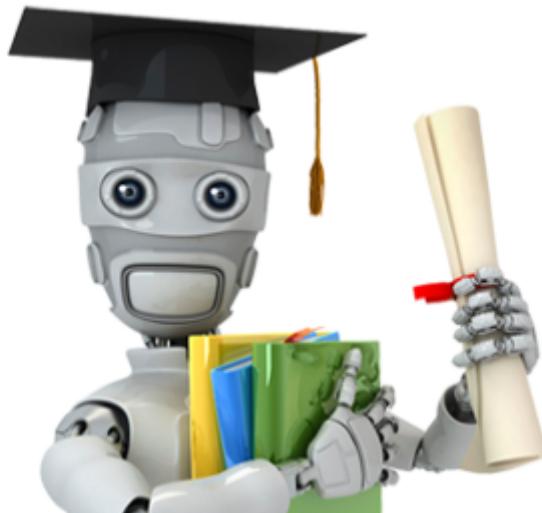


Learning Machine Le...

Learning Machine Learning

There are many good free online college level classes for Machine line. The course offered by Stanford Professor, Andrew Ng is the one I took from Coursera. I am continue suppliment my learning from other online classes and YouTube video and blog. This note is to share the learning experience.



Instructors:

Andrew Ng

Associate Professor, Stanford University; Chief Scientist, Baidu; Chairman and Co-founder, Coursera

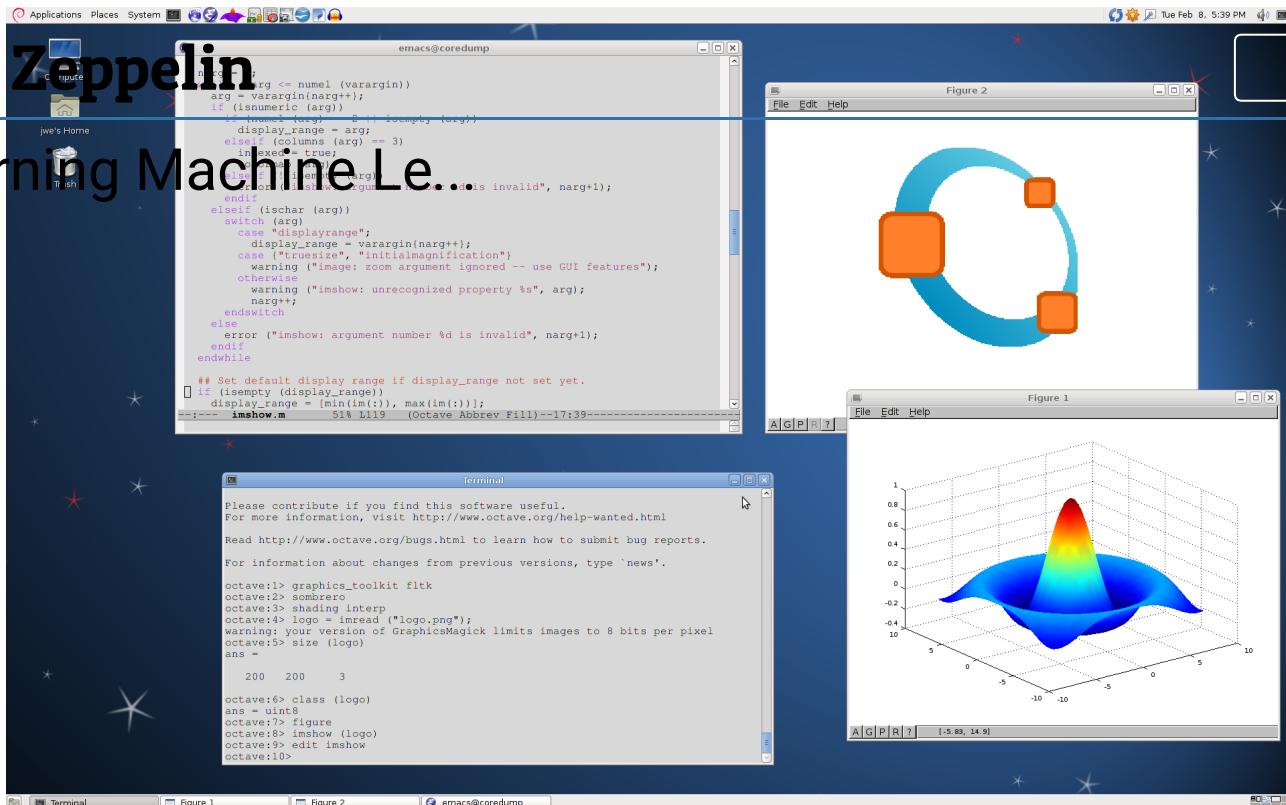
Courses

Course Contents (<https://www.coursera.org/learn/machine-learning>)

Pre-requisites

(They will be reviewed in class):

- Linear Algebra (<https://www.khanacademy.org/math/linear-algebra>)
- Octave (http://wiki.octave.org/Video_tutorials)



History of Machine Learning - Winter 1

Winter 1: Machine Translation

1960s

1954

Georgetown experiment demonstrates successful translation of 60 Russian sentences into English

701 Translator

IBM Press release, January 8, 1954

New York, January 7..... Russian was translated into English by an electronic "brain" today for the first time.
Brief statements about politics, law, mathematics, chemistry, metallurgy, communications and politics again were submitted in Russian by linguists of the Georgetown University Institute of Languages and Linguistics to the famous 701 computer of the International Business Machines Corporation. And the giant computer, within a few seconds, turned the sentences into easily readable English.

A girl who didn't understand a word of the language of the Soviets punched out the Russian messages on IBM cards. The "brain" dashed off its English translations on an automatic printer at the breakneck speed of two and a half lines per second.

1966

Automated Language Processing Advisory Committee says progress is slow, kills machine translation funding for a decade

English: "The spirit is willing, but the flesh is weak."

After English>Russian, Russian>English translation:
"The whisky is strong, but the meat is rotten."

History of Machine Learning - Winter 2

Winter 2: Golden Age

1970s

1959

- “Reasoning as search”— Newell & Simon Geometry Theorem prover, Arthur Samuels checkers player

1968-1970

- Micro-worlds—SHRDLU Terry Winograd, MIT



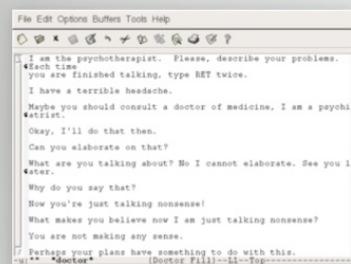
1974

- Whitehill report:

“...utter failure of AI to achieve its grandiose objectives”

1964-1968

- Pattern matching—ELIZA Joseph Weizenbaum, MIT



1970

- Marvin Minsky (in LIFE magazine)

“In from three to eight years we will have a machine with the general intelligence of an average human being.”

History of Machine Learning - Winter 3

Winter 3: Expert Systems

1980s

1965

- Edward Feigenbaum, Carl Djerassi and others write DENDRAL to identify chemical compounds from spectrometer data



1980

- John McDermott's XCON (Expert Configurator) automatically selected VAX minicomputer components based on user needs

1987

- Collapse of Symbolics



1972

- Edward Shortliffe at Stanford writes MYCIN to diagnose infectious blood diseases

Startups

- Symbolics Lisp Machines, IntelliCorp, Aion, Thinking Machines

1990

- IBM: The Integrated Reasoning Shell

HD

New Period of Machine Learning

2012

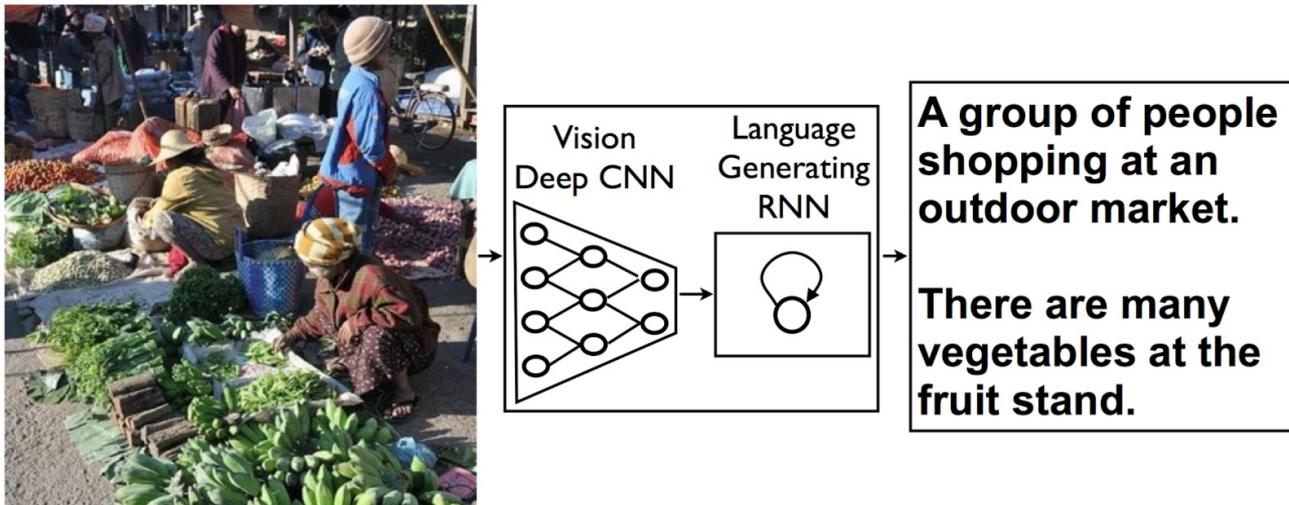
Andrew Ng and his team made major breakthrough in face detector from unlabeled images: from the article Building High-level Features Using Large Scale Unsupervised Learning ([http://www robotics stanford.edu/~ang/papers/icml12-HighLevelFeaturesUsingUnsupervisedLearning.pdf](http://www robotics stanford edu/~ang/papers/icml12-HighLevelFeaturesUsingUnsupervisedLearning.pdf))

We train a 9-layered locally connected sparse autoencoder with pooling and local contrast normalization on a large dataset of images (the model has 1 billion connections, the dataset has 10 million 200x200 pixel images downloaded from the Internet). We train this network using model parallelism and asynchronous SGD on a cluster with 1,000 machines (16,000 cores) for three days.

Contrary to what appears to be a widely-held intuition, our experimental results reveal that it is possible to train a face detector without having to label images as containing a face or not. Control experiments show that this feature detector is robust not only to translation but also to scaling and out-of-plane rotation. We also find that the same network is sensitive to other high-level concepts such as cat faces and human bodies. Starting with these learned features, we trained our network to obtain 15.8% accuracy in recognizing 20,000 object categories from ImageNet, a leap of 70% relative improvement over the previous state-of-the-art.

2015

Stanford University worked on Neural Network which can explain unlabeled photo. For example:



Source: Google

What is Machine Learning?

Machine Learning is concerned with the development, the analysis, and the application of algorithms that allow computers to learn.

Learning:

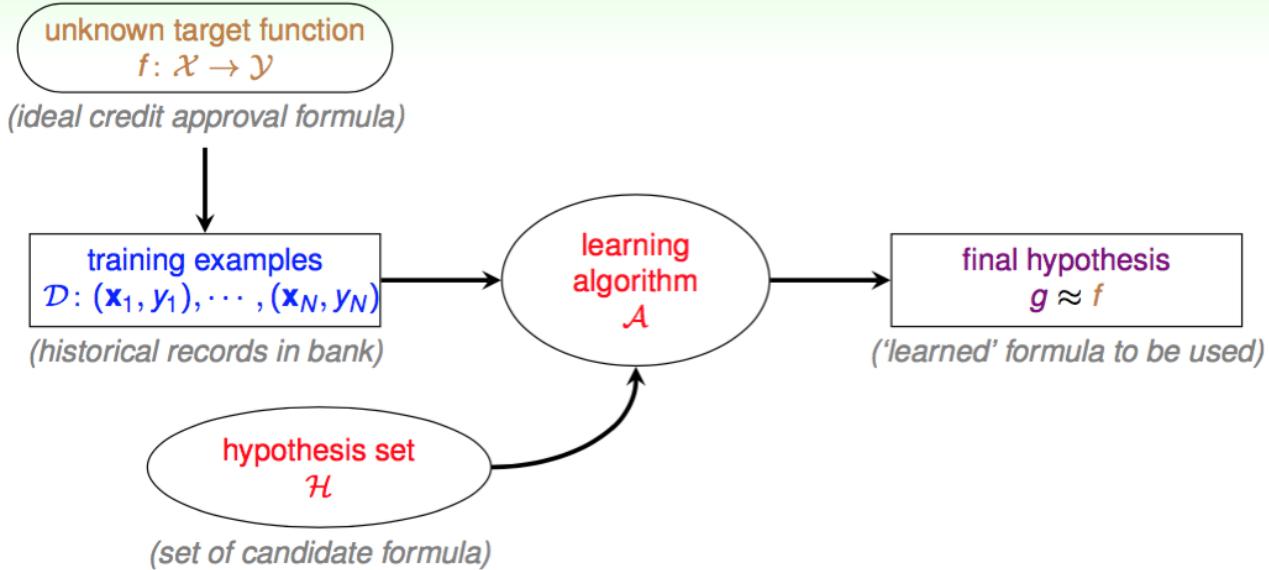
- A computer program is said to learn from experience (E) with some class of tasks (T) and a performance measure (P) if its performance at tasks in T as measured by P improves with E. (i.e. by collecting data)
- Extracting a model of a system from the sole observation (or the simulation) of this system in some situations.
 - A model = some relationships between the variables used to describe the system.
- Two main goals: make prediction and better understand the system.

Definition of Machine Learning

The Learning Problem

Components of Machine Learning

Practical Definition of Machine Learning



machine learning:
use data to compute hypothesis g
that approximates target f

Components of Machine Learning

problem: unknown target function, f

- to find out the pattern for approving the credit card that benefit to a bank. A target function f , which maps applicant X (information about different application) that leads to outcome of Y (different outcomes).

training examples, D

- input: information of each applicant, x : age, salary, exist debts, etc
- output: outcome of each applicant, y : good or bad for bank/late payment/default
- collected data, D : $\{(x_1, y_1), (x_2, y_2), \dots (x_n, y_n)\}$

hypothesis set, H

- There is a set of h in H , we like to find a specific h , good skill, hopefully have good performance. We select the best h , we call it g
- function g , is part of $H = \{h_k\}$, that can map $X \rightarrow Y$ with good accuracy

learning algorithm, A

- Use data to compute the best hypothesis, g , which approximates to f

target function, g

- will be used to forecast future applicants.

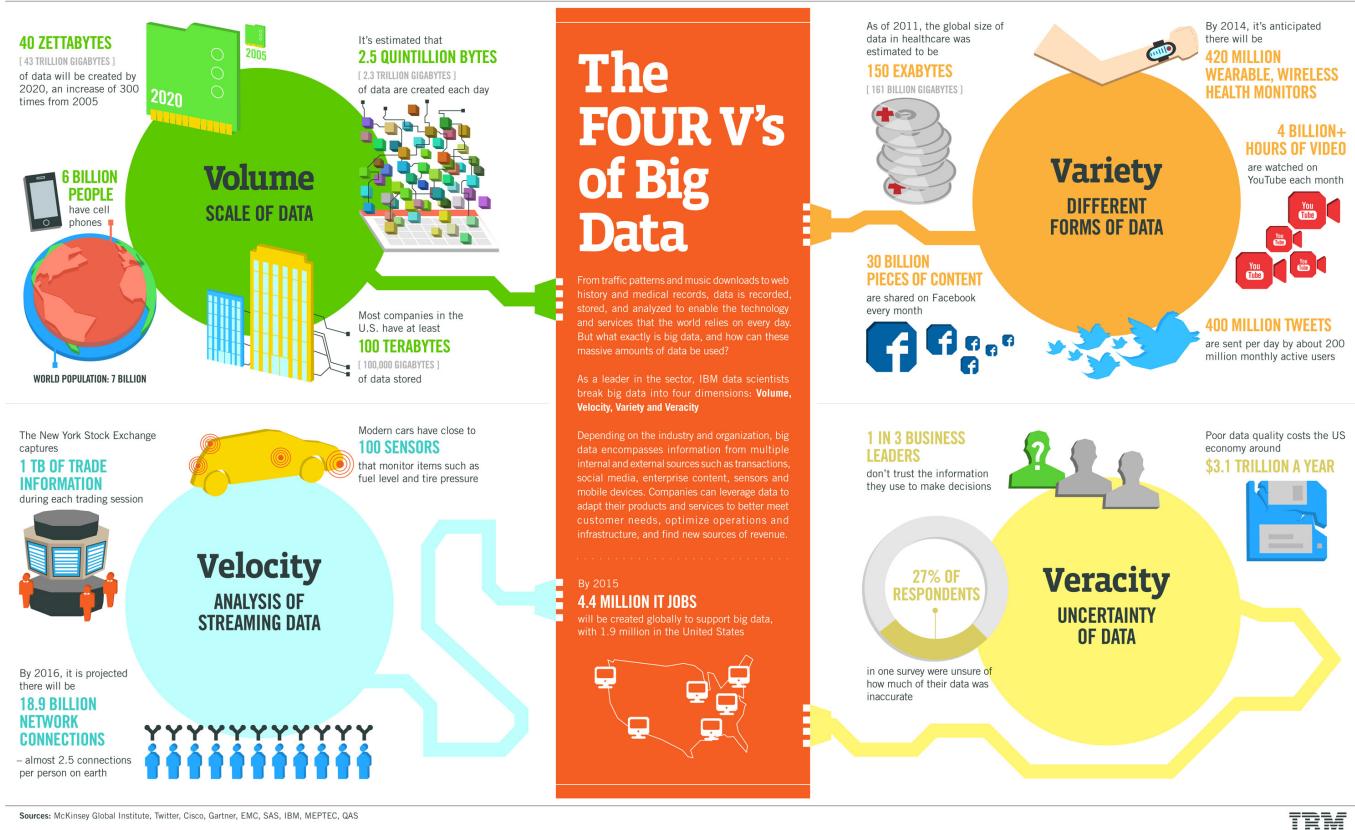
Learning Model

learning algorithm, A and hypothesis set, H

Why Now?

More Data

- Increase of data **Volume, Variety, Velocity, and Veracity**
- awesome-public-datasets (<https://github.com/caesar0301/awesome-public-datasets>)

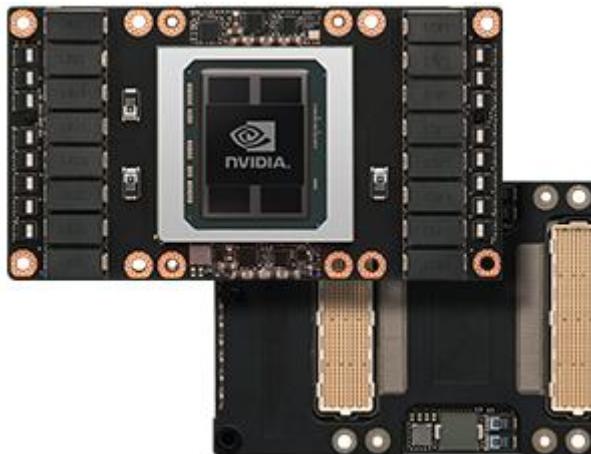


Dedicate Machine Learning Hardware

- Increase of computing power with dedicate hardware, Deep Learning Supercomputer in a box.

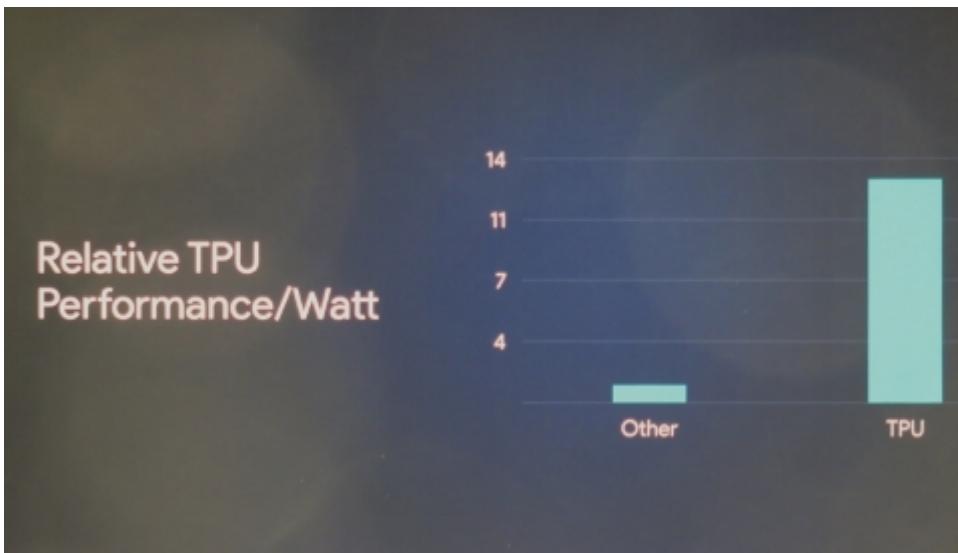
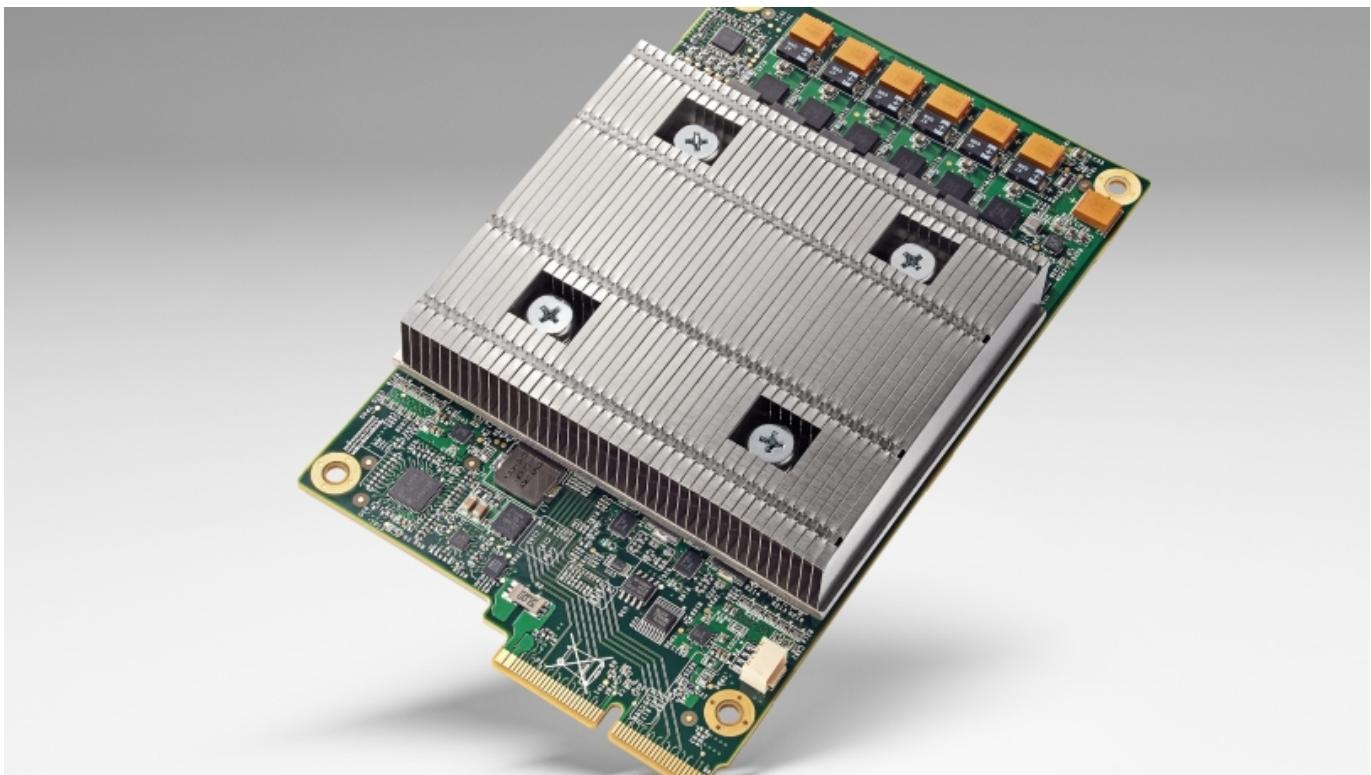


- Nvidia Tesla P100 (<https://www.technologyreview.com/s/601195/a-2-billion-chip-to-accelerate-artificial-intelligence/>)



Google TPU

TPU runs using 8-bit integer math, instead of the higher-precision floating-point math for which most modern CPUs and GPUs are designed. Most machine-learning algorithms can get by fine with lower resolution data, which means the chip can handle more operations in a given area and tackle more complex models efficiently.

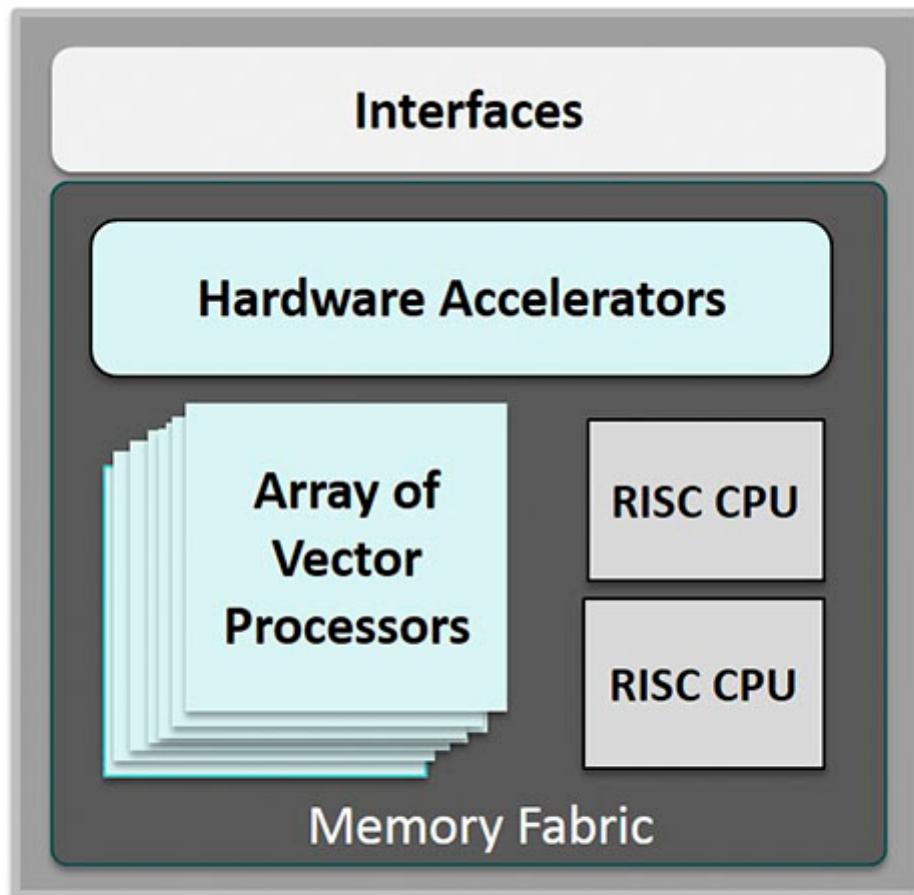


More AI Related Chips

- MIT's 168-core chip could give big brains to mobile devices and robots (<http://www.pcworld.com/article/3029972/components-processors/mits-168-core-chip-could-make-mobile-devices-robots-smarter.html>). MIT says its chips would require a fraction of the resources, and is 10 times more power efficient than a mobile graphics processor. It would be possible to use the chip in wearables, smartphones and battery-operated robots.

Eyeriss will bring self-contained AI capabilities to devices with most of the processing happening locally on a device. Wi-Fi or cellular connections won't be needed to tap into cloud services or servers for image or object recognition.

- A chip startup Movidius (<http://www.movidius.com/>) makes low-power chips it calls vision processing units (or VPUs), which can be part of mobile device.



Myriad 2 Vision Processor Unit (VPU)

- Intel purchased Altera, a leading maker of FPGAs (field-programmable gate arrays), which are somewhere in the middle; they're not as general purpose as GPUs or as specifically designed for TensorFlow as Google's chip, but can be programmed to do a variety of tasks. Microsoft has been experimenting with Altera FPGAs for deep learning.
- IBM is developing its TrueNorth Neurosynaptic chip designed specifically for neural nets, which has recently begun to be used in a variety of applications.
- Cadence (Tensilica), Freescale, and Synopsys are pushing their DSPs (digital-signal processors) to run these models;
- Mobileye and NXP recently announced chips designed specifically for ADAS and self-driving cars;
- Nervana have announced plans for chips specifically designed for AI.

More Algorithm Development and Communities

- More machine learning algorithms and theories are developed by researchers.
- More industry support.
 - OpenAI (<https://openai.com/blog/openai-technical-goals/>)

When to Used?

- We cannot fully predict the problem and human expertise does not exist (navigating on Mars).
- Humans are unable to explain their expertise (speech recognition, play chess or go).

- Solution changes in time (routing on a computer network).
- Solution needs to be adapted to particular cases (user biometrics, recommendations).
- ...

Computer Language for Big Data and Machine Learning

There is a quora discussion notes (<https://www.quora.com/What-is-the-best-language-to-use-while-learning-machine-learning-for-the-first-time>)

A performance table from Julia website (<http://julialang.org/>) can be used as reference.

High-Performance JIT Compiler

Julia's LLVM-based just-in-time (JIT) compiler combined with the language's design allow it to approach and often match the performance of C. To get a sense of relative performance of Julia compared to other languages that can or could be used for numerical and scientific computing, we've written a small set of micro-benchmarks in a variety of languages: [C](#), [Fortran](#), [Julia](#), [Python](#), [Matlab/Octave](#), [R](#), [JavaScript](#), [Java](#), [Lua](#), [Go](#), and [Mathematica](#). We encourage you to skim the code to get a sense for how easy or difficult numerical programming in each language is. The following micro-benchmark results were obtained on a single core (serial execution) on an Intel(R) Xeon(R) CPU E7-8850 2.00GHz CPU with 1TB of 1067MHz DDR3 RAM, running Linux:

	Fortran	Julia	Python	R	Matlab	Octave	Mathematica	JavaScript	Go	LuaJIT	Java
	gcc 5.1.1	0.4.0	3.4.3	3.2.2	R2015b	4.0.0	10.2.0	V8 3.28.71.19	go1.5	gsl-shell 2.3.1	1.8.0_45
<code>fib</code>	0.70	2.11	77.76	533.52	26.89	9324.35	118.53	3.36	1.86	1.71	1.21
<code>parse_int</code>	5.05	1.45	17.02	45.73	802.52	9581.44	15.02	6.06	1.20	5.77	3.35
<code>quicksort</code>	1.31	1.15	32.89	264.54	4.92	1866.01	43.23	2.70	1.29	2.03	2.60
<code>mandel</code>	0.81	0.79	15.32	53.16	7.58	451.81	5.13	0.66	1.11	0.67	1.35
<code>pi_sum</code>	1.00	1.00	21.99	9.56	1.00	299.31	1.69	1.01	1.00	1.00	1.00
<code>rand_mat_stat</code>	1.45	1.66	17.93	14.56	14.52	30.93	5.95	2.30	2.96	3.27	3.92
<code>rand_mat_mul</code>	3.48	1.02	1.14	1.57	1.12	1.12	1.30	15.07	1.42	1.16	2.36

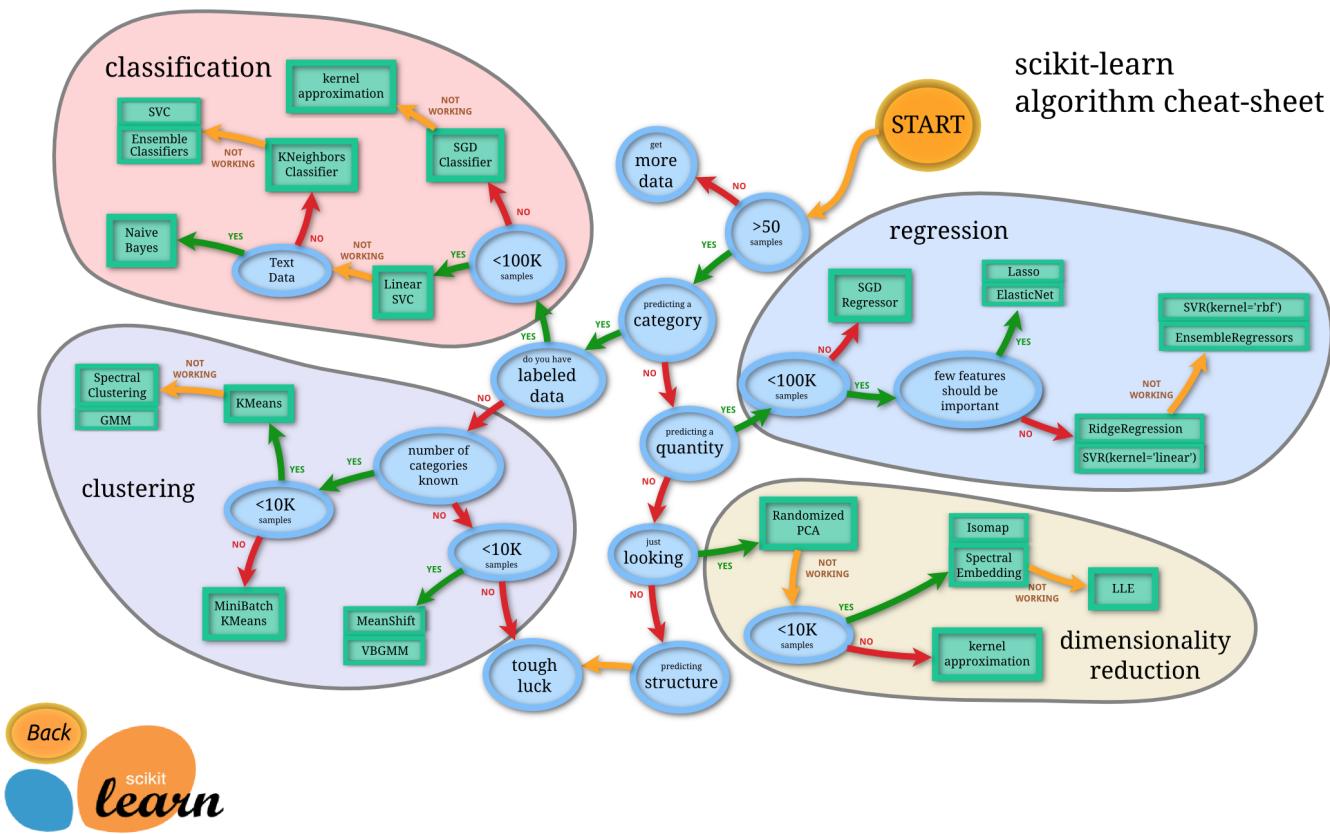
Figure: benchmark times relative to C (smaller is better, C performance = 1.0).

C and Fortran compiled by gcc 5.1.1, taking best timing from all optimization levels (-O0 through -O3). C, Fortran, Go, and Julia use [OpenBLAS](#) v0.2.14.

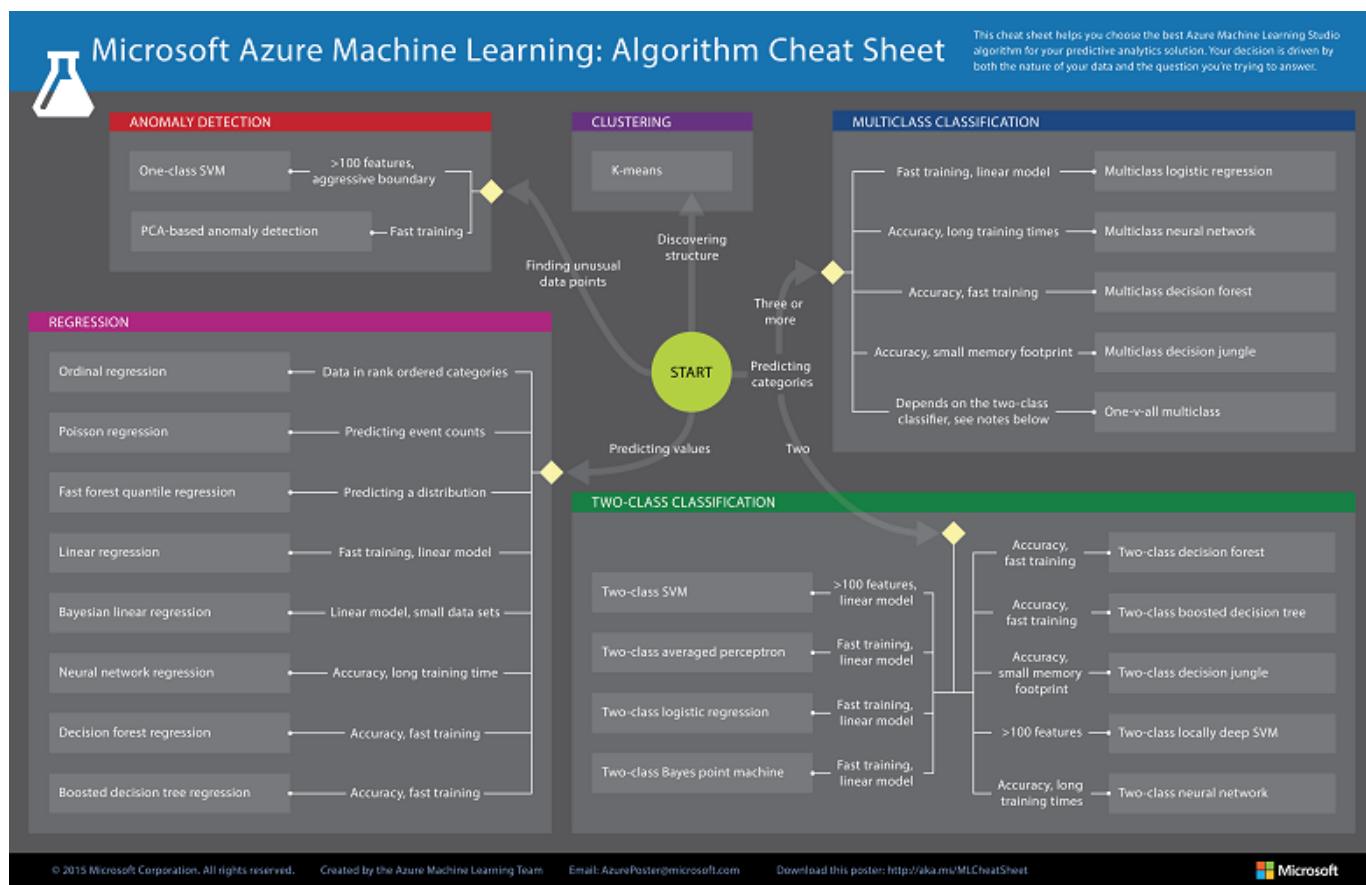
Python 3 was installed from the [Anaconda distribution](#). The Python implementations of `rand_mat_stat` and `rand_mat_mul` use NumPy (v1.9.2) functions; the rest are pure Python implementations.

Benchmarks can also be seen [here as a plot](#) created with [Gadfly](#).

Algorithm - scikit-learn

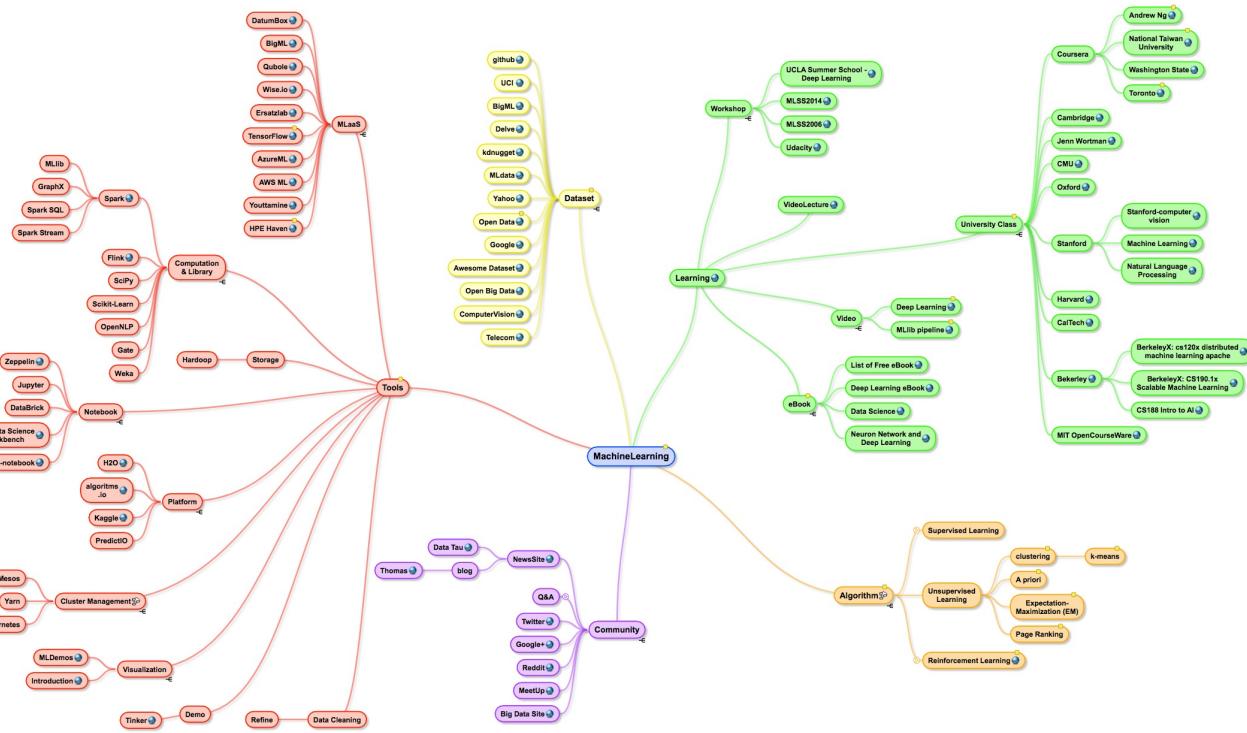


Algorithm - Microsoft Azure Machine Learning



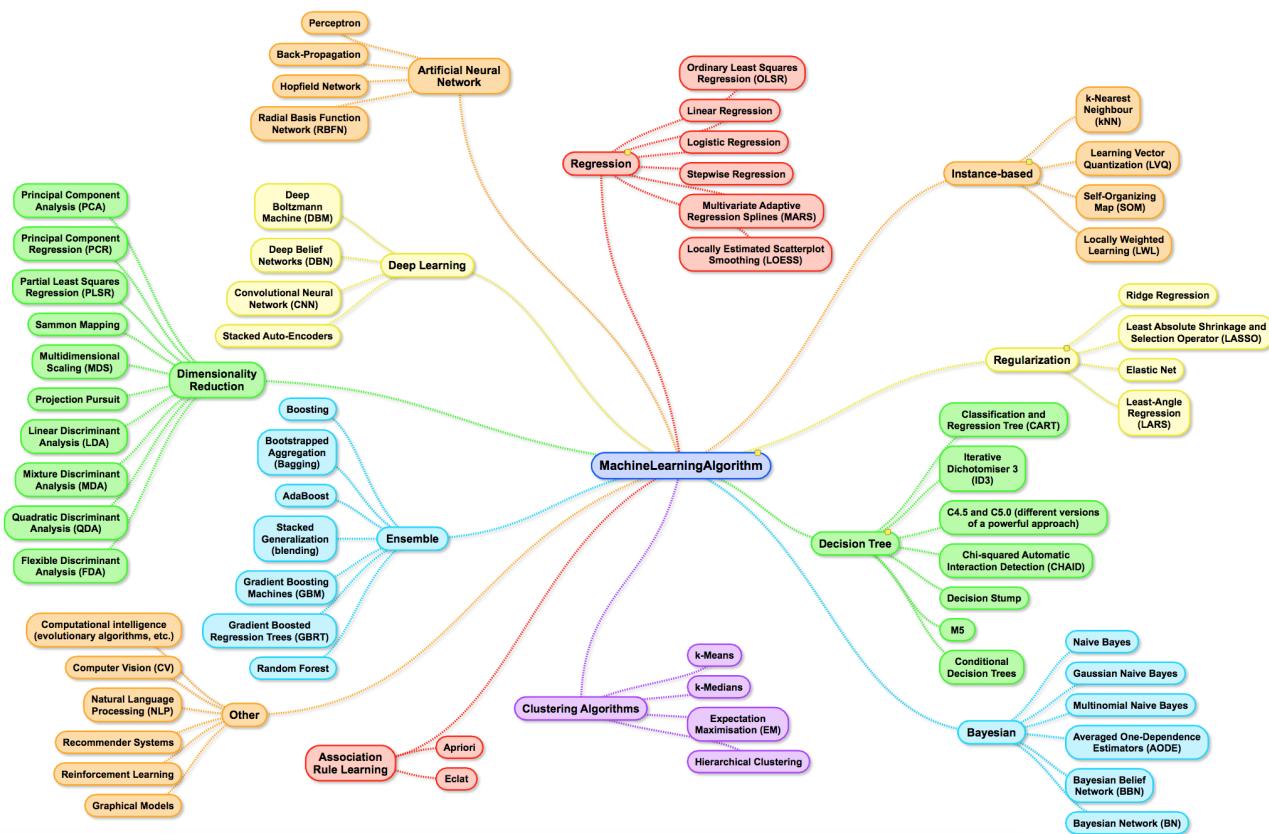
Learning Machine Learning

An ecosystem for learning machine learning.



Algorithm

A subset of machine learning algorithm.



eBook

machine learning ebooks

(https://github.com/rasbt/pattern_classification/blob/master/resources/machine_learning_ebooks.md), deep learning (<http://www.deeplearningbook.org/>)

data science books (<http://www.wzchen.com/data-science-books/>)

neural network and deep learning (<http://neuralnetworksanddeeplearning.com/chap1.html>)

hacker lists - free machine learning books (<https://hackerlists.com/free-machine-learning-books/>)

Vectorization

Vectorization ([https://en.wikipedia.org/wiki/Vectorization_\(mathematics\)](https://en.wikipedia.org/wiki/Vectorization_(mathematics))) makes coding easier and more readable. Many languages

(https://en.wikipedia.org/wiki/List_of_programming_languages_by_type#Array_language) support Array Programming (https://en.wikipedia.org/wiki/Array_programming).

Vectorization example.

$$\rightarrow h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

$$= \theta^T x$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$$

Unvectorized implementation

```

→ prediction = 0.0;
→ for j = 1:n+1,
    prediction = prediction + theta(j) * x(j)
end;

```

Vectorized implementation

```

→ prediction = theta' * x;

```



3:58 / 13:48



Functional Programming and Computing Nodes

We can treat algorithm like a graph, with each node is computation and each link is data flow. This concept is described in Directed Acyclic Graph ([https://en.wikipedia.org/wiki/Directed_acyclic_graph]) (DAG)

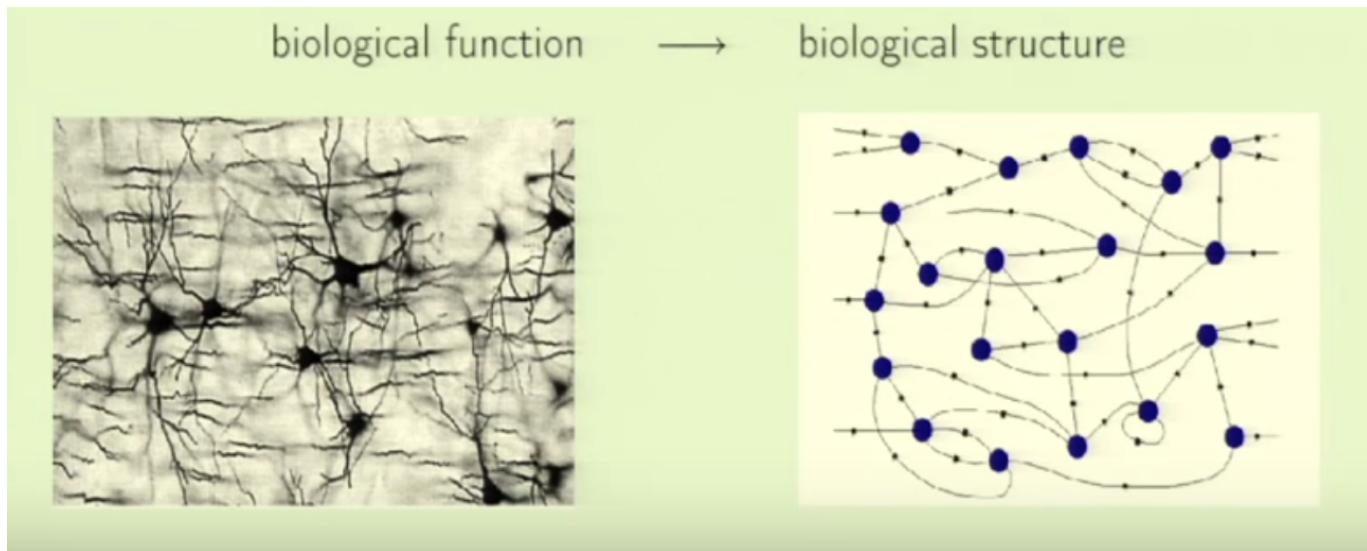
Learning from Nature

biological function

→

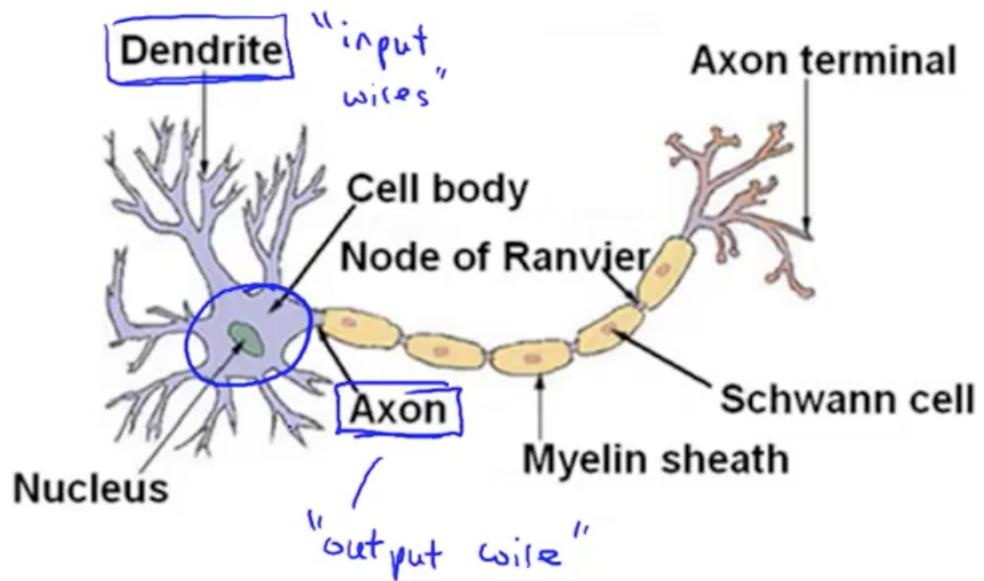
biological structure





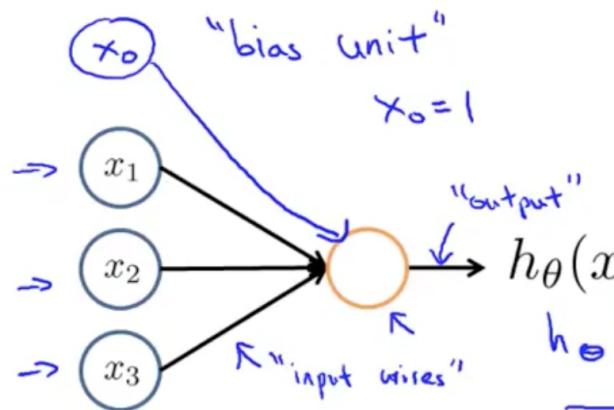
Bio Neural Network

Neuron in the brain



Artificial Neuron Wiki (https://en.wikipedia.org/wiki/Artificial_neuron#Comparison_to_biological_neurons)

Neuron model: Logistic unit



$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

"weights" ↗
↙ parameters

Sigmoid (logistic) activation function.

$$g(z) = \frac{1}{1 + e^{-z}}$$

The first layer can be considered as a (trained) feature layer, then second layer is feature-feature layer,...

Logistic Regression Model

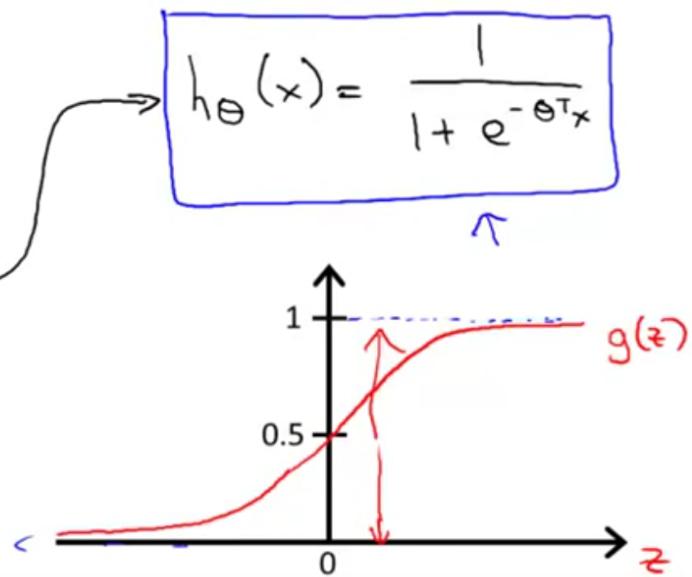
Want $0 \leq h_\theta(x) \leq 1$

$$h_\theta(x) = g(\theta^T x)$$

$$\rightarrow g(z) = \frac{1}{1 + e^{-z}}$$

$\theta^T x$

↳ Sigmoid function
↳ Logistic function



Activation Function

There are different types of activation functions.

Activation functions

Logistic

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$



$$\frac{\partial \sigma}{\partial z}(z) = \sigma(z)(1 - \sigma(z))$$

Hyperbolic Tangent

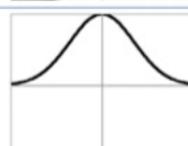
$$\sigma(z) = \frac{1 - \exp(-2z)}{1 + \exp(-2z)}$$



$$\frac{\partial \sigma}{\partial z}(z) = 1 - (\sigma(z))^2$$

Gaussian

$$\sigma(z) = \exp(-z^2/2)$$



$$\frac{\partial \sigma}{\partial z}(z) = -z\sigma(z)$$

Linear

$$\sigma(z) = z$$



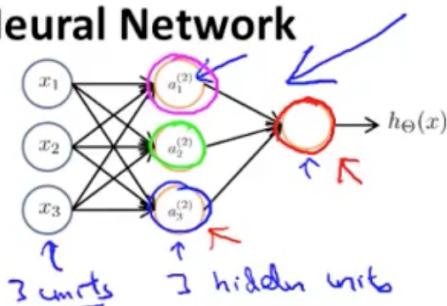
$$\frac{\partial \sigma}{\partial z}(z) = 1$$

And others...



Most of those functions have “mathematics convenience”, which simplifies the downstream equation calculation.

Neural Network



$\rightarrow a_i^{(j)}$ = “activation” of unit i in layer j

$\rightarrow \Theta^{(j)}$ = matrix of weights controlling function mapping from layer j to layer $j + 1$

$$h_{\Theta}(x)$$

$$\rightarrow a_1^{(2)} = g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3)$$

$$\rightarrow a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3)$$

$$\rightarrow a_3^{(2)} = g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3)$$

$$\rightarrow h_{\Theta}(x) = a_1^{(3)} = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

$$(2)$$

\rightarrow If network has s_j units in layer j , s_{j+1} units in layer $j + 1$, then $\Theta^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$.

$$s_{j+1} \times (s_j + 1)$$

Artificial Neural Network (https://en.wikipedia.org/wiki/Artificial_neural_network)

Deep Learning (https://en.wikipedia.org/wiki/Deep_learning#Deep_neural_networks)

Cost function

Logistic regression:

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

i - index for training data m
s_l - number of units in layer l
L - total number of layers
l - index for layer L
n - number of features
m - number of test data
K - number of output/class

for each test data
 find cost for each class
 +
 for each layer
 calculate THETA of each node

Neural network:

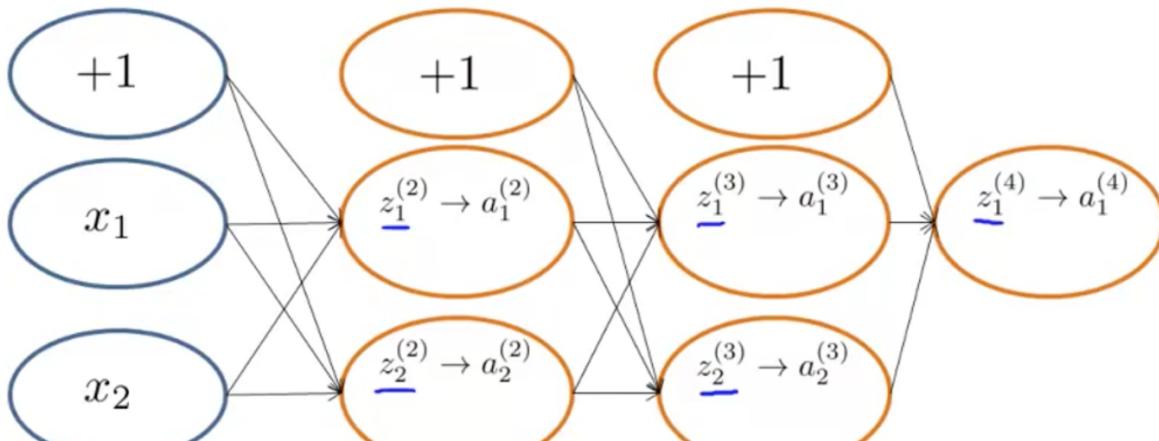
THETA is next layer s_l x previous layer s_{l-1} + bias $\rightarrow h_\Theta(x) \in \mathbb{R}^K \quad (h_\Theta(x))_i = i^{th}$ output

$$\rightarrow J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\Theta(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_\Theta(x^{(i)}))_k) \right]$$

$$\boxed{\frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2}$$

Andrew Ng

Forward Propagation



$\rightarrow \delta_j^{(l)}$ = "error" of cost for $a_j^{(l)}$ (unit j in layer l).

Formally, $\delta_j^{(l)} = \frac{\partial \text{cost}(i)}{\partial z_j^{(l)}}$ (for $j \geq 0$), where

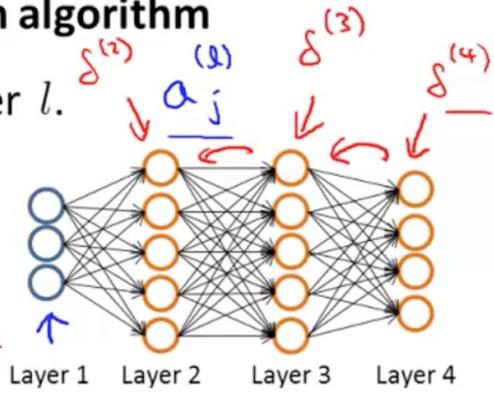
$$\text{cost}(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\Theta(x^{(i)}))$$

Gradient computation: Backpropagation algorithm

Intuition: $\delta_j^{(l)}$ = "error" of node j in layer l .

For each output unit (layer $L = 4$)

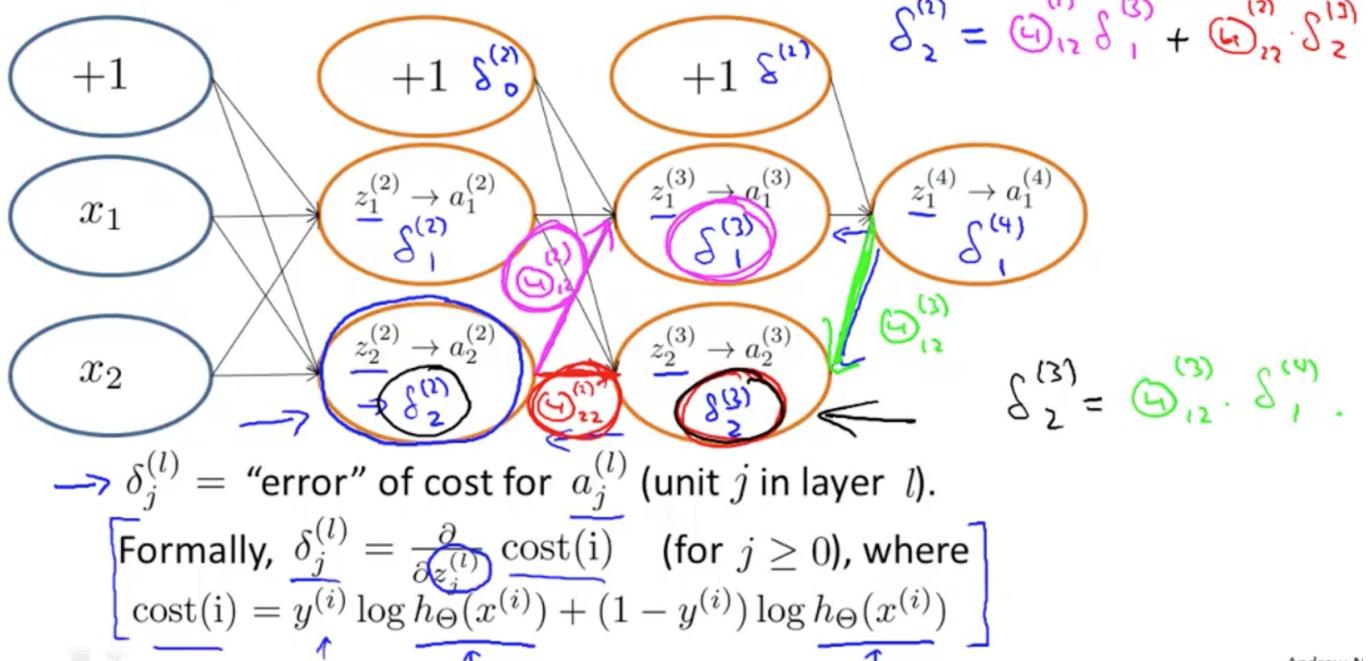
$$\delta_j^{(4)} = \underline{a_j^{(4)}} - \underline{y_j} \quad (\underline{h_\Theta(x)})_j \quad \underline{\delta^{(4)}} = \underline{a} - \underline{y}$$



$$\begin{aligned} \rightarrow \delta^{(3)} &= (\Theta^{(3)})^T \delta^{(4)} \cdot g'(z^{(3)}) \\ \rightarrow \delta^{(2)} &= (\Theta^{(2)})^T \delta^{(3)} \cdot g'(z^{(2)}) \\ (\text{No } \delta^{(1)}) & \frac{\partial}{\partial \Theta_{ij}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)} \quad (\text{ignoring } \lambda; \text{ if } \lambda = 0) \end{aligned}$$

See MIT 6.034 lecture-12 for derivation of gradient descent formula; $a_3 \cdot (1 - a_3)$

Backward Propagation

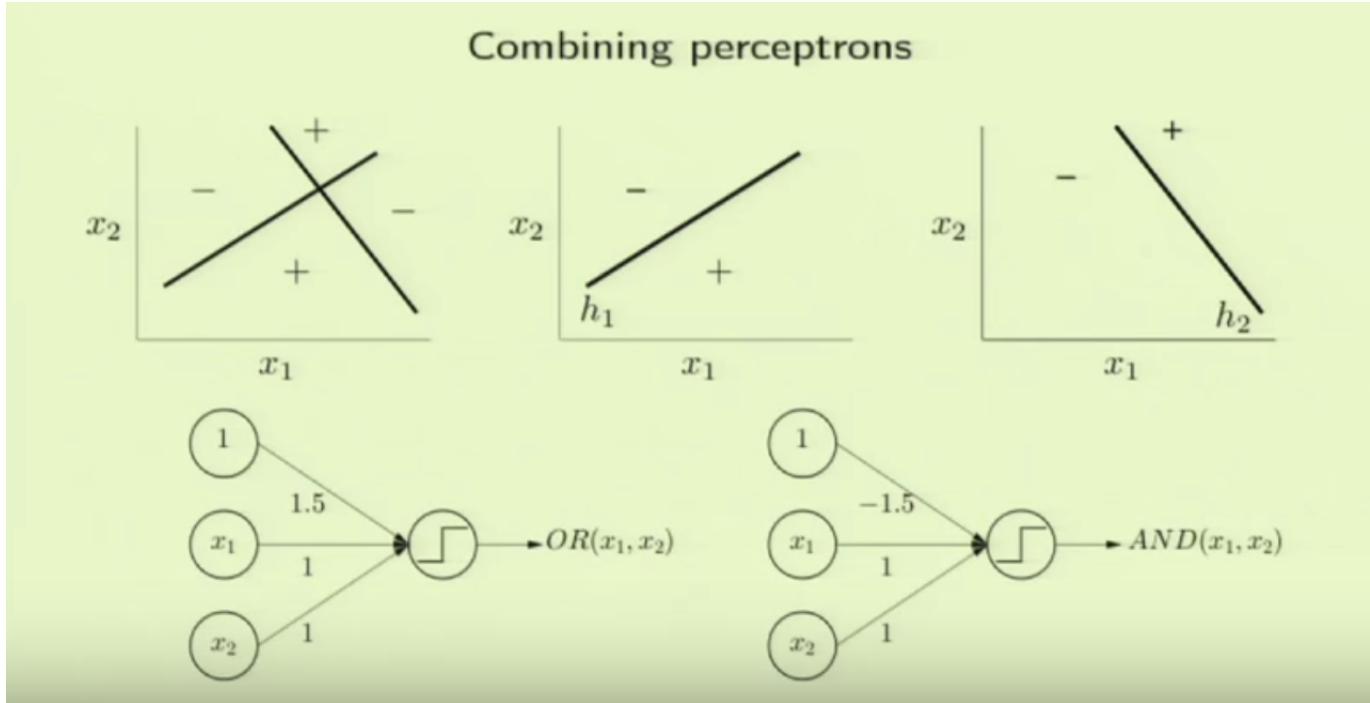


Andrew N

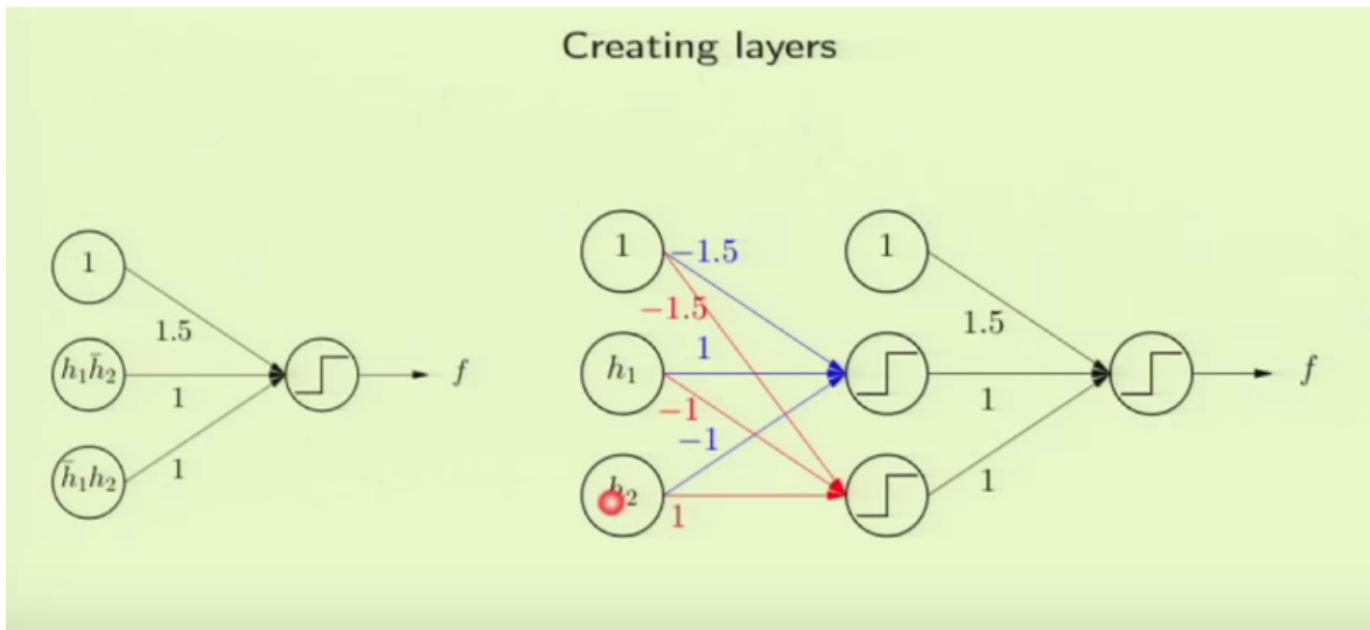
More detailed computation steps.

Caltech CS 156 Machine Learning - Learning from Data

lecture-10 (<http://work.caltech.edu/telecourse.html>)

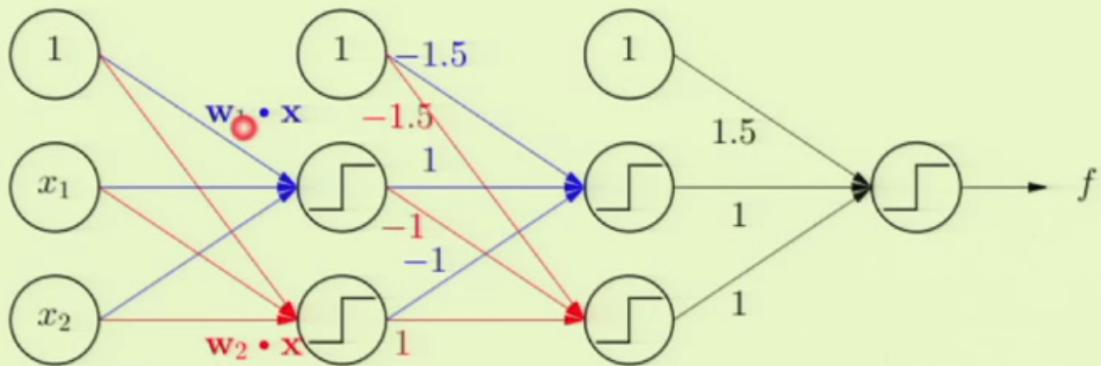


One simple logistic regression can not separate the testing data. There are many different activation functions can be used here, $\tanh(s)$ is here. See selection approach
(<http://stats.stackexchange.com/questions/101560/tanh-activation-function-vs-sigmoid-activation-function>)



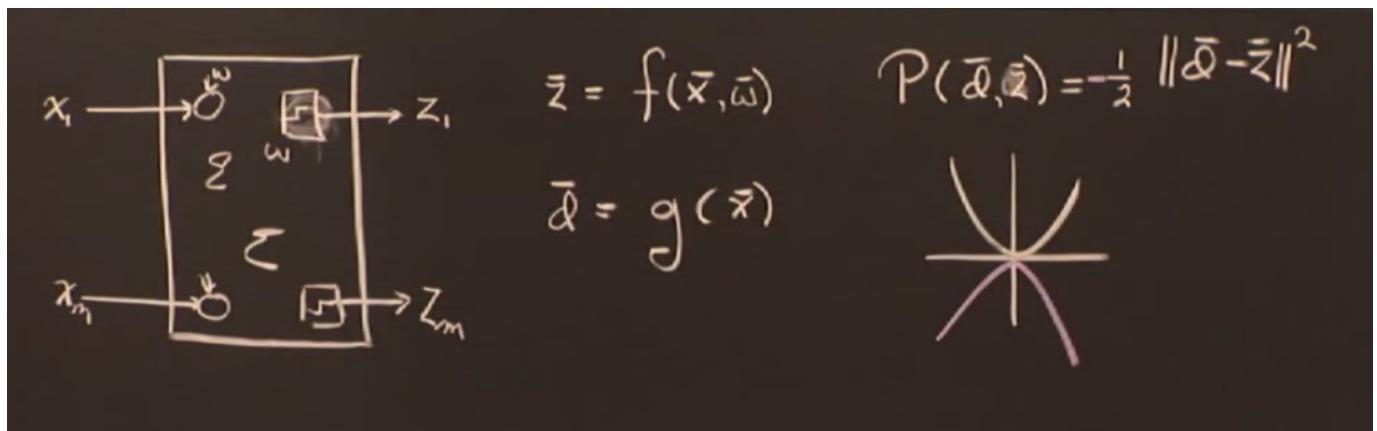
We have to use two separate nodes to cover the problem space.

The multilayer perceptron



This is two features ($n=2$), two hidden layers ($L=2$), one classification ($K=1$)

MIT Course Number 6.034 lecture-12
[\(https://www.youtube.com/watch?v=q0pm3BrIUFo\)](https://www.youtube.com/watch?v=q0pm3BrIUFo)



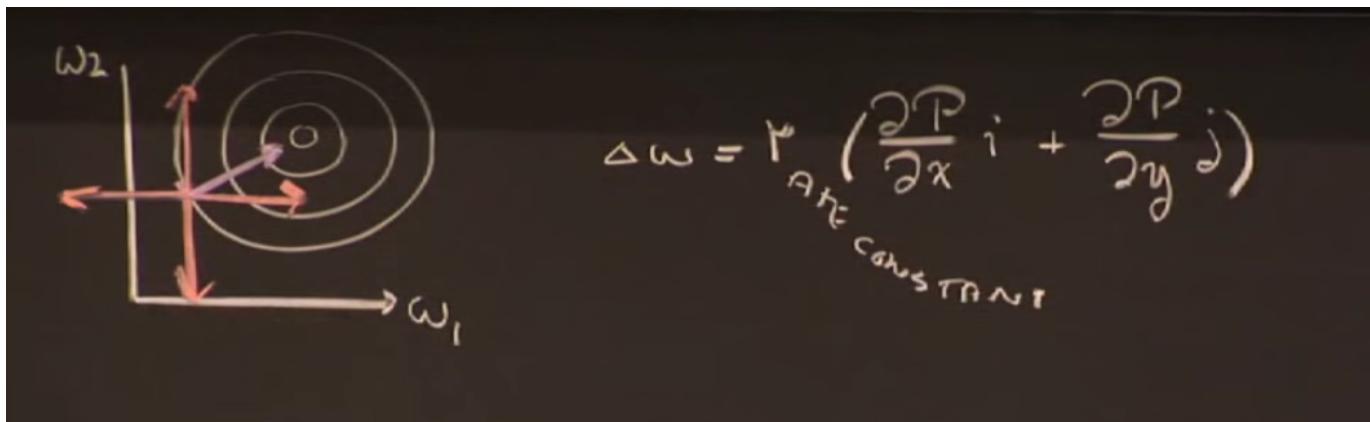
x_1, \dots, x_n is input,

z_1, \dots, z_n is output, which is equivalent to y_1, \dots, y_n in Prof. Ng's lecture, and (x_1, z_1) is a pair.

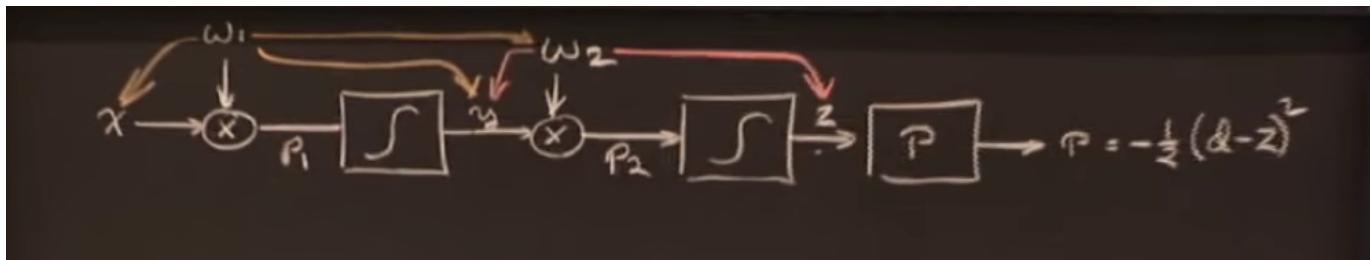
d is \hat{y} , it is the value calculated based on hypothesis.

P is performance, is a cost function.

y is a_2 in Prof. Ng's notes.



There is a typo in picture, x and y should be w1 and w2.



w1 is input layer, consider x is a single variable or a vector.

w2 is hidden layer,

z is output layer,

p is error, cost function.

This model is set for proving backpropagation.

$$\begin{aligned}\frac{\partial P}{\partial w_2} &= \frac{\partial P}{\partial z} \frac{\partial z}{\partial w_2} = (d-z) \frac{\partial z}{\partial p_2} \frac{\partial p_2}{\partial w_2} = (d-z) \frac{\partial z}{\partial p_2} y \\ \frac{\partial P}{\partial w_1} &= \frac{\partial P}{\partial z} \frac{\partial z}{\partial w_1} = (d-z) \frac{\partial z}{\partial p_2} \frac{\partial p_2}{\partial w_1} = (d-z) \frac{\partial z}{\partial p_2} \frac{\partial y}{\partial w_1} \\ &= (d-z) \frac{\partial z}{\partial p_2} w_2 \frac{\partial y}{\partial p_1} \times y(1-y)\end{aligned}$$

This exercise proves the performance improvement is local dependency, e.g. for $\partial(p/w_2)$ is dependent on $(d-z)$, y , and $\partial(z/p_2)$.

$$\partial(z/p_2) = z^*(1-z)$$

Use Cases

- equipment failure prediction

- facial recognition
- speech recognition
- text classification
- self-driving car
- smart home
- surveillance and security
- medical image and diagnostic
- spam discovery and filtering
- predictive maintenance
- ...

Versions

A study note about Learning Machine Learning,

- v.0.0.1, April-10 2016, initial version, Richard Kuo, at La Boulanger, Mountain View, CA
- v.0.0.2, June-20 2016, add history and update mindmap, Richard Kuo, at La Boulanger, Mountain View, CA