

Application of Optimization to Machine learning

Look at artificial neural networks (ANNs)

Aim: N measurements

$$\begin{matrix} \underline{x}_1 \\ \vdots \\ \underline{x}_N \end{matrix} \quad \left\{ \text{Inputs} \right. \quad \begin{matrix} \underline{y}_1 \\ \vdots \\ \underline{y}_N \end{matrix} \quad \left\{ \text{Outputs} \right.$$

Unknown functional relationship between inputs and outputs. Aim of ML is to "learn" the functional relationship.

What this means: A general input \underline{x} , a general output \underline{y} . Approximate the functional relationship as

$$\underline{y} \approx f(\underline{x}, \underline{\rho})$$

where the ρ 's are adjustable parameters and f is a function to be determined.

NEW NOTATION FOR THIS PART!

- ρ for parameters
- f for function TBC (not cost function!)

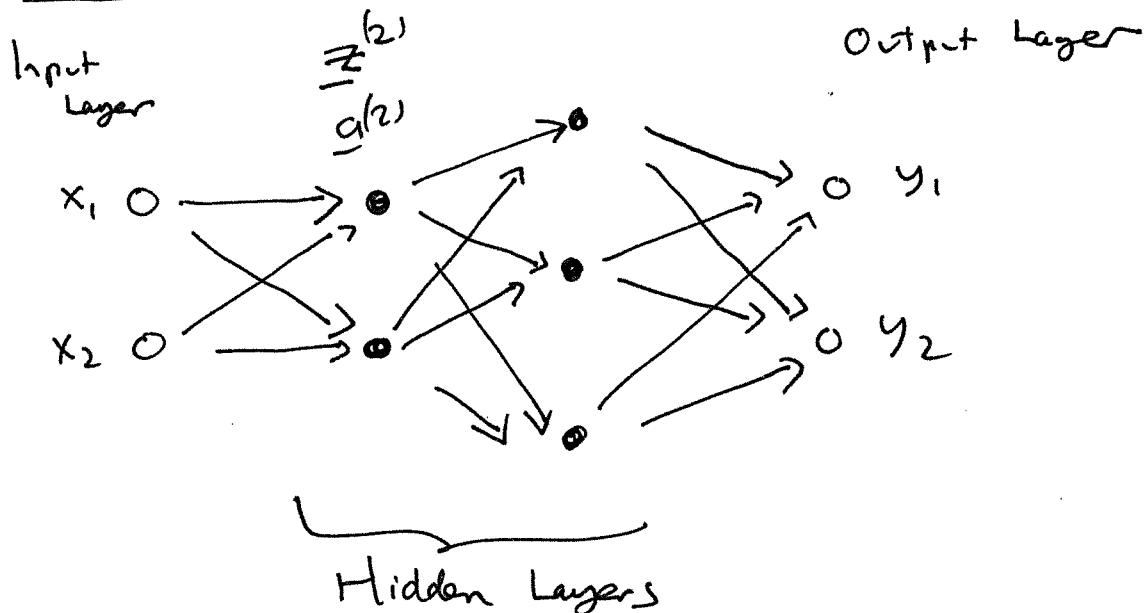
In Artificial Neural Networks (ANNs) $f(\underline{x}, \underline{\rho})$ is approximated by a composition of:

- o Activation functions
- o Linear functions — described by their weights and biases.

$$\underbrace{\hspace{1cm}}_{\underline{\rho}}$$

The idea is inspired by how the human brain works → a network of nodes which are either on or off — neurons or ~~connected~~ and connections between nodes — synapses. The first ANN was created by Frank Rosenblatt in 1957 and was implemented on an IBM 704 machine. He called his ANN the "perceptron".

Example: A small neural network:



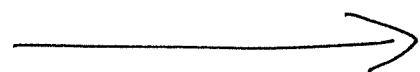
- The input layer: $\underline{a}^{(1)} = \underline{x}$
- Gets transformed according to

$$\underline{\Sigma}^{(2)} = W^{(2)} \underline{a}^{(1)} + \underline{b}^{(2)}$$
- Gets activated according to:

$$\underline{a}^{(2)} = \sigma(\underline{\Sigma}^{(2)})$$
 POINTWISE
- Gets transformed:

$$\underline{\Sigma}^{(3)} = W^{(3)} \underline{a}^{(2)} + \underline{b}^{(3)}$$
- Gets activated:

$$\underline{a}^{(3)} = \sigma(\underline{\Sigma}^{(3)})$$



- Gets transformed to the output layer: III

$$\underline{z}^{(4)} = W^{(4)} \underline{a}^{(3)} + \underline{b}^{(4)}$$

- Gets activated:

$$\underline{a}^{(4)} = \sigma(\underline{z}^{(4)}) \simeq \underline{y}$$

Write as a composition of functions:

$$\begin{aligned} \underline{y} = f(\underline{x}) &\simeq \sigma(W^{(4)} \underline{a}^{(3)} + \underline{b}^{(4)}) \\ &= \sigma\left[W^{(4)}\left(W^{(3)} \underline{a}^{(2)} + \underline{b}^{(3)}\right) + \underline{b}^{(4)}\right] \\ &= \sigma\left\{W^{(4)}\sigma\left[W^{(3)}\sigma\left(W^{(2)} \underline{a}^{(1)} + \underline{b}^{(2)}\right) + \underline{b}^{(3)}\right] + \underline{b}^{(4)}\right\} \\ &= \sigma\left\{W^{(4)}\sigma\left[W^{(3)}\sigma\left(W^{(2)} \underline{x} + \underline{b}^{(2)}\right) + \underline{b}^{(3)}\right] + \underline{b}^{(4)}\right\} \end{aligned}$$

Parameters: $P = \{W^{(1)}, W^{(2)}, W^{(3)}, W^{(4)}, b^{(1)}, b^{(2)}, b^{(3)}, b^{(4)}\}$

Counting:

- $\underline{z}^{(n)}$
 $\underline{z}^{(2)} \in \mathbb{R}^2, \underline{z}^{(3)} \in \mathbb{R}^3, \underline{z}^{(4)} \in \mathbb{R}^2$
- $W^{(n)}$
 $W^{(2)} \in \mathbb{R}^{2 \times 3}, W^{(3)} \in \mathbb{R}^{3 \times 2}, W^{(4)} \in \mathbb{R}^{2 \times 2}$ 16
- $\underline{b}^{(n)}$
 $\underline{b}^{(2)} \in \mathbb{R}^2, \underline{b}^{(3)} \in \mathbb{R}^3, \underline{b}^{(4)} \in \mathbb{R}^2$ 7

23 parameters.

IV

Training data: $\begin{cases} x_1, \dots, x_N \\ y_1, \dots, y_N \end{cases}$

Cost function ("loss function"):

$$C(\rho) = \sum_{i=1}^N \frac{1}{2} \| f(x_i, \rho) - y_i \|_2^2$$
$$= \sum_{i=1}^N C_i(\rho)$$

Steepest Descent:

$$\rho^{I+1} = \rho^I - \eta \nabla C(\rho^I)$$

Problem is computational cost of computing gradient (23 dimensions $\times N$ = # observations)

Instead, we stochastic gradient descent:

- At iteration I , choose $i \in \{1, 2, \dots, N\}$ at random.
- Update ρ^I according to:

$$\rho^{I+1} = \rho^I - \eta \nabla_{\rho} C_i(\rho^I)$$

Back-propagation : Nature of approximating $f \approx$

(compositions) means that the gradients $\nabla_p C_i$ can be computed ANALYTICALLY, via back-propagation.

For counter = 1 up to Niter

Stochastic gradient descent

Choose an integer k uniformly at random from $\{1, 2, \dots, N\}$

$\underline{x}^{(k)}$ is set as the current training point.

Set $\underline{a}^{(1)} = \underline{x}^{(k)}$.

Feed forward step :

For $l = 2$ up to L # l and L label layers

$$\underline{z}^{(l)} = W^{(l)} \underline{a}^{(l-1)} + b^{(l)}$$

$$\underline{a}^{(l)} = \sigma(\underline{z}^{(l)})$$

$$D^{(l)} = \text{diag}(\sigma'(\underline{z}^{(l)}))$$

end

Back propagation steps :

Initialize with $\delta^{(L)} = \frac{\partial C}{\partial z^{(L)}}$:

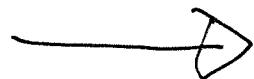
$$\delta^{(L)} = D^{(L)} (\underline{a}^{(L)} - \underline{y}^{(k)})$$

For $l = L-1$ down to 2

$$\delta^{(l)} = D^{(l)} (W^{(l+1)})^T \delta^{(l+1)}$$

end

Here, $\delta^{(l)} = \frac{\partial C}{\partial z^{(l)}}$.



Steepest-descent update

For $\ell = L$ down to 2

$$W^{(\ell)} \rightarrow W^{(\ell)} - \eta \delta^{(\ell)} (a^{(\ell-1)})^T$$

Here: $\frac{\partial C}{\partial W^{(\ell)}} = \delta^{(\ell)} (a^{(\ell-1)})^T$,

outer product

Also, $\frac{\partial C}{\partial b^{(\ell)}} = \delta^{(\ell)}$:

$$b^{(\ell)} \rightarrow b^{(\ell)} - \eta \delta^{(\ell)}$$

end

end

Example :

$$\left\{ \begin{array}{l} x_1, \dots, x_N = \text{points in } \mathbb{R}^2 \\ y_1, \dots, y_n \in [0,1] \end{array} \right.$$

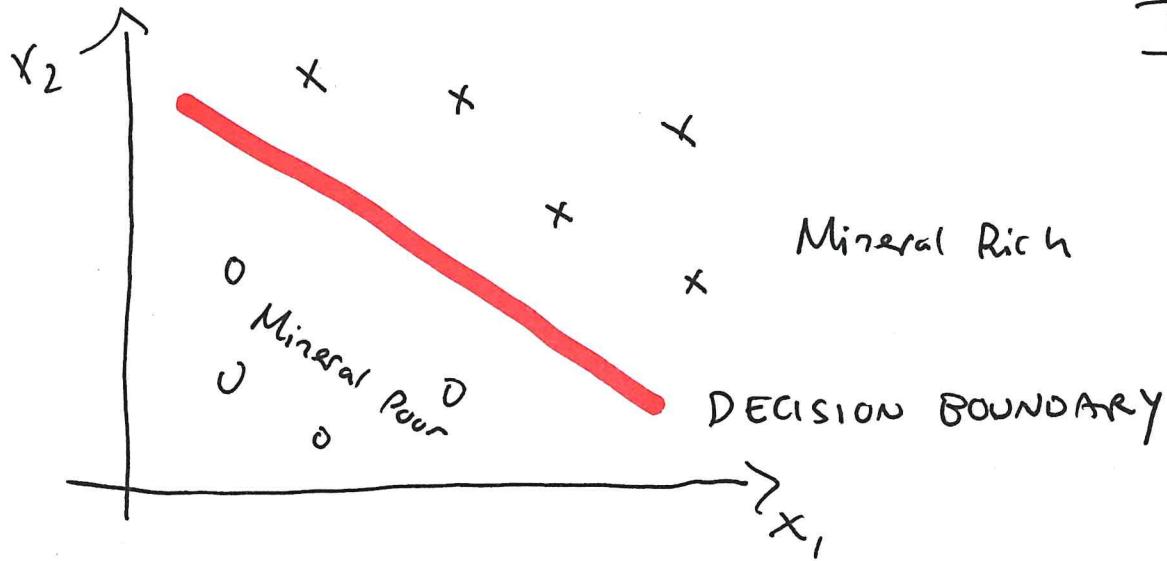
Could be samples from various places where:

- x_i are latitude / longitude

- $y_i = \begin{cases} 1 & \text{if } x_i \text{ is a mineral-rich region} \\ 0 & \text{if } x_i \text{ is mineral poor} \end{cases}$

We are looking for the decision boundary

separating mineral-poor and mineral-rich regions. This then tells us where to explore for minerals. 



Instead of writing our own code, we have used MATLAB's `feedforwardnet` function to create an ANN with two hidden layers :

- First hidden layer has 2 nodes
- Second " " " 3 nodes

Structure of input / output layers set by training data :

- Input layer has 2 nodes (x_1, x_2)
- Output layer has 1 node (y)

- o ANN is supplied with synthetic training data
- o ANN is trained using `train` function in Matlab
- o Result is plotted and decision boundary is constructed as the 0.5 contour of $f(x, \theta)$.

LINK TO MATLAB CODE

DEMO ON LAPTOP

Summary

VIII

- Whistling tour of Continuous optimization
- Described problem of training an artificial neural network in mathematical terms
- Showed how an ANN can be used to construct a decision boundary.

Key problems in ANN

- Non-convex optimization problem — not clear if global min. is reached
- Understanding best network architecture (# nodes, # layers) for the problem in hand
- Understanding the best choice of activation function for each layer
- Understanding which cost function is best for the problem in hand
- Requires good knowledge of Mathematics to be a responsible user of ANNs.

