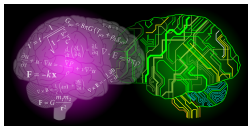


ML and Physics

Lec I: Introduction to ML



Ghadir Jafari

University of Cape Town, Cape Town, South Africa

Farhangian University, Tehran, Iran

The Objectives

- ① Motivate to learn ML
- ② Physical intuition for ML concepts
- ③ Introduce some research lines
- ④ Physics and ML connections

Outline

- ① An Introduction to ML for Physicist
- ② ML for Physics
- ③ Physics for ML
- ④ Physics and ML

Where is ML?

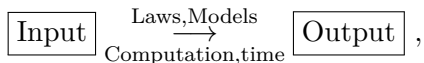
- Search engines, Spam filter in email, Recommender systems, Social networks
- Computer Vision, Speech recognition,
- Medical diagnosis, Medicine
- Industry, Economics,
- Theorem proving, Climate science,
- Self Driving Cars
- ...

Technology and ML

- Advances in computational hardware, the widespread use of GPUs
- Moore's law
- “Big Data” revolution
- Developments in theoretical aspects and programming
- Artificial neural networks

Physics And ML

- **Information:** Physics and Machine Learning
- Gathering and analyzing Data
- Finding patterns, Design models
- Predict the behaviour of (complex) systems
- Learning: passing information from input to outputs



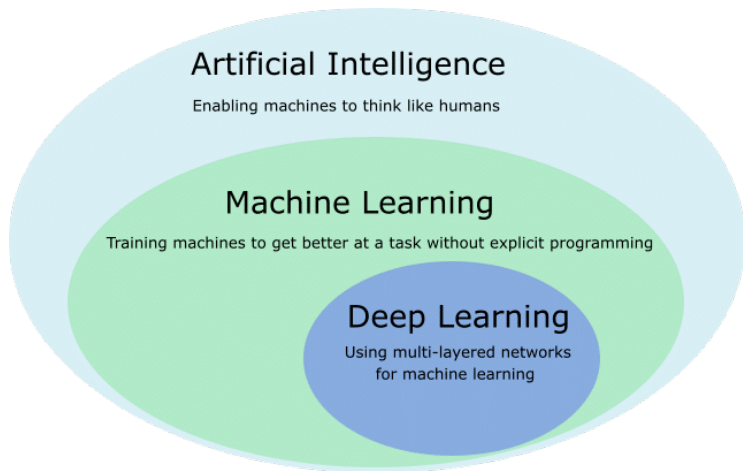
Physics Vs Machine Learning

- Physicists using their own knowledge, intelligence and intuition to inform their models
- Machine Learning: models are agnostic and the machine provides the intelligence by extracting it from data.

Machine Learning

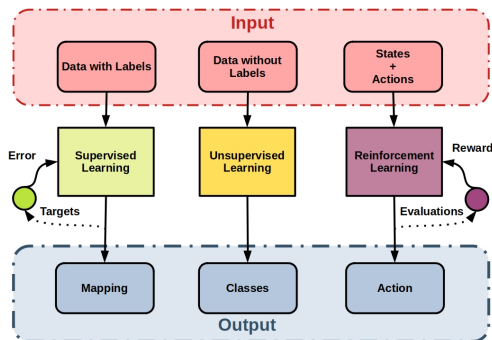
- Machine Learning is the field of study that gives computers the ability to learn without being explicitly programmed.
- Task easy for humans to do, but difficult for humans to describe in terms of elementary operations.
- Problems for which existing solutions require a lot of fine-tuning or long lists of rules
- Complex problems for which using a traditional approach yields no good solution
- Getting insights about complex problems and large amounts of data

AI, ML and DL



Approaches to ML

- Supervised learning: Data inputs and outputs (labels)
- Unsupervised learning: No labels, Clustering, find structure in its input
- Reinforcement learning: dynamic input space, maximize a reward function



Main Elements In ML

- ① Data Preparation
- ② Model Selection,
- ③ Training, Optimization
- ④ Validation, Test

From Data to Machines

- A Collection of Pairs (x_i, y_i)

$$\text{Data} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

- x and y could be anything
- The goal: find (make) a function (machine) f

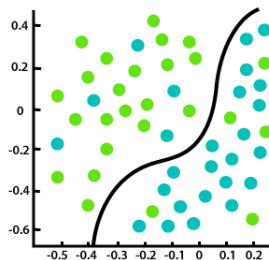
$$x_i \longrightarrow f(x_i) = y_i$$

- Such that for new data $(x_{\text{new}}, y_{\text{new}})$

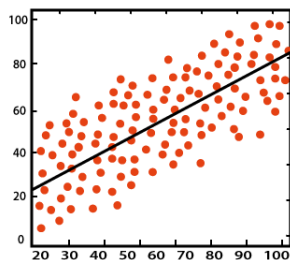
$$y_{\text{new}} = f(x_{\text{new}})$$

Classification Vs Regression

- Classification: the y take discrete values (y : labels)
- Regression: y continuous



Classification



Regression

- Classification \rightarrow Nonlinear regression

Modeling: Function Representation (Approximation)

- Polynomial Representation

$$f_a(x) \approx \sum_{n=1}^N a_n x^n$$

- Fourier Representation:

$$f_c(x) \approx \sum_{n=-N}^N c_n e^{inx}$$

- e.g. $y^i \sim \sum_j^N w^{ij} x_j + b_j$
- Any Representation $\Rightarrow f_\theta(x)$
- A set of parameters which parameterize the space of functions

$$\theta = \{\theta_1, \theta_1, \dots, \theta_N\}$$

Training (Optimization)

- Finding suitable $\{\theta_i^*\}$ such that

$$y_i \approx f_{\theta^*}(x_i)$$

- Define a cost $\mathcal{C}_i(f_\theta, y_i) \Rightarrow \mathcal{C}_i \approx 0$ when $y_i \approx f(x_i)$

$$\text{i.e. } \mathcal{C}_i = (y_i - f(x_i))^2$$

- The Cost Function

$$\mathcal{C}(\theta) = \frac{1}{N} \sum_i^N \mathcal{C}_i$$

- Want: $\theta^* \rightarrow \mathcal{C}(\theta^*) \approx 0$
- \Rightarrow Minimizing $\mathcal{C}(\theta)$ in function space

Gradient Descent Method

- Change in \mathcal{C} by $d\theta$:

$$d\mathcal{C}(\theta) = d\theta \cdot \nabla_{\theta}\mathcal{C},$$

- If

$$d\theta = -\eta\nabla_{\theta}\mathcal{C}$$

- Then:

$$d\mathcal{C}(\theta) = -\eta|\nabla_{\theta}\mathcal{C}|^2 \leq 0$$

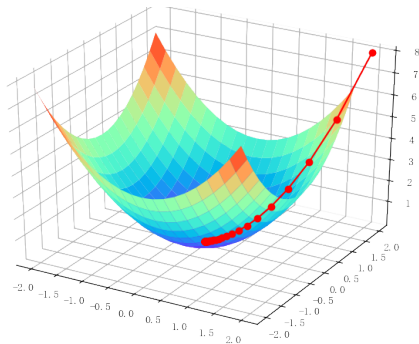
- Steps toward the local minimum:

$$\theta \rightarrow \theta - \eta\nabla_{\theta}\mathcal{C}$$

$$\mathcal{C} \rightarrow \mathcal{C} - \eta|\nabla_{\theta}\mathcal{C}|^2$$

Gradient Descent Method

- η : Learning rate, n Number of steps (epochs)



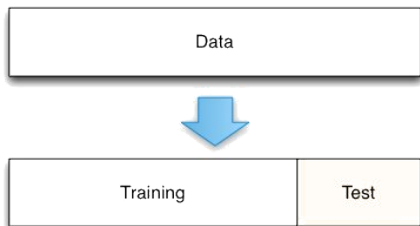
- Particle moving in a potential

Testing

- New Data $\{x_{\text{new}}, y_{\text{new}}\}$

$$y_{\text{new}} = f(x_{\text{new}})?$$

- Training and Test Split:



Example: Linear Regression

- data: $\{x_i, y_i\}$
- objective: predict $y = f(x)$
- model: $f(x; \vec{\theta}) = mx + b$, $\vec{\theta} = (m, b)$
- loss function:

$$\mathcal{C}(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2$$

- model training: optimal parameters

$$\hat{\vec{\theta}} = \arg \min \mathcal{C}(\vec{\theta})$$

Example: Learning The Digits

- Problem: A program to recognize the images of handwritten digits
- Input:

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

- Output: Digit letters $\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$
- Impossible algorithmic programming
- Machine Learning: computers learn (doing tasks) without being explicitly programmed

Example: Learning The Digits

- The Data:

{**0** → 0, **0** → 0, **0** → 0, **0** → 0, **0** → 0, **0** → 0, **0** → 0, **0** → 0, **0** → 0, **0** → 0},
{**1** → 1, **1** → 1, **1** → 1, **1** → 1, **1** → 1, **1** → 1, **1** → 1, **1** → 1, **1** → 1, **1** → 1},
{**2** → 2, **2** → 2, **2** → 2, **2** → 2, **2** → 2, **2** → 2, **2** → 2, **2** → 2, **2** → 2, **2** → 2},
{**3** → 3, **3** → 3, **3** → 3, **3** → 3, **3** → 3, **3** → 3, **3** → 3, **3** → 3, **3** → 3, **3** → 3},
{**4** → 4, **4** → 4, **4** → 4, **4** → 4, **4** → 4, **4** → 4, **4** → 4, **4** → 4, **4** → 4, **4** → 4},
{**5** → 5, **5** → 5, **5** → 5, **5** → 5, **5** → 5, **5** → 5, **5** → 5, **5** → 5, **5** → 5, **5** → 5},
{**6** → 6, **6** → 6, **6** → 6, **6** → 6, **6** → 6, **6** → 6, **6** → 6, **6** → 6, **6** → 6, **6** → 6},
{**7** → 7, **7** → 7, **7** → 7, **7** → 7, **7** → 7, **7** → 7, **7** → 7, **7** → 7, **7** → 7, **7** → 7},
{**8** → 8, **8** → 8, **8** → 8, **8** → 8, **8** → 8, **8** → 8, **8** → 8, **8** → 8, **8** → 8, **8** → 8},
{**9** → 9, **9** → 9, **9** → 9, **9** → 9, **9** → 9, **9** → 9, **9** → 9, **9** → 9, **9** → 9, **9** → 9}

Example: Learning The Digits

- Each image a matrix:

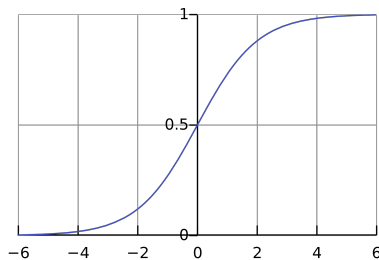
```
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 0 0 0 1 1 0 0 0 0 0 1 1 1 1 1
1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
```

- Outputs are labels (discrete values)

20×20 Matrices $([0, 1]^{400}) \rightarrow \{1, 2, \dots, 10\}$

Example: Learning The Digits

- How to model?
- Discrete output: classification problem
- A function to $[0, 1]^{10}$
- A linear model $f(x_n) = w^{in}x_n + b^i$
- Using the logistic sigmoid function : $\sigma(x) = \frac{1}{1+e^{-x}}$



- The model $f(x_n) = \sigma(w^{in}x_n + b^i)$

Example: Learning The Digits

- Optimization: Find a minimum for

$$\mathcal{C}(w, b) = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2$$

- No much precision!
- Is the model simple?
- How to generalize?
- $x_n \rightarrow w^{\text{in}}x_n + b^i \rightarrow \sigma(w^{\text{in}}x_n + b^i)$
- The number of parameters : $10 * 400 + 10$

Neural Networks

$$z_i^{(1)} \equiv b_i^{(1)} + \sum_{j=1}^{n_0} W_{ij}^{(1)} x_j, \quad \text{for } i = 1, \dots, n_1$$

$$z_i^{(1)} \rightarrow \sigma(z_i^{(1)})$$

$$z_i^{(2)} \equiv b_i^{(2)} + \sum_{j=1}^{n_1} W_{ij}^{(2)} \sigma(z_j^{(1)}), \quad \text{for } i = 1, \dots, n_2$$

⋮

$$z_i^{(\ell)}(x) \equiv b_i^{(\ell)} + \sum_{j=1}^{n_{\ell-1}} W_{ij}^{(\ell)} \sigma(z_j^{(\ell-1)}), \quad \text{for } i = 1, \dots, n_\ell$$

⋮

for $\ell = 1, \dots, L$

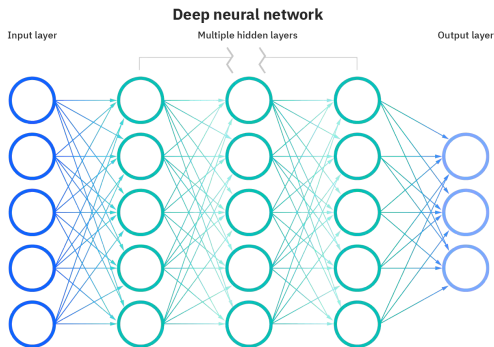
$$z_i^{(L)} \rightarrow \sigma(z_i^{(L)})$$

$$\Rightarrow y = f_{\{w,b\}}(x) = \sigma(z_i^{(L)})$$

Neural Networks

- A neural network is a composite nonlinear function :

$$y = f_L(w^L f_{L-1}(w^{L-1} \dots f_1(w^1 x) \dots))$$



- Basic Ingredients: Neurons
- Parameters $\{w_{i,j}^\ell, b_i^\ell\}$

Learning The Digits

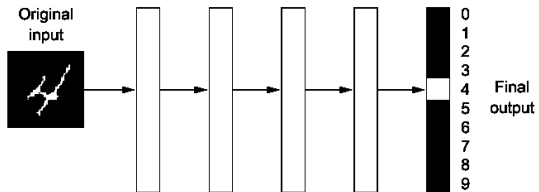
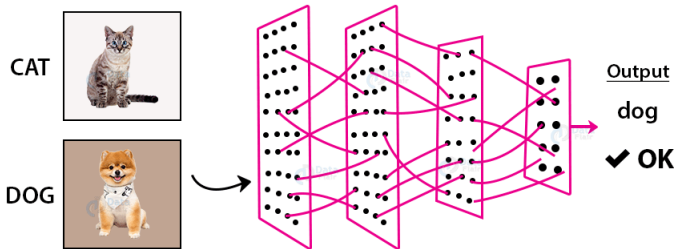
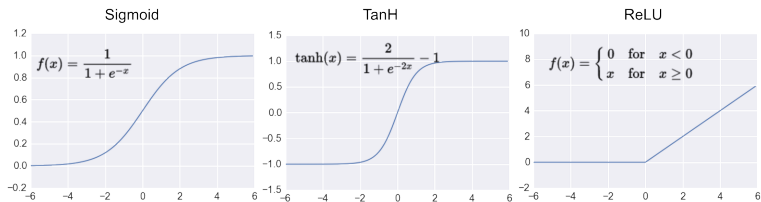


Image Recognition



Neural Networks Architecture

- **The Hyper-parameters:** number of layers L , The number of neurons in each layer n_ℓ , Initialization of Parameters, Learning Rate η , The Number of Epochs, ...
- The activation function $\sigma(x)$ e.g.:



NN as Universal Approximation

- Neural Networks can represent any functions
- For every y, d, D and ϵ

$$y : \mathbb{R}^d \rightarrow \mathbb{R}^D$$

- There exist a network $f_{w,b}$ such that

$$\|y(x) - f_{w,b}(x)\| < \epsilon$$

- \Rightarrow Neural network: A flexible representation for high-dim functions
- Importance: Every process is a function computation

The backpropagation Algorithm

- Training: the calculation of $\nabla_{\theta} \mathcal{C}$ at each step

$$\mathcal{C}(\vec{\theta}) = \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i))^2$$

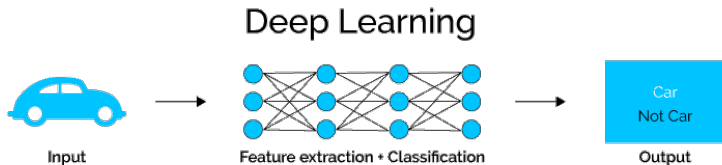
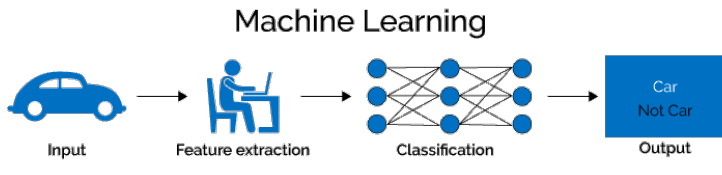
$$\vec{\theta} = \{w_{i,j}^{\ell}, b_i^{\ell}\}$$

- How to compute for NN?
- Suppose $w_{i,j}^{\ell} \rightarrow w_{i,j}^{\ell} + \delta w_{i,j}^{\ell}$
- Chain Rule:

$$\nabla_{w^{\ell}} \mathcal{C} = [(f^{\ell})' \circ (w^{\ell+1})^T \dots \circ (w^{L-1})^T \cdot (f^{L-1})' \circ (w^L)^T \cdot (f^L)' \circ \nabla_{z^{\ell}} \mathcal{C}] (z^{\ell-1})^T$$

Deep Learning

- Large number of layers in NN
- \Rightarrow More efficient



The Black Box?

- No Theoretical explanations of deep learning
- Why the output $f_{\theta^*}(x)$ is given on an input x
- Hyper-Parameters: trial-and-error
- The activation functions, the widths of the layers and the depth
- which learning problems are computationally tractable?
- \Rightarrow Black box



A Theory for Deep Learning?

- Complexity \Rightarrow Statistical explanation

$$P(\mathbf{x}|f_{\theta}(\mathbf{x})) \quad \text{Or/And} \quad P(f_{\theta}(\mathbf{x})|\mathbf{x})$$

- Statistical Theory beneath experimental thermodynamic rules
- Like statistical mechanics explains how the macroscopic laws of thermodynamics describing the steam engines.
- \Rightarrow A theory for machines

Unsupervised Learning

- No outputs (labels)

$$\mathcal{D} = \{\mathbf{x}_i\}$$

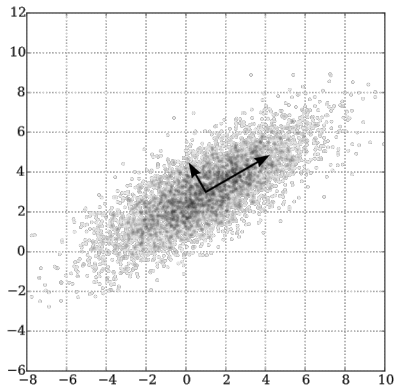
- Depend on the specific details of the task
- We assume that there exists some probability distribution $P(\mathbf{x})$
- The model $Q_{\theta}(\mathbf{x})$
- Make them similar
- Reduce the relative entropy

$$D_{\text{KL}}(P \parallel Q) = \sum_{\mathbf{x} \in \mathcal{D}} P(\mathbf{x}) \log \left(\frac{P(\mathbf{x})}{Q_{\theta}(\mathbf{x})} \right).$$

- Both modeling and optimization \Rightarrow Deep Learning help

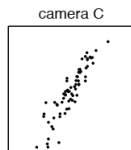
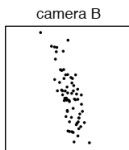
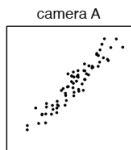
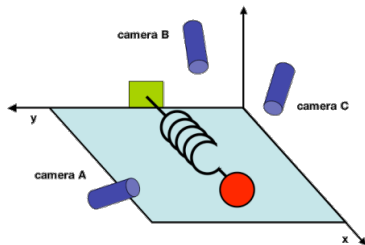
Principal component analysis

- Compress the data to loss information as less as possible
- Dimension Reduction



Simple Example¹

- Data from an oscillating spring



¹[arXiv:1404.1100](https://arxiv.org/abs/1404.1100)

Which Problems could be the subject of ML?

- ML: Direct relation of input and outputs,
- Data, Modeling, Optimization
- \Rightarrow In Principle any problem

Some Textbooks

- Pattern recognition and machine learning , Bishop, C M (2006),
- Deep Learning, Ian Goodfellow, Yoshua Bengio, Aaron Courville ¹
- Neural Networks and Deep Learning, Michael Nielsen ²

¹www.deeplearningbook.org

²neuralnetworksanddeeplearning.com

Programming

- Most Common languages: Python and R
- Libraries: Tensorflow, Scikit-learn, Keras, PyTorch
- Prerequisites: NumPy, pandas, and Matplotlib
- Hands-On Machine Learning with Scikit-Learn and TensorFlow, Aurélien Geron¹
- Wolfram Mathematica

¹github.com/ageron/handson-ml2

Summery

- ML deal with modeling and predicting similar to physics
- Main elements include **modeling** and **optimization**
- Neural networks provide a flexible modeling
- With the cost of being black box