





## Machines, please behave!

Understanding the behavior of machine learning algorithms

Matan Abraham





## Supervised learning.

- Function that maps inputs to outputs based on 'learning' from labelled training data of examples of input-output pairs.
- This is called 'training' the model.
- Common examples:
  - Regression
  - Decision tress
  - Neural networks







## Unsupervised learning.

- No explicit target outputs or environmental evaluations associated with inputs.
- Patterns are modeled using measures of similarity defined by metrics such as Euclidean and probabilistic distance.
- Common examples:
  - Clustering
  - Anomaly detection







### Neural networks.



### Neural network function:

Inputs: 3 Outputs: 1 (the prediction) Parameters: 16 weights/biases

### "the algorithm"

### Cost function:

Inputs: 16 Outputs: 1 (the cost) Parameters: All our training examples

## "the measurement of performance"



## Deep learning.



Each node in the hidden layer has an activation function.

It is not easily understood what is happening in hidden layer.

This layer defines the complexity of the neural network.



f(y)

f(y) = 0

Inputs: 10 Outputs: 4 Parameters: 159





## Neural network intuition.

### Hidden layer





## Applications of neural networks.







Image analysis and interpretation

Signal analysis and interpretation

Drug development

## Neural network – Insight health-risk predictor.





# You don't know what you don't know.

"Black box" machine learning limits visibility into the decisionmaking process.

Studying the algorithms ensures:

- Confidence predictions are as accurate as expected to be.
- Results mean what users think they mean.
- Predictions don't rely on unwanted information.



"Information theory of deep learning"

"Explainable AI"



## Aligning with reality.

Training requires a lot of data.

Training on a particular dataset which is applied to make predictions on future data – the scope of which cannot be predicted. "Transfer learning"

"Lean data learning"

"Synthesizing new data through simulations"

## Baked-in bias.

Algorithms can learn biases that are undesirable as predictions.

There need to be extra constraints imposed on an algorithm beyond just "accuracy".

This becomes more challenging the more complex the algorithm.



## False assumptions.

It is imperative that the assumptions made in designing training the algorithms are appropriate.

Examples:

- Missing data
- Data collection
- Distributions
- Parameterization

# The self-fulfilling prophecy.

"If the algorithm tries to exploit what it learned, without leaving any room for exploration, it will keep reinforcing what it already knows and it will not learn new things, eventually becoming useless."







## Artificial Intelligence vs. Human Intelligence









