neural Networks

Introduction

Mathematical definition

Deep Neural Networks

Convolutional Neural Networks

Long Short Term Neural Networks

▶ Simple Recurrent Neural Network

$$h_t = \theta_h f(h_{t-1}) + \theta_x x_t$$

$$\hat{y}_t = \theta_y f(h_t)$$

▶ RNNs suffer from the vanishing/exploding gradient problem.

- * Set the cost function to be

$$\mathcal{Q}_T(\theta) = \frac{1}{T} \sum_{t=1}^T (y_t - \hat{y}_t)^2$$

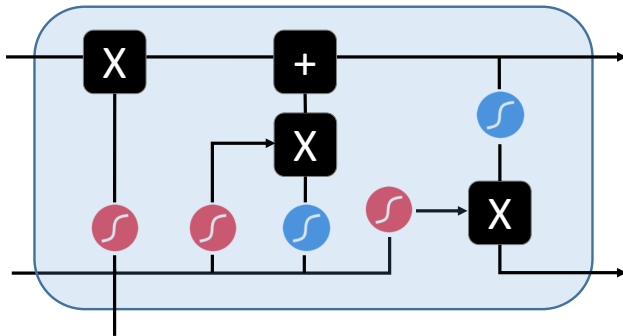
- * $\frac{\partial \mathcal{Q}_T(\theta)}{\partial \theta}$ can be very small or diverge.

▶ Solution: LSTM



Long Short Term Memory Networks

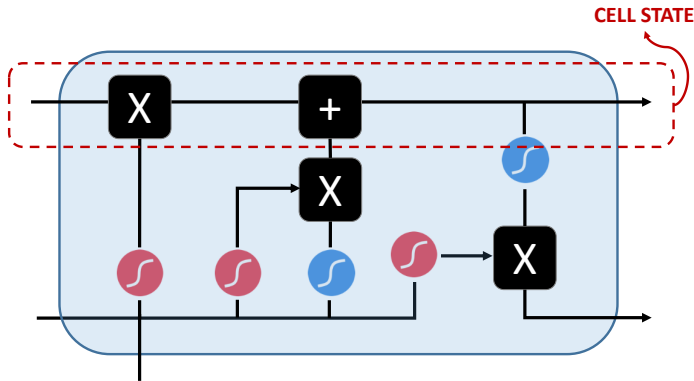
The LSTM Cell



- ▶ Red circles: logistic functions
- ▶ Blue circles: hyperbolic tangent functions

Long Short Term Memory Networks

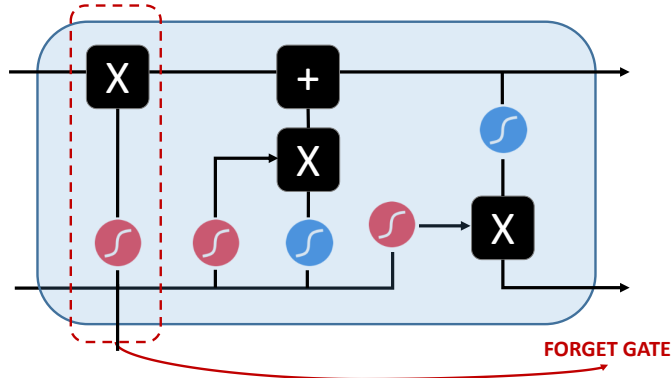
The LSTM Cell: The Cell State



- ▶ The cell state: a bit of memory to the LSTM to “remember” the past.
- ▶ LSTM learns to keep only relevant information and forget nonrelevant data.

Long Short Term Memory Networks

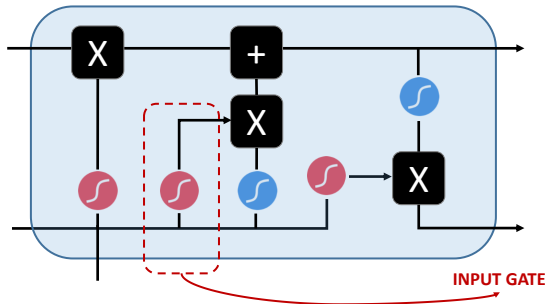
The LSTM Cell: The Forget Gate



- ▶ The forget gate tells which information to throw away from the cell state.
- ▶ It is composed of a logistic function

Long Short Term Memory Networks

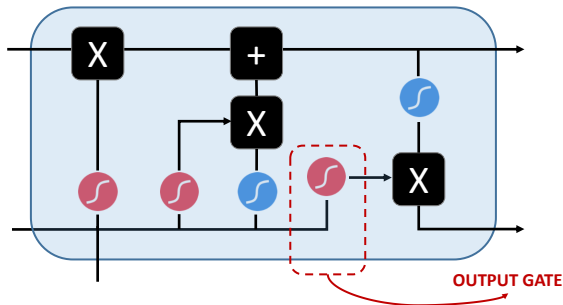
The LSTM Cell: The Input Gate



- ▶ The input gate tells which new information should be stored in the cell state.
- ▶ It is composed of a logistic function

Long Short Term Memory Networks

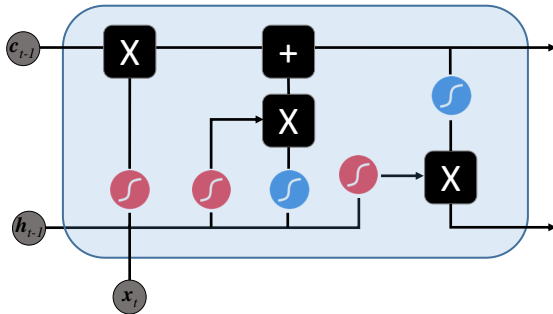
The LSTM Cell: The Output Gate



- ▶ The output gate provides the activation to the final output of the LSTM block at time t .
- ▶ It is composed of a logistic function

Long Short Term Memory Networks

The Information Flow

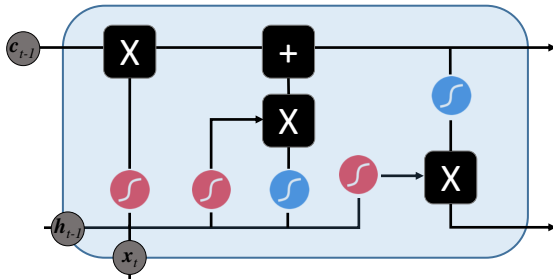


► At time t :

1. Inputs: $x_t \in \mathbb{R}^p$ and past *hidden state* $h_{t-1} \in \mathbb{R}^q$
2. Running state cell: $c_{t-1} \in \mathbb{R}^q$

Long Short Term Memory Networks

The Information Flow

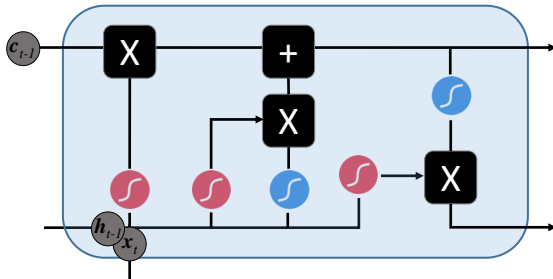


► At time t :

1. Inputs: $x_t \in \mathbb{R}^p$ and past *hidden state* $h_{t-1} \in \mathbb{R}^q$
2. Running state cell: $c_{t-1} \in \mathbb{R}^q$

Long Short Term Memory Networks

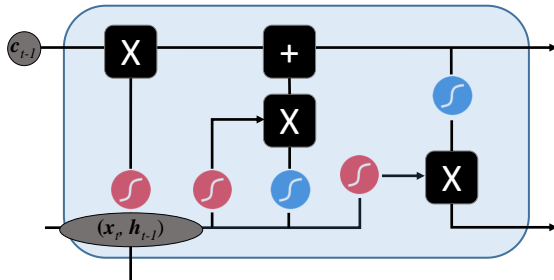
The Information Flow



► At time t :

1. Inputs: $x_t \in \mathbb{R}^p$ and past hidden state $h_{t-1} \in \mathbb{R}^q$
2. Running state cell: $c_{t-1} \in \mathbb{R}^q$

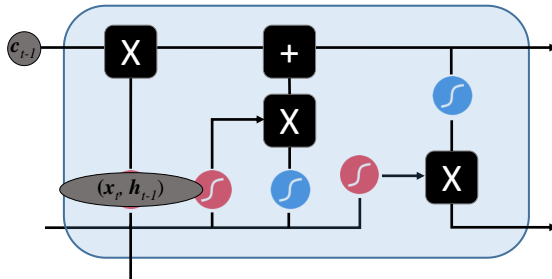
The Information Flow



- ▶ At time t :
 1. Inputs concatenate: $z_t = (x'_t, h'_{t-1})'$
 2. Running state cell: $c_{t-1} \in \mathbb{R}^q$

Long Short Term Memory Networks

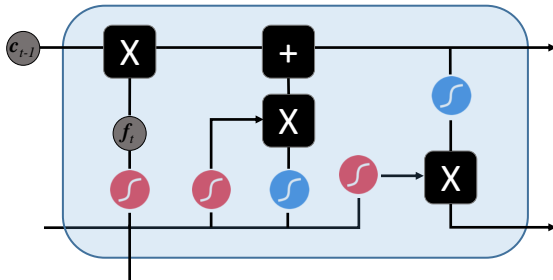
The Information Flow



► At time t :

1. Logistic activation: $f_t = \text{logistic}(\Gamma_f z_t + \gamma_{0f}) \in \mathbb{R}^q$
2. Running state cell: $c_{t-1} \in \mathbb{R}^q$

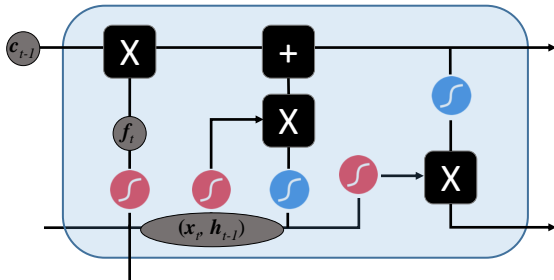
The Information Flow



- At time t :
 1. Logistic activation: $\mathbf{f}_t = \text{logistic}(\mathbf{\Gamma}_f \mathbf{z}_t + \gamma_{0f}) \in \mathbb{R}^q$
 2. Running state cell: $\mathbf{c}_{t-1} \in \mathbb{R}^q$

Long Short Term Memory Networks

The Information Flow

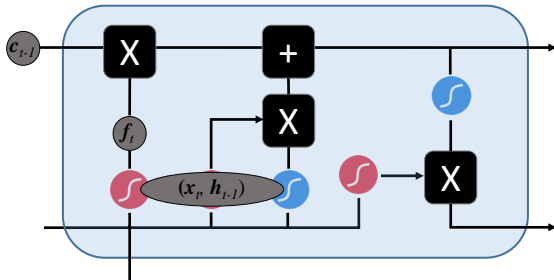


► At time t :

1. Logistic activation: $i_t = \text{logistic}(\Gamma'_i z_t + \gamma_{0i}) \in \mathbb{R}^q$
2. Running state cell: $c_{t-1} \in \mathbb{R}^q$

Long Short Term Memory Networks

The Information Flow

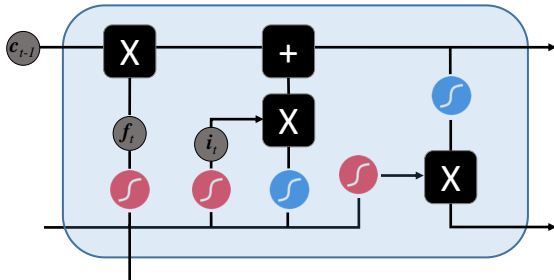


► At time t :

1. Logistic activation: $i_t = \text{logistic}(\Gamma'_i z_t + \gamma_{0i}) \in \mathbb{R}^q$
2. Running state cell: $c_{t-1} \in \mathbb{R}^q$

Long Short Term Memory Networks

The Information Flow

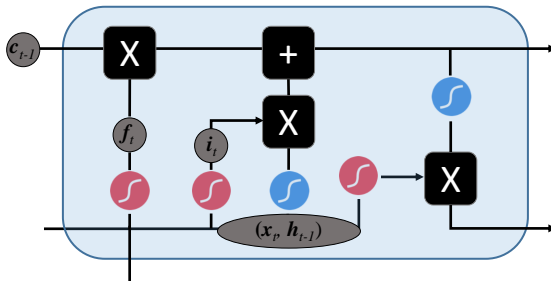


► At time t :

1. Logistic activation: $i_t = \text{logistic}(\Gamma'_i z_t + \gamma_{0i}) \in \mathbb{R}^q$
2. Running state cell: $c_{t-1} \in \mathbb{R}^q$

Long Short Term Memory Networks

The Information Flow

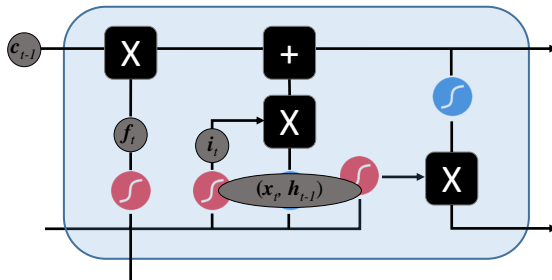


► At time t :

1. Logistic activation: $i_t = \text{logistic}(\Gamma'_i z_t + \gamma_{0i}) \in \mathbb{R}^q$
2. Tanh activation: $p_t = \text{tanh}(\Gamma'_p z_t + \gamma_{0p}) \in \mathbb{R}^q$ (potential state cell)

Long Short Term Memory Networks

The Information Flow

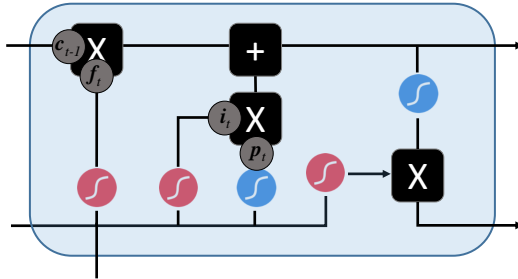


► At time t :

1. Logistic activation: $i_t = \text{logistic}(\Gamma'_i z_t + \gamma_{0i}) \in \mathbb{R}^q$
2. Tanh activation: $p_t = \tanh(\Gamma'_p z_t + \gamma_{0p}) \in \mathbb{R}^q$ (potential state cell)

Long Short Term Memory Networks

The Information Flow

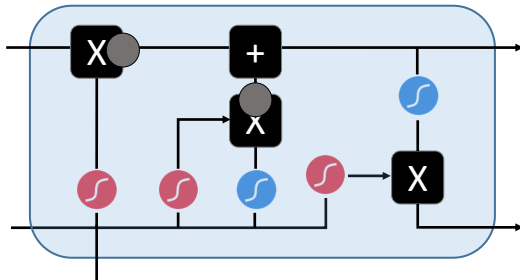


► At time t :

1. $c_{t-1} \odot f_t$
2. $i_t \odot p_t$

Long Short Term Memory Networks

The Information Flow

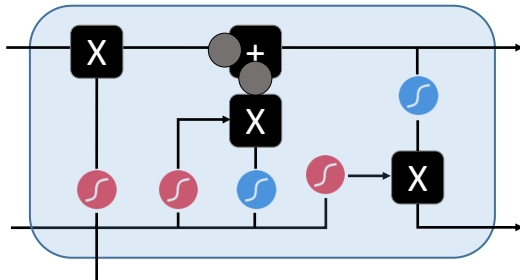


► At time t :

1. $c_{t-1} \odot f_t$
2. $i_t \odot p_t$

Long Short Term Memory Networks

The Information Flow

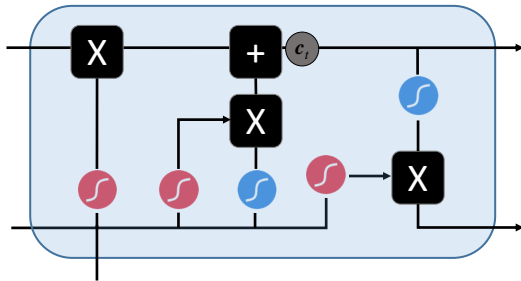


► At time t :

1. $c_{t-1} \odot f_t$
2. $i_t \odot p_t$

Long Short Term Memory Networks

The Information Flow

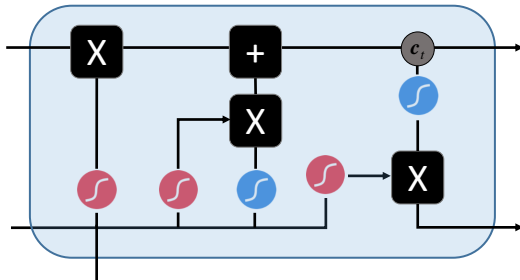


► At time t :

1. $c_t = c_{t-1} \odot f_t + i_t \odot p_t$

Long Short Term Memory Networks

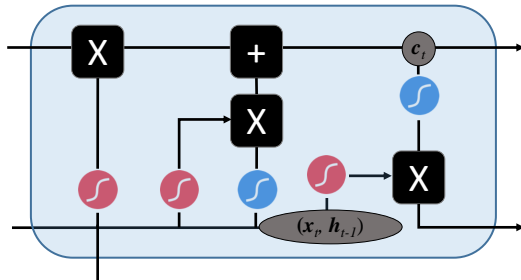
The Information Flow



► At time t :

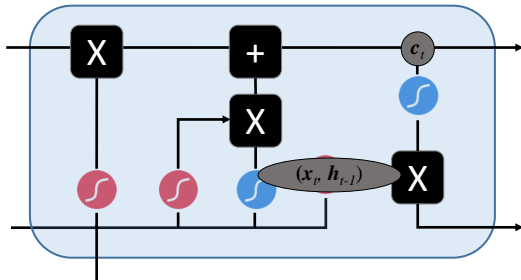
1. $c_t = c_{t-1} \odot f_t + i_t \odot p_t$

The Information Flow



- At time t :
 1. $\mathbf{c}_t = \mathbf{c}_{t-1} \odot \mathbf{f}_t + \mathbf{i}_t \odot \mathbf{p}_t$
 2. Logistic activation: $\mathbf{o}_t = \text{logistic}(\mathbf{\Gamma}'_o \mathbf{z}_t + \gamma_{0o}) \in \mathbb{R}^q$

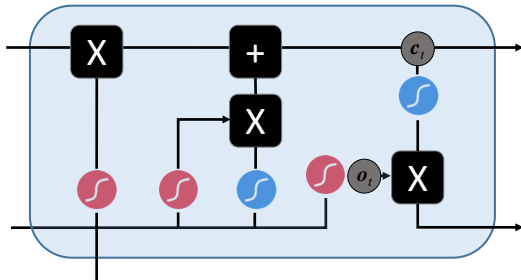
The Information Flow



- At time t :
 1. $\mathbf{c}_t = \mathbf{c}_{t-1} \odot \mathbf{f}_t + \mathbf{i}_t \odot \mathbf{p}_t$
 2. Logistic activation: $\mathbf{o}_t = \text{logistic}(\mathbf{\Gamma}'_o \mathbf{z}_t + \gamma_{0o}) \in \mathbb{R}^q$

Long Short Term Memory Networks

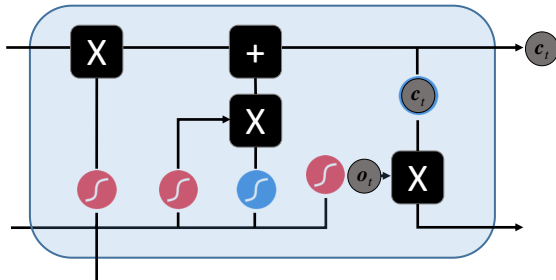
The Information Flow



► At time t :

1. $c_t = c_{t-1} \odot f_t + i_t \odot p_t$
2. Logistic activation: $o_t = \text{logistic}(\Gamma'_o z_t + \gamma_{0o}) \in \mathbb{R}^q$

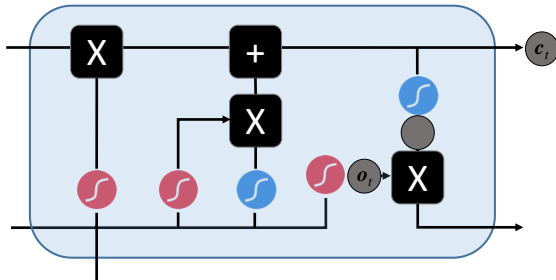
The Information Flow



- At time t :
 1. $\mathbf{c}_t = \mathbf{c}_{t-1} \odot \mathbf{f}_t + \mathbf{i}_t \odot \mathbf{p}_t$
 2. Logistic activation: $\mathbf{o}_t = \text{logistic}(\mathbf{\Gamma}'_o \mathbf{z}_t + \gamma_{0o}) \in \mathbb{R}^q$

Long Short Term Memory Networks

The Information Flow

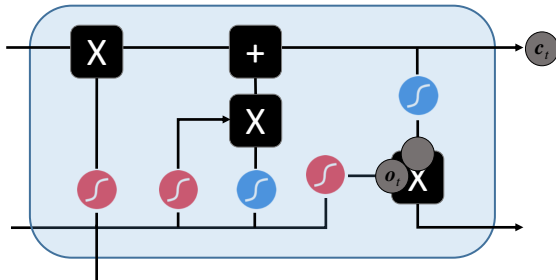


► At time t :

1. $c_t = c_{t-1} \odot f_t + i_t \odot p_t$
2. Logistic activation: $o_t = \text{logistic}(\Gamma'_o z_t + \gamma_{0o}) \in \mathbb{R}^q$

Long Short Term Memory Networks

The Information Flow

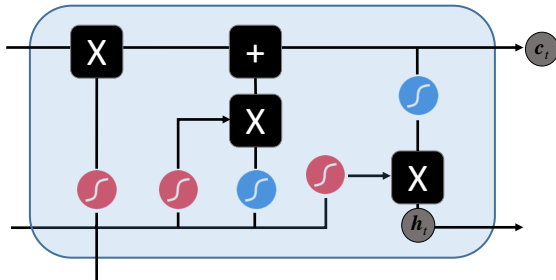


► At time t :

1. $h_t = o_t \odot \tanh(c_t)$

Long Short Term Memory Networks

The Information Flow

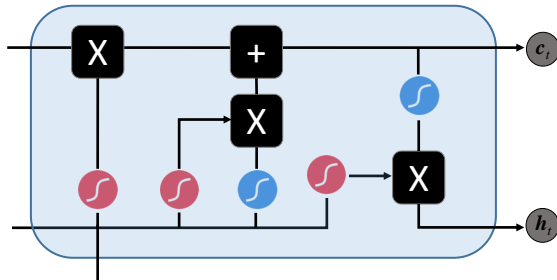


► At time t :

1. $h_t = o_t \odot \tanh(c_t)$

Long Short Term Memory Networks

The Information Flow



► At time t :

1. $h_t = o_t \odot \tanh(c_t)$
 $y_t = \gamma_y' h_t + \gamma_{0y}$

LSTM Network

Formal (Mathematical) representation

- ▶ Initiate with $\mathbf{c}_0 = 0$ and $\mathbf{h}_0 = 0$.

LSTM Network

Formal (Mathematical) representation

- ▶ Initiate with $\mathbf{c}_0 = 0$ and $\mathbf{h}_0 = 0$.
- ▶ Given input \mathbf{x}_t do for $t \in \{1, \dots, T\}$:

LSTM Network

Formal (Mathematical) representation

- ▶ Initiate with $\mathbf{c}_0 = 0$ and $\mathbf{h}_0 = 0$.
- ▶ Given input \mathbf{x}_t do for $t \in \{1, \dots, T\}$:

$$\mathbf{f}_t = \text{Logistic}(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f)$$

LSTM Network

Formal (Mathematical) representation

- ▶ Initiate with $\mathbf{c}_0 = 0$ and $\mathbf{h}_0 = 0$.
- ▶ Given input \mathbf{x}_t do for $t \in \{1, \dots, T\}$:

$$\mathbf{f}_t = \text{Logistic}(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f)$$

$$\mathbf{i}_t = \text{Logistic}(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i)$$

LSTM Network

Formal (Mathematical) representation

- ▶ Initiate with $\mathbf{c}_0 = 0$ and $\mathbf{h}_0 = 0$.
- ▶ Given input \mathbf{x}_t do for $t \in \{1, \dots, T\}$:

$$\mathbf{f}_t = \text{Logistic}(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f)$$

$$\mathbf{i}_t = \text{Logistic}(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i)$$

$$\mathbf{o}_t = \text{Logistic}(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o)$$

LSTM Network

Formal (Mathematical) representation

- ▶ Initiate with $\mathbf{c}_0 = 0$ and $\mathbf{h}_0 = 0$.
- ▶ Given input \mathbf{x}_t do for $t \in \{1, \dots, T\}$:

$$\mathbf{f}_t = \text{Logistic}(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f)$$

$$\mathbf{i}_t = \text{Logistic}(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i)$$

$$\mathbf{o}_t = \text{Logistic}(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o)$$

$$\mathbf{p}_t = \text{Tanh}(\mathbf{W}_c \mathbf{x}_t + \mathbf{U}_c \mathbf{h}_{t-1} + \mathbf{b}_c)$$

LSTM Network

Formal (Mathematical) representation

- ▶ Initiate with $\mathbf{c}_0 = 0$ and $\mathbf{h}_0 = 0$.
- ▶ Given input \mathbf{x}_t do for $t \in \{1, \dots, T\}$:

$$\mathbf{f}_t = \text{Logistic}(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f)$$

$$\mathbf{i}_t = \text{Logistic}(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i)$$

$$\mathbf{o}_t = \text{Logistic}(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o)$$

$$\mathbf{p}_t = \text{Tanh}(\mathbf{W}_c \mathbf{x}_t + \mathbf{U}_c \mathbf{h}_{t-1} + \mathbf{b}_c)$$

$$\mathbf{c}_t = (\mathbf{f}_t \odot \mathbf{c}_{t-1}) + (\mathbf{i}_t \odot \mathbf{p}_t)$$

LSTM Network

Formal (Mathematical) representation

- ▶ Initiate with $\mathbf{c}_0 = 0$ and $\mathbf{h}_0 = 0$.
- ▶ Given input \mathbf{x}_t do for $t \in \{1, \dots, T\}$:

$$\mathbf{f}_t = \text{Logistic}(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f)$$

$$\mathbf{i}_t = \text{Logistic}(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i)$$

$$\mathbf{o}_t = \text{Logistic}(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o)$$

$$\mathbf{p}_t = \text{Tanh}(\mathbf{W}_c \mathbf{x}_t + \mathbf{U}_c \mathbf{h}_{t-1} + \mathbf{b}_c)$$

$$\mathbf{c}_t = (\mathbf{f}_t \odot \mathbf{c}_{t-1}) + (\mathbf{i}_t \odot \mathbf{p}_t)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \text{Tanh}(\mathbf{c}_t)$$

LSTM Network

Formal (Mathematical) representation

- ▶ Initiate with $\mathbf{c}_0 = 0$ and $\mathbf{h}_0 = 0$.
- ▶ Given input \mathbf{x}_t do for $t \in \{1, \dots, T\}$:

$$\mathbf{f}_t = \text{Logistic}(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f)$$

$$\mathbf{i}_t = \text{Logistic}(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i)$$

$$\mathbf{o}_t = \text{Logistic}(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o)$$

$$\mathbf{p}_t = \text{Tanh}(\mathbf{W}_c \mathbf{x}_t + \mathbf{U}_c \mathbf{h}_{t-1} + \mathbf{b}_c)$$

$$\mathbf{c}_t = (\mathbf{f}_t \odot \mathbf{c}_{t-1}) + (\mathbf{i}_t \odot \mathbf{p}_t)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot \text{Tanh}(\mathbf{c}_t)$$

$$\mathbf{y}_t = \mathbf{W}_y \mathbf{h}_t + \mathbf{b}_y$$

Obrigado