

robot doggo can learn how to walk???

An introductory survey of reinforcement learning and its applications to  
robotic systems.

Sky Hong

April 28, 2025

## **Abstract**

something idk, figure out final formatting later

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Machine Learning</b>	<b>4</b>
2.1	Introduction to Machine Learning . . . . .	4
2.1.1	Parameters . . . . .	4
2.1.2	Loss Functions . . . . .	5
2.1.3	Optimization . . . . .	5
2.2	Neural Networks . . . . .	7
2.2.1	Activation Functions and the Universal Approximation Theorems . .	7
2.2.2	Layers . . . . .	8
2.3	Backpropagation . . . . .	9
2.4	Overfitting . . . . .	9
2.4.1	Cross-validation . . . . .	10
2.4.2	Regularization . . . . .	10
2.4.3	Dropout . . . . .	10
2.5	Techniques for Improving Optimization with Gradient Descent . . . . .	10
2.5.1	Data Batching . . . . .	10
2.5.2	Stochastic Gradient Descent . . . . .	10
2.5.3	Momentum . . . . .	10
2.5.4	Adaptive Learning Rate . . . . .	10
2.5.5	Learning Rate Scheduling . . . . .	10
2.5.6	Feature Scaling and Batch Normalization . . . . .	10
2.6	Classification . . . . .	10
2.6.1	Softmax . . . . .	10
2.6.2	Cross Entropy Loss . . . . .	10
2.7	Reinforcement Learning . . . . .	10
2.7.1	Policy Gradient . . . . .	10
2.7.2	Actor-Critic . . . . .	10
2.7.3	Reward Shaping . . . . .	10
2.7.4	Inverse Reinforcement Learning . . . . .	10
<b>3</b>	<b>Robotic Systems and Control</b>	<b>11</b>
3.1	Hutter, Lee, Hwangbo et al. and the ANYmal by ANYbotics . . . . .	11
3.2	Luo et al. and Exoskeletons . . . . .	11

3.3 more examples will go here . . . . .	11
<b>Glossary</b>	<b>12</b>
<b>Bibliography</b>	<b>13</b>

# Chapter 1

## Introduction

# Chapter 2

## Machine Learning

### 2.1 Introduction to Machine Learning

The field of machine learning is fundamentally about finding functions that model data. For example, a speech recognition model takes audio as an input and outputs a text, and AlphaZero takes a chessboard state as an input and outputs a move. These functions may be vastly complex, with far too many parameters and relations for any human to reasonably articulate and program, even if the task comes naturally to our evolved biology. Machine learning offers methods for computers to achieve these tasks, without the need for a human to give explicit instructions on *how* it should be accomplished.

This section will give a very introductory overview of machine learning, the specific details of which will be delved more deeply in following sections. The general framework of creating a machine learning model is as follows:

1. Identification of the properties of a model. What should this model be able to do? What are its inputs and outputs? What kind of model architecture is best suited to its problem?
2. Defining a loss function, some metric to measure how well a model performs.
3. Optimizing the parameters of a model to minimize the loss function and maximize performance.

#### 2.1.1 Parameters

As an example, consider a simple model consisting of a linear relationship between a data set and the output of the model.

$$f(\mathbf{x}) = b + \sum_i w_i x_i \tag{2.1}$$

Parameters that directly multiply values,  $w_i$ , are called **weights**, and parameters that offset

values,  $b$ , are called **biases**. The inputs from the data set,  $\mathbf{x}$ , are called **features**. This can be generalized to vector outputs with matrices, and in more complex models, the notation may be extended to tensors.

$$f = \mathbf{b} + \mathbf{W}\mathbf{x} \quad (2.2)$$

The parameters of a model are usually collectively referred to as a vector  $\boldsymbol{\theta}$ , with components of the individual weights and biases of the model. For complex models,  $\boldsymbol{\theta}$  may be millions, billions, or even trillions of parameters long, populated by numerous parameter tensors.

$$\boldsymbol{\theta} = \begin{bmatrix} W_{11} \\ W_{12} \\ \vdots \\ b_1 \\ b_2 \\ \vdots \end{bmatrix} \quad (2.3)$$

For virtually all non-trivial problems, linear relationships are far too reductive to completely capture the complexity of a problem. Such inherent limitations of a model due to its architecture known as **model bias**. Model architectures will be expounded upon in more detail in Section 1.2.

### 2.1.2 Loss Functions

A **loss function** (also sometimes called a **cost function**), typically denoted  $L(\boldsymbol{\theta})$  or just  $L$ , is the measure of how “bad” a set of parameters  $\boldsymbol{\theta}$  for a model is. A common definition is the deviation between a model’s prediction and an actual result. For example, when building a speech recognition model, the loss may be defined as the error rate between its output transcription and the correct transcription. The features of the training data that are provided to the model, in this case, the correct transcriptions, are identified with **labels** that tell the model what it should train towards.

Losses over all labels in the training data are aggregated into an value for the overall loss function, the most common method of which is the **mean square error (MSE)**:

$$L = \frac{1}{N} \sum (y - \hat{y})^2. \quad (2.4)$$

$L$  is associated with an **error surface** that can be understood as the plot of  $L$  in a  $|\boldsymbol{\theta}|$ -dimensional parameter space, though it is usually far too complex to be directly visualized.

### 2.1.3 Optimization

Since  $L$  is continuous, there must exist some set of parameters  $\boldsymbol{\theta}^*$  that minimize  $L$  and model the training data best.

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L \quad (2.5)$$

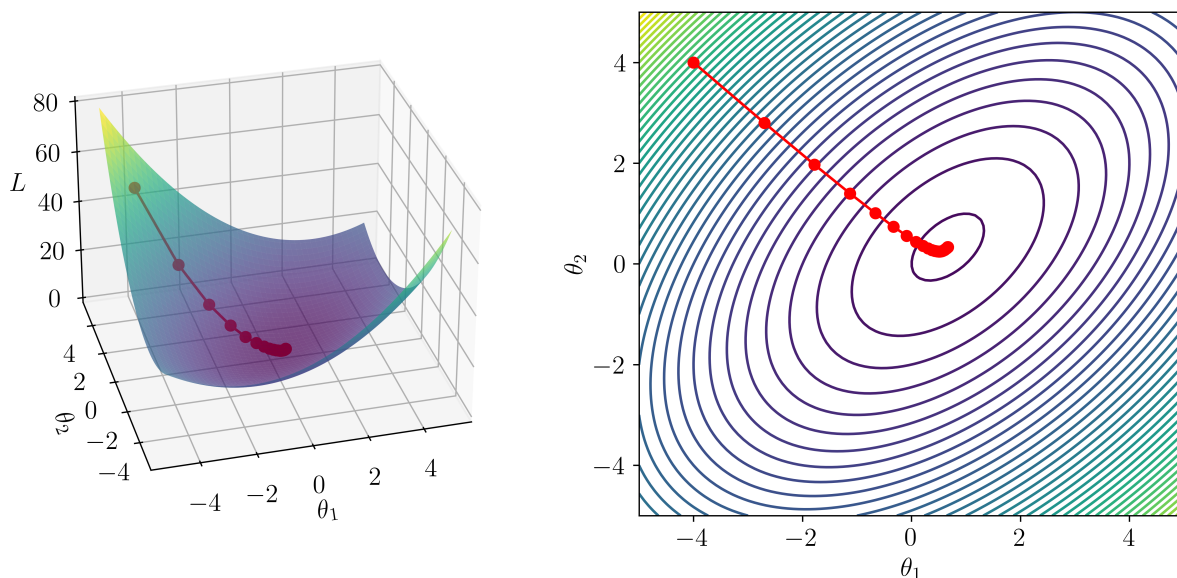
The process of improving the model by finding values for  $\theta$  that lower the loss function is what it really means for a model to “learn” or “train”.

## Gradient descent

Since brute-forcing the error surface of  $L$  is infeasible with a large number of parameters, the standard algorithmic approach to find minima of  $L$  is **gradient descent**. In gradient descent,  $L$  is iteratively lowered by stepping  $\theta$  against the gradient of  $L$ .

$$\theta^{(t+1)} = \theta^{(t)} - \eta \nabla L(\theta^{(t)}) \quad (2.6)$$

The **learning rate**,  $\eta$ , determines the scaling for the size of each step.



**Figure 2.1:** Visualization of naive gradient descent on an error surface with two parameters ( $L = \theta_1^2 + \theta_2^2 - \theta_1\theta_2 - \theta_1$ ) by a 3D plot (left), and a contour plot (right). The red line shows 100 iterations of gradient descent, starting from initial parameters  $\theta^{(0)} = [-4, 4]$ , with  $\eta = 0.1$ . (Original)

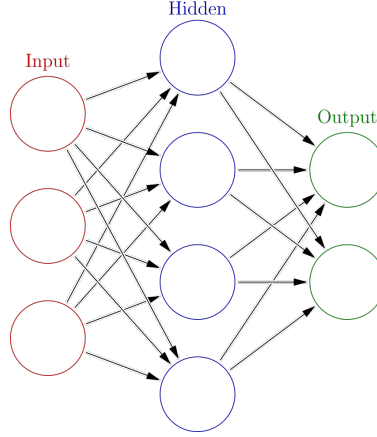
In practice, some conditions may be defined to determine when to cease training, such as setting an upper limit on the number of iterations, setting a time limit, or setting a threshold for a satisfactory value of  $L$ . These values that modify aspects of a model’s learning are called **hyperparameters**. Hyperparameter optimization is itself a complex topic that is beyond the scope of this section.

When  $\theta$  reaches a critical point ( $\nabla L = 0$ ), the updates to  $\theta$  vanish, and the algorithm has reached a stable ending point. In an ideal case, this would be the global minima of  $L$ , the best that the model could possibly reach given the training data. In reality, reaching the global minima is highly improbable,  $\theta$  is much more likely to be stuck at a local minima, or a saddle point. We will discuss methods to overcome these challenges in following sections.



## 2.2 Neural Networks

**Artificial neural networks (ANNs)** often shortened to **neural networks (NNs)**, are one of the largest classes of models used in machine learning. NNs are composed of **nodes** that are arranged in layers. The first layer is called the **input layer**, and the last layer is called the **output layer**. All layers in between are called **hidden layers**. NNs are usually fully connected, meaning that each node in a layer is considered by every node in the next layer.



**Figure 2.2:** A schematic of a neural network with a input layer with three nodes, one hidden layer with four nodes, and an output layer with two nodes [1].

Each node takes a biased weighted sum of the values of the previous layer, applies an activation function  $\phi$ , and outputs the value to the next layer.

$$a_j^{(l)} = \phi \left( b_j^{(l)} + \sum_i w_{ij}^{(l-1)} a_i^{(l-1)} \right) \quad (2.7)$$

Or, using matrix notation, where we define  $\phi$  to act component-wise on a vector,

$$\mathbf{a}^{(l)} = \phi \left( \mathbf{b}^{(l)} + \mathbf{W}^{(l-1)} \mathbf{a}^{(l-1)} \right). \quad (2.8)$$

In this article, we will use the notation  $a_i^{(l)}$  to refer to the output of node  $i$  in layer  $l$ , and  $\mathbf{a}^{(l)}$  to refer to the vector of outputs of all nodes in layer  $l$ . The weights and biases applied onto nodes from layer  $l$  to be inputted into layer  $l + 1$  will be denoted as  $\mathbf{W}^{(l)}$  and  $\mathbf{b}^{(l)}$ , respectively. Note that this superscript notation is different from the one we use with  $\boldsymbol{\theta}^{(t)}$ , which is used to denote the iteration of the gradient descent algorithm.

### 2.2.1 Activation Functions and the Universal Approximation Theorems

For virtually all non-trivial problems, linear relationships are far too reductive to completely capture the complexity of the task at hand. Therefore, a variety of non-linear **activation**

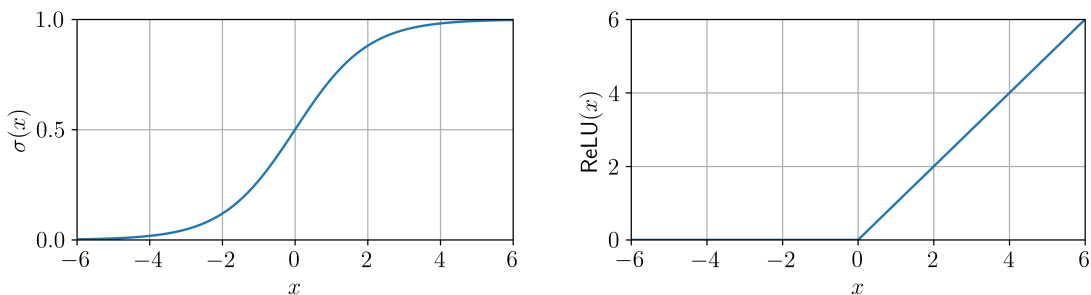
**functions** are applied at each node to eliminate this kind of model bias to allow the model to learn more non-linear relationships.

**Sigmoid functions** are ubiquitous in this role, and they are also useful since they normalize the output of a node to be between 0 and 1. This is useful for many applications, such as when the output of a node is interpreted as a probability.

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (2.9)$$

Another very common activation function is the **rectified linear unit (ReLU)**, which has the advantage of being computationally easier to calculate.

$$\text{ReLU}(x) = \max(0, x) \quad (2.10)$$



**Figure 2.3:** Plot of  $\sigma(x)$  (left) and  $\text{ReLU}(x)$  (right). (Original)

Any continuous relationship between two variables can be approximated by a *sufficiently large* linear combination of sigmoid functions [2].

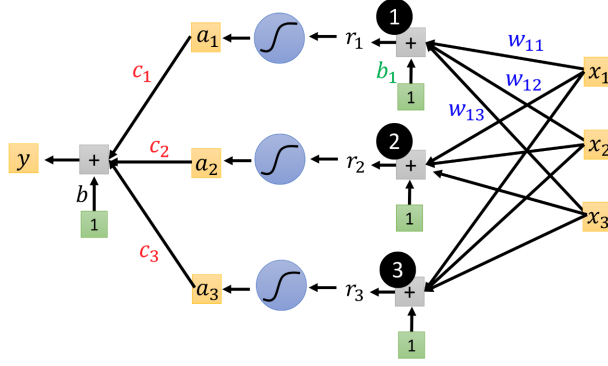
$$f \approx b + \mathbf{c} \cdot \sigma(\mathbf{b} + \mathbf{W}\mathbf{x}). \quad (2.11)$$

Equation 2.11 is equivalent to a neural network with one hidden layer and scalar output (omitting the final application of an activation function).

This result can be generalized to any non-polynomial activation function, including ReLUs, in the **universal approximation theorems**, implying that NNs that use such activation functions may theoretically learn any relationship given *sufficiently large* amounts nodes and layers [4, 5, 6]. It is important to note that this result is only theoretical, and in practice, computational and data constraints usually frustrate this ideal.

## 2.2.2 Layers

The addition of layers to a model increases its complexity, the number of which is referred to as its **depth**. The number of nodes in a layer is referred to as the layer's **width**. The width



**Figure 2.4:** Reprinted from [3].

and depth of a model are hyperparameters, and there is a balance to be made between model bias, computational demands, and overfitting, a phenomenon we will discuss in Section 2.4.

## 2.3 Backpropagation

**Backpropagation** is the textbook algorithm for computing the gradient of the loss function. Essentially, it is the application of the chain rule to neural networks.

$$\nabla L = \frac{\partial L}{\partial \theta_i} = \frac{\partial L}{\partial a_j} \frac{\partial a_j}{\partial z_j} \frac{\partial z_j}{\partial \theta_i} \quad (2.12)$$

## 2.4 Overfitting

**Overfitting** is a phenomenon that occurs when a model learns the training data too well, and is unable to generalize to new data. Generally, models with more parameters and degrees of freedom are more prone to overfitting. This may be mitigated by simply using a smaller model or by stopping training earlier, but these methods are not always feasible.

2.4.1 Cross-validation

2.4.2 Regularization

2.4.3 Dropout

## 2.5 Techniques for Improving Optimization with Gradient Descent

2.5.1 Data Batching

2.5.2 Stochastic Gradient Descent

2.5.3 Momentum

Nesterov Momentum

2.5.4 Adaptive Learning Rate

Adagrad

RMSprop

Adam

2.5.5 Learning Rate Scheduling

2.5.6 Feature Scaling and Batch Normalization

## 2.6 Classification

2.6.1 Softmax

2.6.2 Cross Entropy Loss

## 2.7 Reinforcement Learning

2.7.1 Policy Gradient

2.7.2 Actor-Critic

2.7.3 Reward Shaping

2.7.4 Inverse Reinforcement Learning

# Chapter 3

## Robotic Systems and Control

**3.1** Hutter, Lee, Hwangbo et al. and the ANYmal by ANYbotics

[7]

**3.2** Luo et al. and Exoskeletons

**3.3** more examples will go here

# Glossary

**ANN** artificial neural network. 7

**bias** Parameters that add/offset to values in a model. Adds to a value add a node in neural networks. Not to be confused with *model bias*. 5

**error surface** The surface that a loss function plots in parameter space. 5

**feature** . 5

**hyperparameter** Values that specify aspects of a model’s inherent architecture or its learning. 6

**label** The identifiers attached to training data that function as an “answer key” for what it is supposed to be. 5

**loss function** Function defined to measure the performance of a model, lower is better. 5

**model bias** Inherent limitations of a model’s performance to its architecture. Not to be confused with *bias*. 5

**MSE** mean square error. 5

**NN** neural network. 7

**node** . 7

**ReLU** rectified linear unit. 8

**weight** Parameters that are direct coefficients to values in a model. Determines strength of connections between nodes in neural networks. 4

# Bibliography

- [1] Wikimedia Commons. *File:Colored neural network.svg* — *Wikimedia Commons, the free media repository*. 2025. URL: [https://commons.wikimedia.org/w/index.php?title=File:Colored\\_neural\\_network.svg&oldid=995727191](https://commons.wikimedia.org/w/index.php?title=File:Colored_neural_network.svg&oldid=995727191).
- [2] G. Cybenko. “Approximation by superpositions of a sigmoidal function”. In: *Mathematics of Control, Signals and Systems* 2.4 (Dec. 1, 1989), pp. 303–314. ISSN: 1435-568X. DOI: 10.1007/BF02551274. URL: <https://doi.org/10.1007/BF02551274>.
- [3] Hung-Yi Lee. “Introduction of Machine / Deep Learning”. 2021. URL: [https://speech.ee.ntu.edu.tw/~hylee/ml/ml2021-course-data/regression%20\(v16\).pdf](https://speech.ee.ntu.edu.tw/~hylee/ml/ml2021-course-data/regression%20(v16).pdf).
- [4] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. “Multilayer feedforward networks are universal approximators”. In: *Neural Networks* 2.5 (1989), pp. 359–366. ISSN: 0893-6080. DOI: [https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8). URL: <https://www.sciencedirect.com/science/article/pii/0893608089900208>.
- [5] Patrick Kidger and Terry Lyons. “Universal Approximation with Deep Narrow Networks”. In: *Proceedings of Thirty Third Conference on Learning Theory*. Ed. by Jacob Abernethy and Shivani Agarwal. Vol. 125. Proceedings of Machine Learning Research. PMLR, July 9, 2020, pp. 2306–2327. URL: <https://proceedings.mlr.press/v125/kidger20a.html>.
- [6] Allan Pinkus. “Approximation theory of the MLP model in neural networks”. In: *Acta Numerica* 8 (1999), pp. 143–195. DOI: 10.1017/S0962492900002919.
- [7] Marco Hutter. “Legged Robots on the way from subterranean”. 2022 IEEE International Conference on Robotics and Automation (ICRA), May 24, 2022. URL: [https://www.youtube.com/watch?v=XwheB2\\_dyMQ](https://www.youtube.com/watch?v=XwheB2_dyMQ).