

Machine Learning

Neural Networks

Carl Henrik Ek - carlhenrik.ek@bristol.ac.uk

November 20, 2017

<http://www.carlhenrik.com>

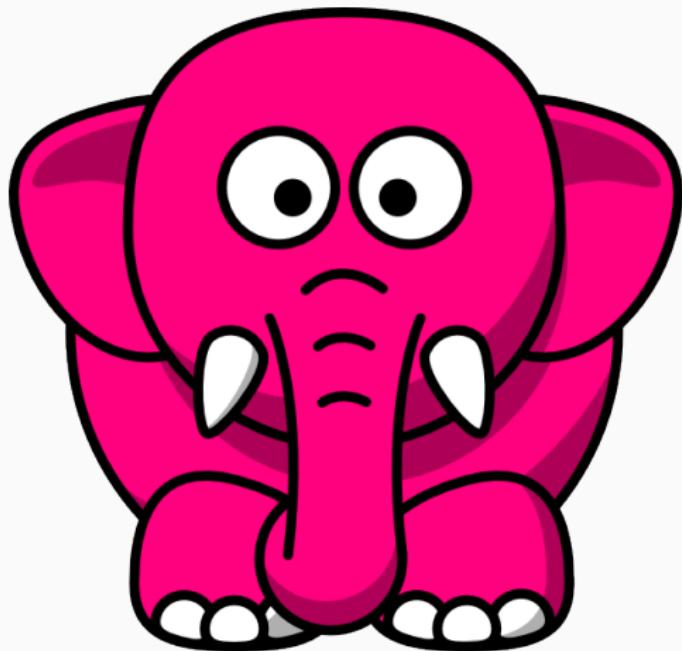
Introduction

- Extra lab
 - Thursday 17-19 2.11
- Coursework 1
 - 60% done
- Tomorrow

Lecture 1 Decisions

Lecture 2 Machine learning project

Today



Darthmouth workshop¹



¹https://en.wikipedia.org/wiki/Dartmouth_workshop

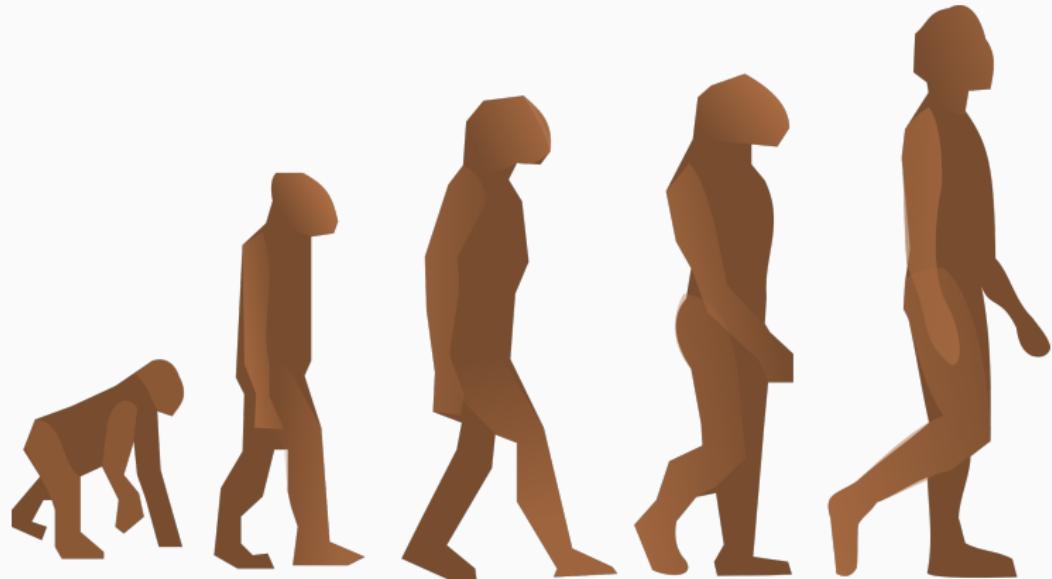
We propose that a 2 month, 10 man study of artificial intelligence be carried out during the summer of 1956 at Dartmouth College in Hanover, New Hampshire. The study is to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it. An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves. We think that a significant advance can be made in one or more of these problems if a carefully selected group of scientists work on it together for a summer.

– August 31st 1955

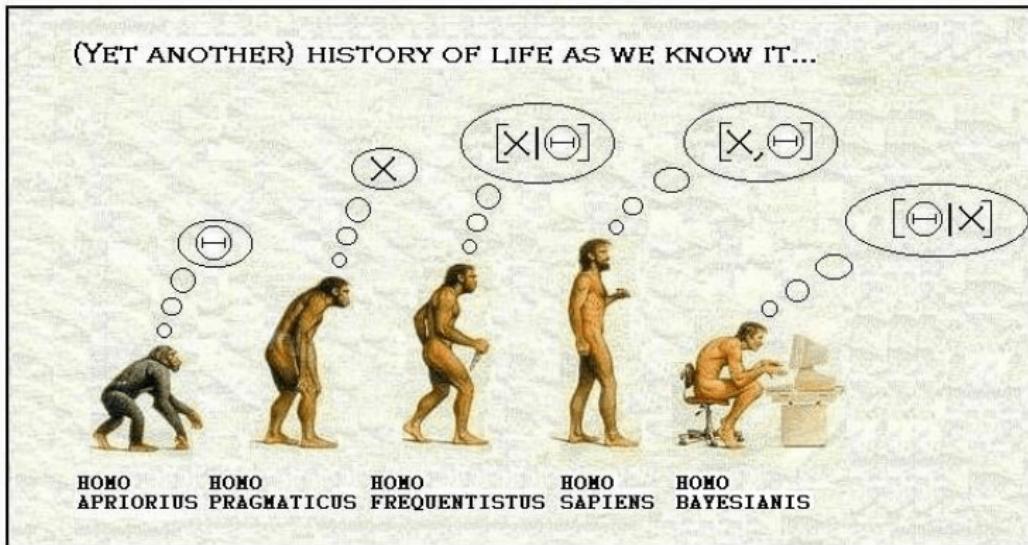
- Artificial Intelligence
 - the idea that intelligence could be described as a set of logical arguments and rules
 - *failed spectacularly*

- Artificial Intelligence
 - the idea that intelligence could be described as a set of logical arguments and rules
 - *failed spectacularly*
- Machine learning
 - it is naive to think that we can understand everything explicitly (Laplace)
 - can we describe rules for how to learn from evidence/data?

The learning machine



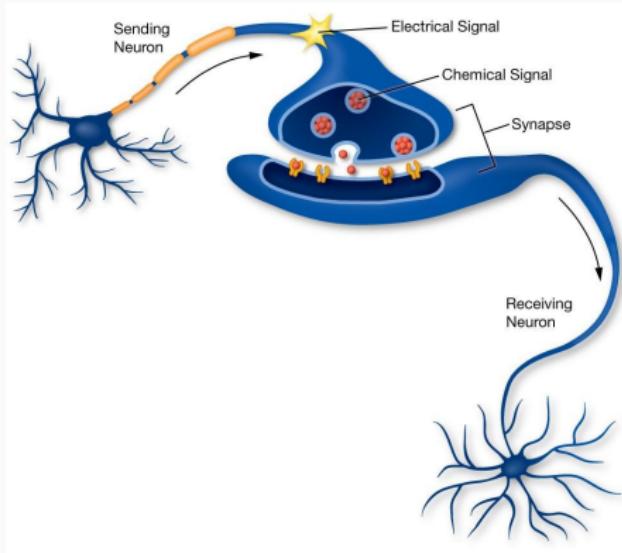
Another view



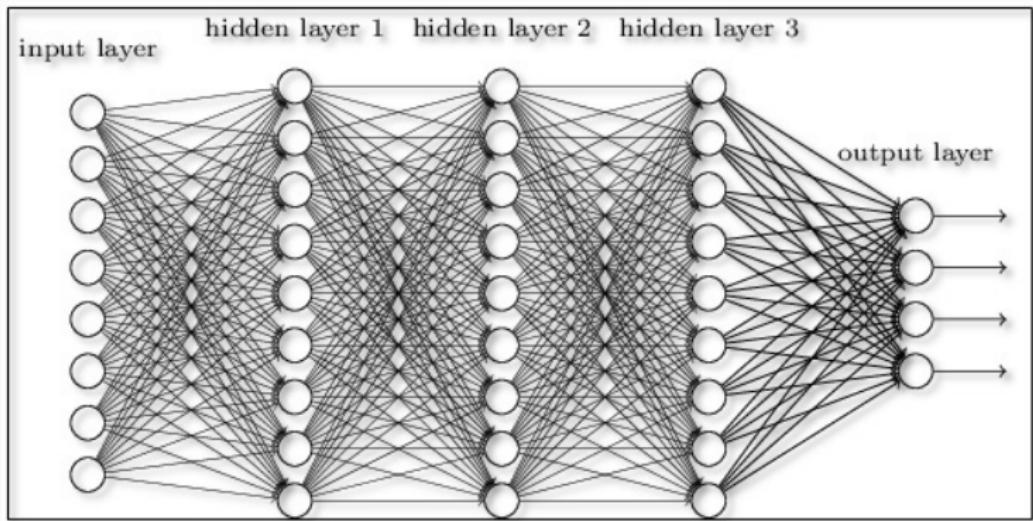
Neural Networks

- First described in 1940 by Walter Pitts
- First computational model called **Perceptron** described in Rosenblatt 1958
- Idea that we should composite a function into several functions
- Usually motivated by "neuroscience"

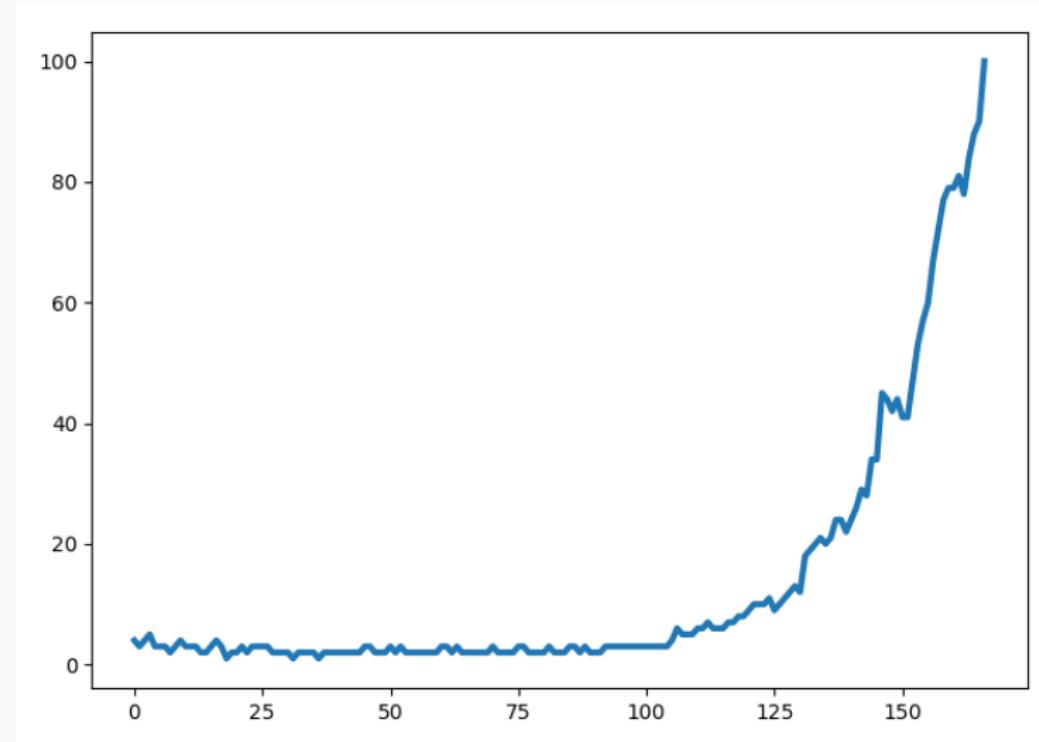
Neurons and spike trains



Neural Networks







Neural Networks

Functions

$$y(\mathbf{x}, \mathbf{w}) = f \left(\sum_{j=1}^M w_j \phi_j(\mathbf{x}) \right)$$

- regression: f is identity function
- classification: f is nonlinear function
- f : an activation function

Neural Networks: hidden layer

$$a_j = \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)}$$

$$z_j = h(a_j)$$

a_j activation

w_{ji} weight

w_{j0} bias

z_j hidden units

$h(\cdot)$ activation function

Neural Networks: output layer

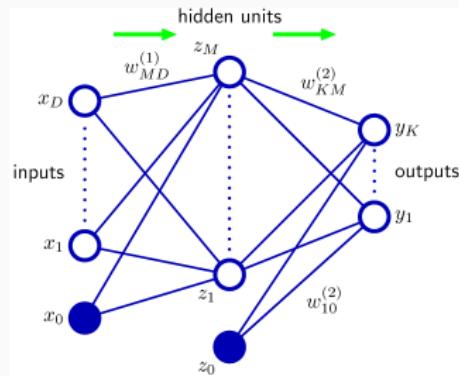
$$a_k = \sum_{j=1}^M w_{kj}^{(2)} z_j + w_{k0}^{(2)}$$

$$y_k = f(a_k) = \sigma(a_k)$$

y_k output variable

$f(\cdot)$ output activation

Neural Network



$$y_k(\mathbf{x}, \mathbf{w}) = \sigma \left(\sum_j^M w_{kj}^{(2)} h \left(\sum_j^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right)$$

Composite Functions

Why are composite functions attractive?

$$y = g(x) = f_K(f_{K-1}(f_{K-2}(\dots f_1(x) \dots)))$$

Linear Algebra

- Kernel of a function

$$\text{Kern}(f_k) = \{(x, x') | f_k(x) = f_k(x')\}$$

Linear Algebra

- Kernel of a function

$$\text{Kern}(f_k) = \{(x, x') | f_k(x) = f_k(x')\}$$

- Image of a function

$$\text{Im}(f_k(x)) = \{y \in Y | y = f_k(x), x \in X\}$$

Linear Algebra

- Rank-Nullity Theorem

$$\dim(\text{Im}(f)) = \dim(V) - \dim(\text{Kern}(f))$$

Linear Algebra

- Rank-Nullity Theorem

$$\dim(\text{Im}(f)) = \dim(V) - \dim(\text{Kern}(f))$$

- Kernel of function

$$\text{Kern}(f_1) \subseteq \text{Kern}(f_{k-1} \circ \dots \circ f_2 \circ f_1) \subseteq \text{Kern}(f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1)$$

Linear Algebra

- Rank-Nullity Theorem

$$\dim(\text{Im}(f)) = \dim(V) - \dim(\text{Kern}(f))$$

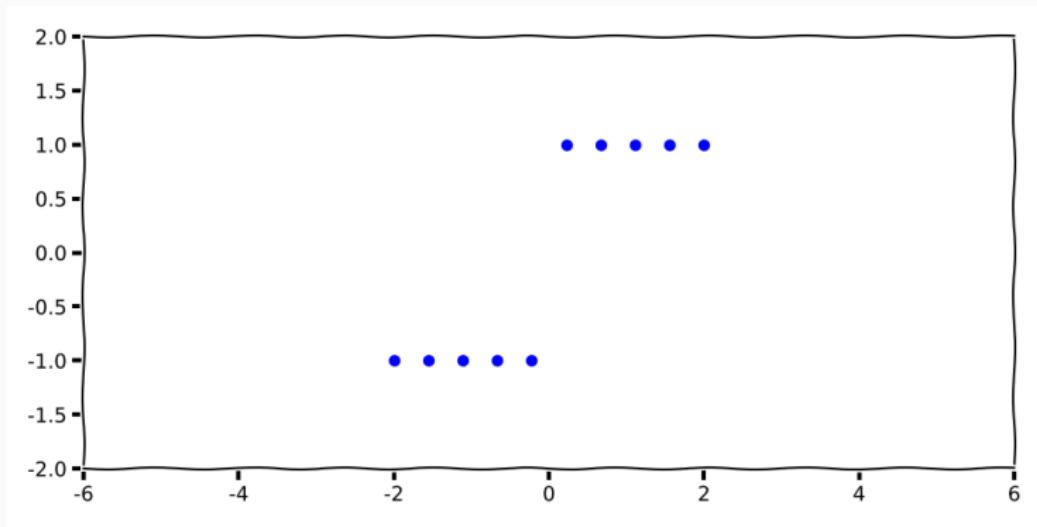
- Kernel of function

$$\text{Kern}(f_1) \subseteq \text{Kern}(f_{k-1} \circ \dots \circ f_2 \circ f_1) \subseteq \text{Kern}(f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1)$$

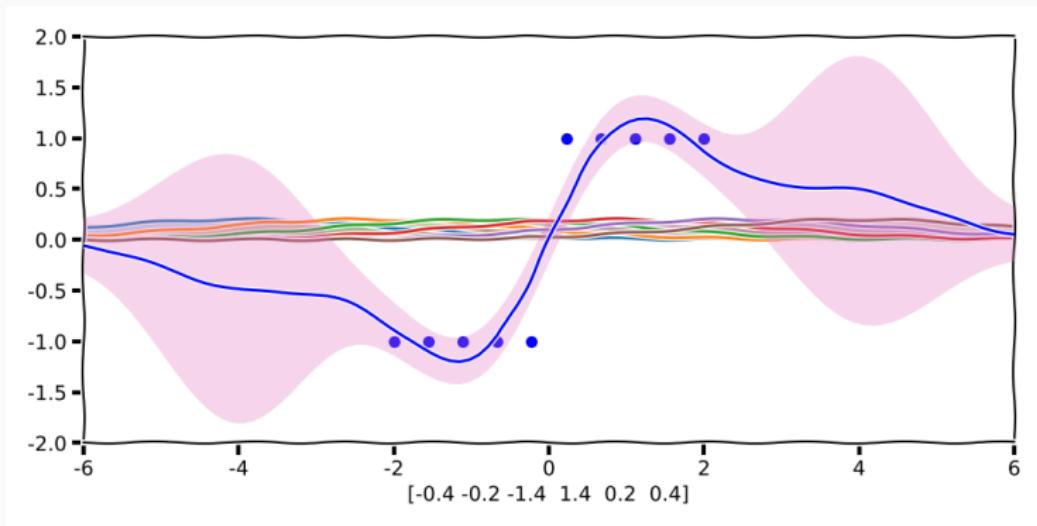
- Image of a function

$$\text{Im}(f_k \circ f_{k-1} \circ \dots \circ f_2 \circ f_1) \subseteq \text{Im}(f_k \circ f_{k-1} \circ \dots \circ f_2) \subseteq \dots \subseteq \text{Im}(f_k)$$

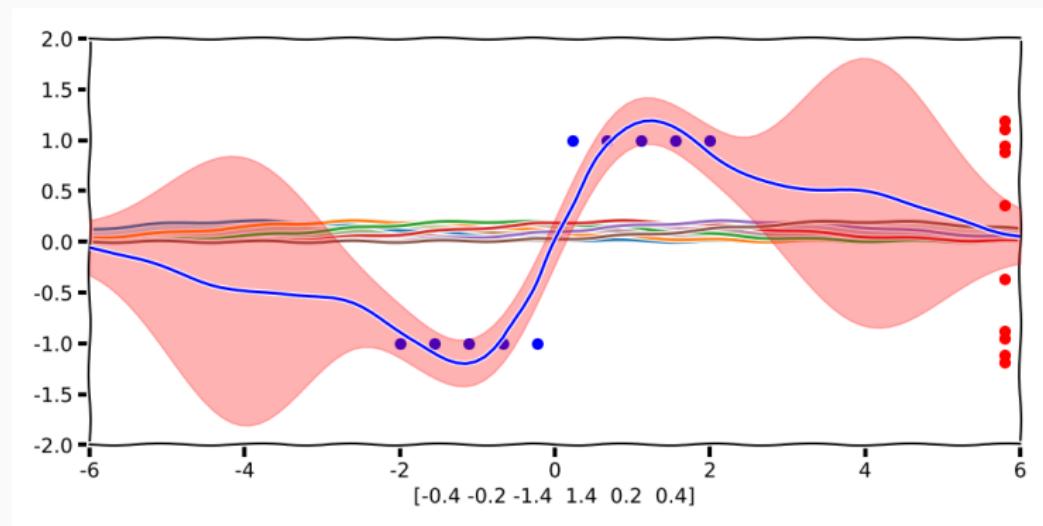
Composite Functions



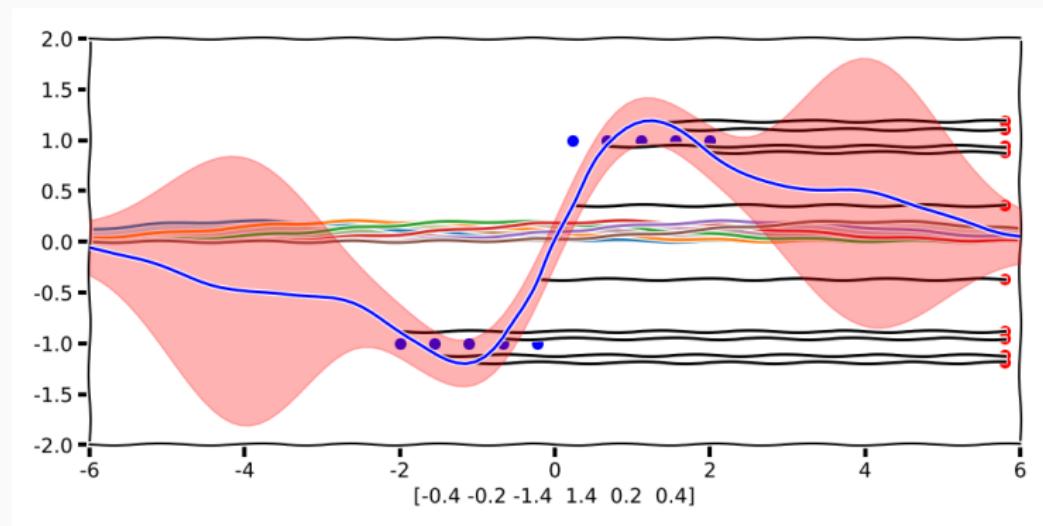
Composite Functions



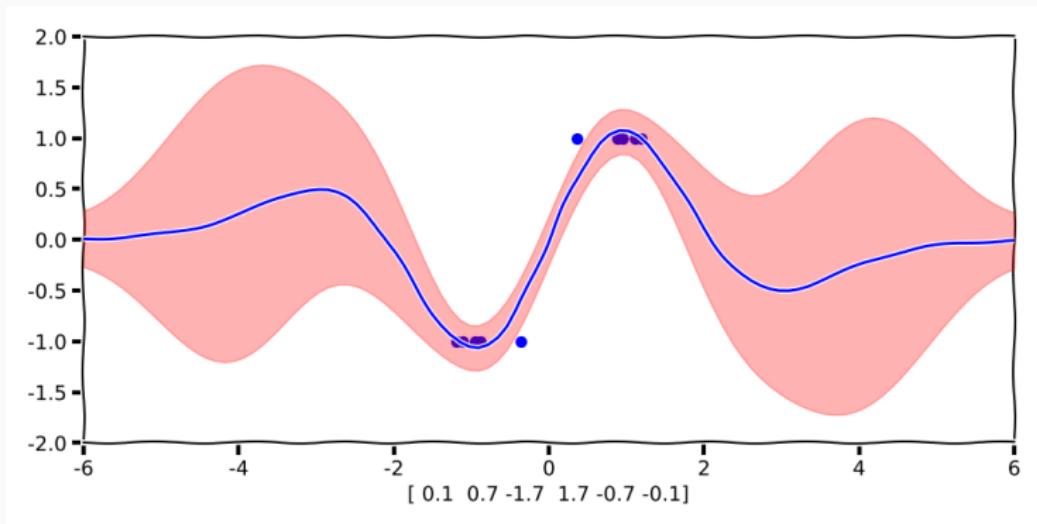
Composite Functions



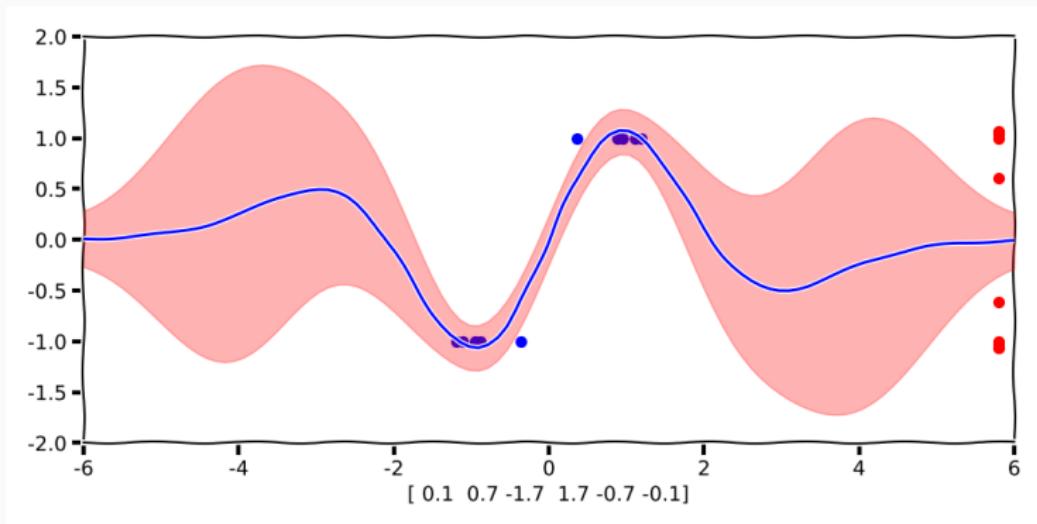
Composite Functions



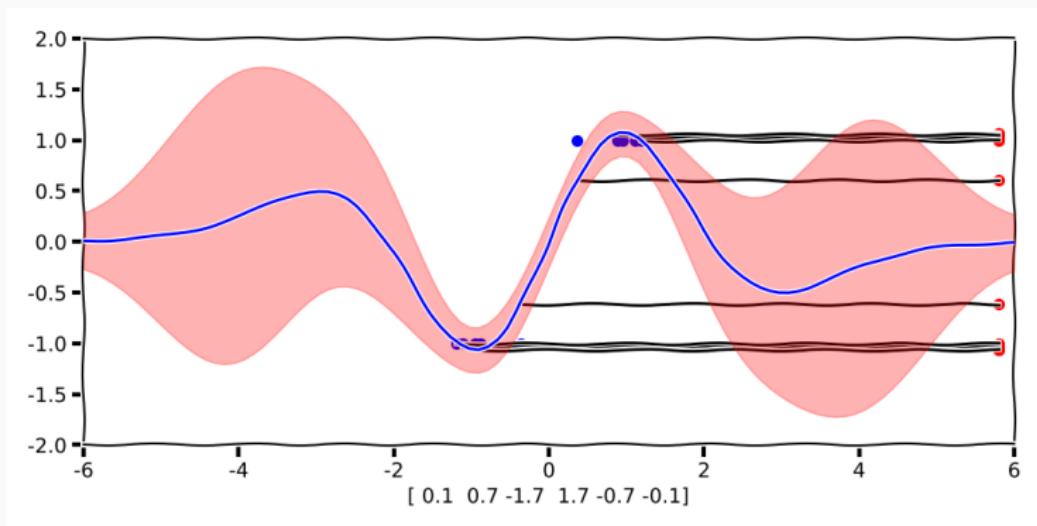
Composite Functions



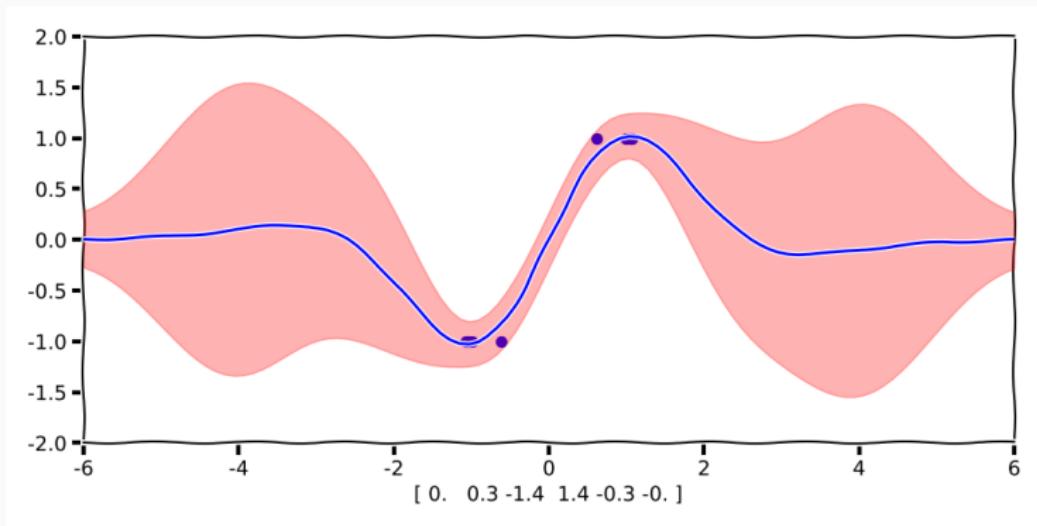
Composite Functions



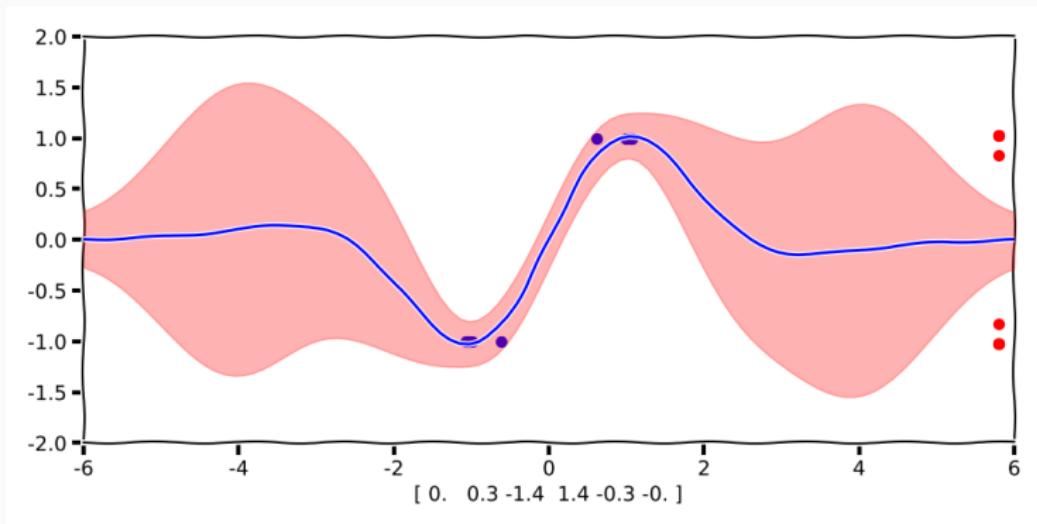
Composite Functions



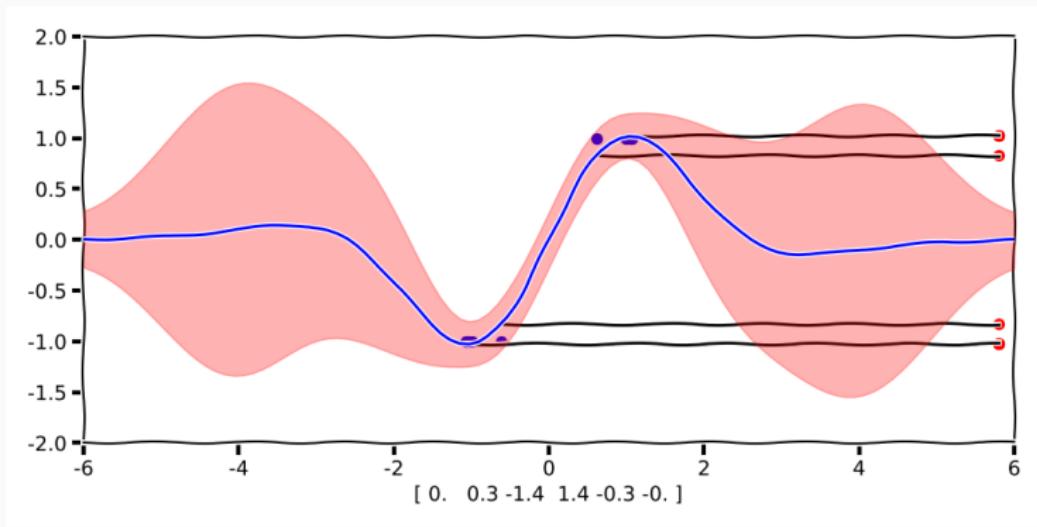
Composite Functions



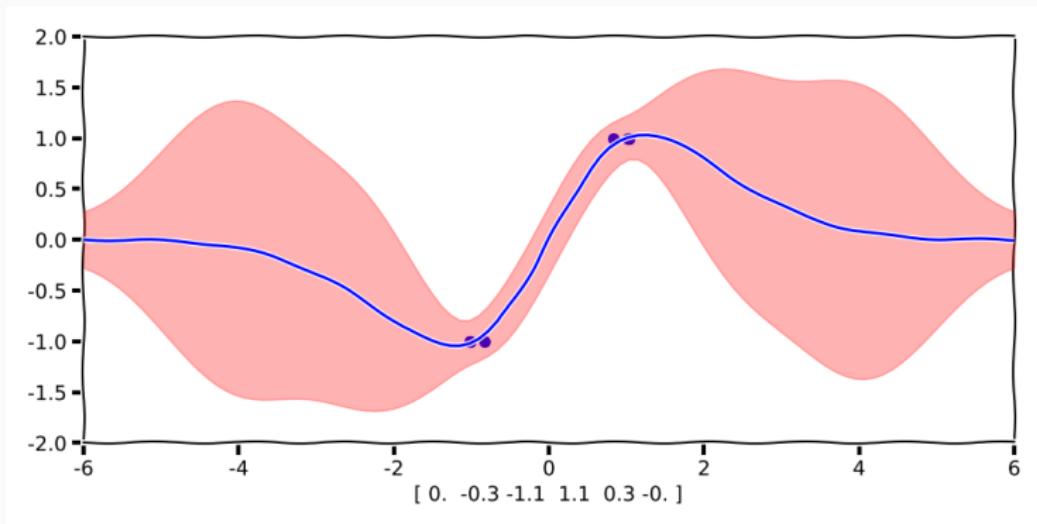
Composite Functions



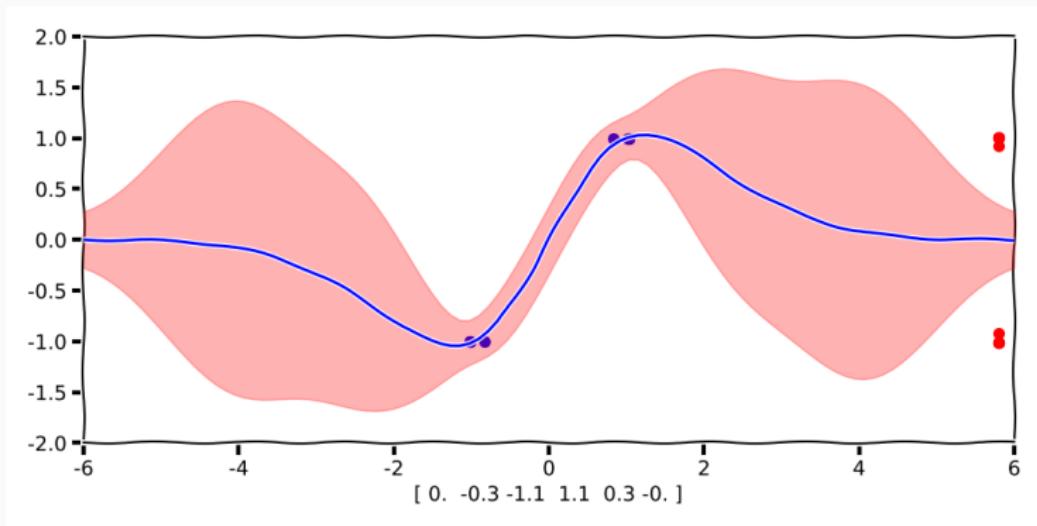
Composite Functions



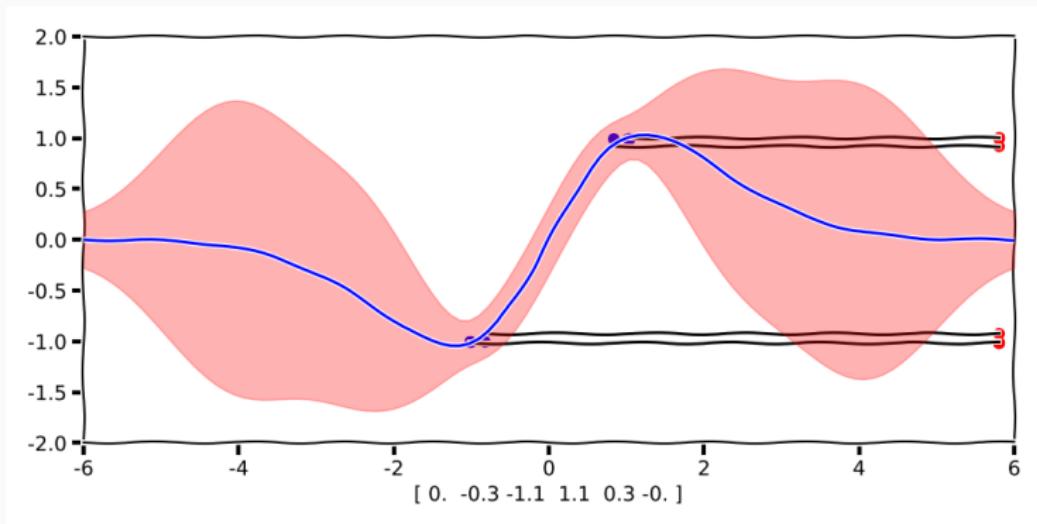
Composite Functions



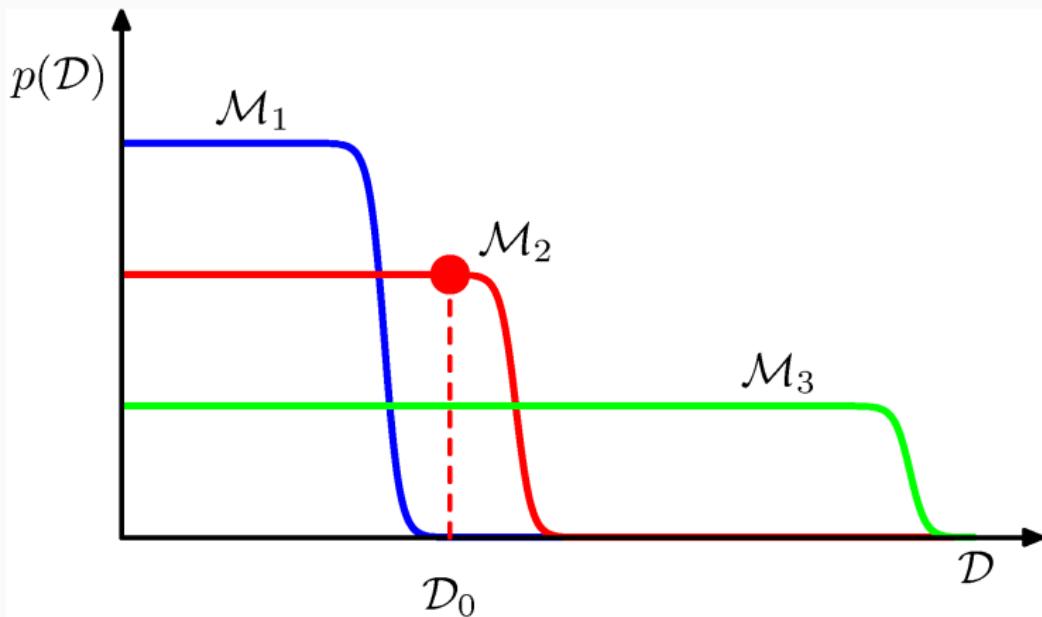
Composite Functions



Composite Functions



MacKay plot



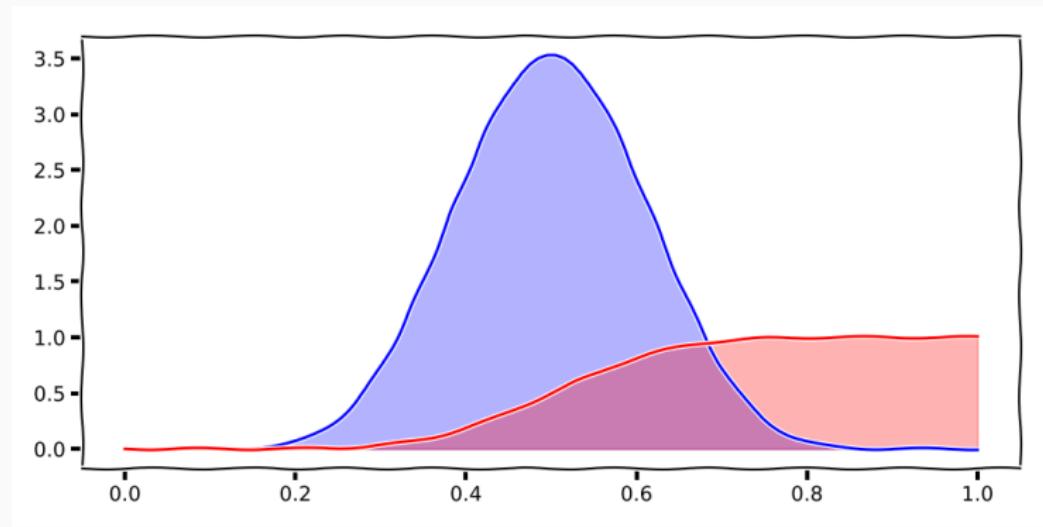
Theorem (Change of Variable)

Let $x \in \mathcal{X} \subseteq \mathbb{R}^n$ be a random vector with a probability density function given by $p_x(x)$, and let $y \in \mathcal{Y} \subseteq \mathbb{R}^n$ be a random vector such that $\psi(y) = x$, where the function $\psi : \mathcal{Y} \rightarrow \mathcal{X}$ is bijective of class of \mathcal{C}^1 and $|\nabla \psi(y)| > 0, \forall y \in \mathcal{Y}$. Then, the probability density function $p_y(\cdot)$ induced in \mathcal{Y} is given by

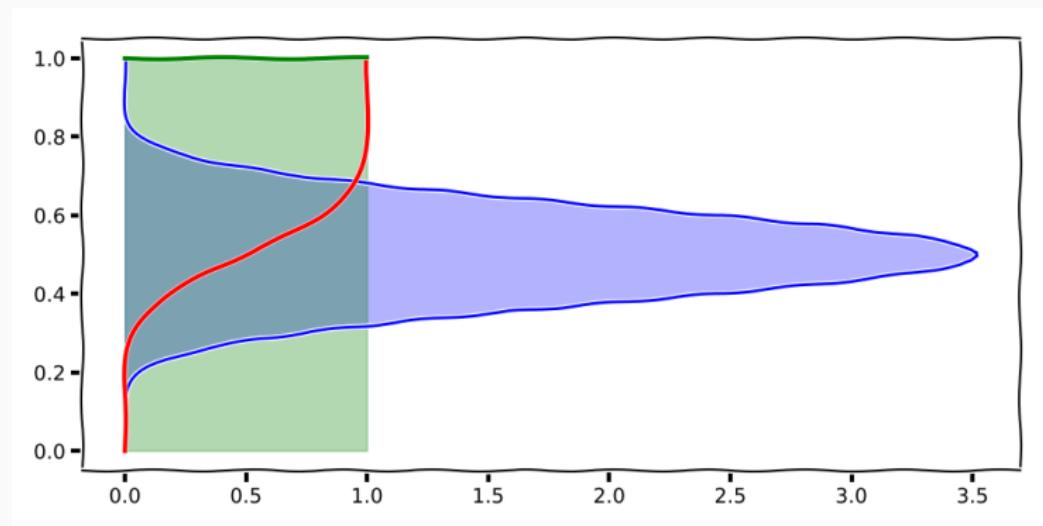
$$p_y(y) = p_x(\psi(y)) |\nabla \psi(y)|$$

where $\nabla \psi(\cdot)$ denotes the Jacobian of $\psi(\cdot)$, and $|\cdot|$ denotes the determinant operator.

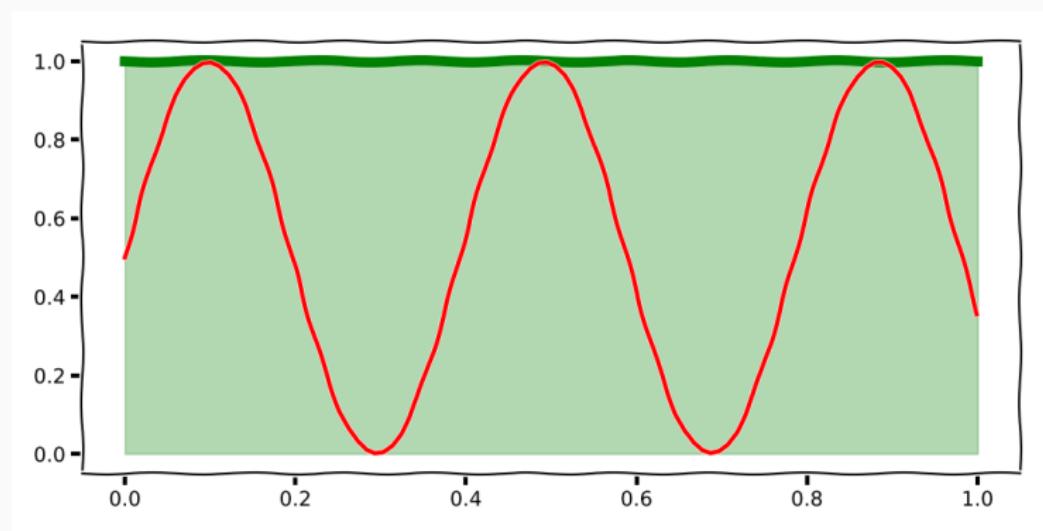
Sampling



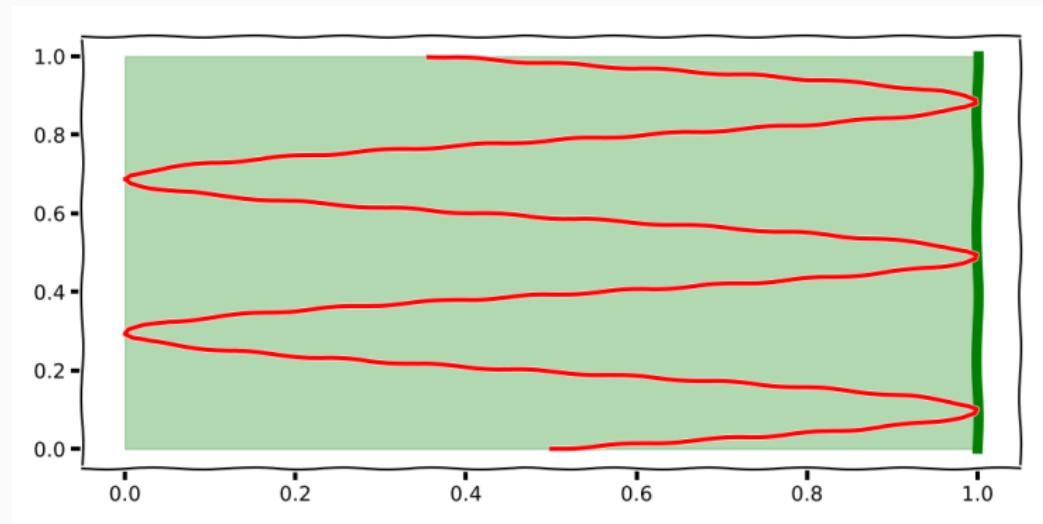
Sampling



Change of Variables



Change of Variables



Learning

Neural Networks

- Specifies a composite function

$$y_i = g(\mathbf{x}_i) = (f_K \circ f_{K-1} \circ \dots \circ f_1)(\mathbf{x}_i)$$

- Traditional form

$$y_k = g(\mathbf{x}, \mathbf{W}) = \sigma \left(\sum_j^M w_{kj}^{(2)} h \left(\sum_j^D w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right)$$

- Linear weights, non-linear activation functions (sigmoid, tanh, etc)

Neural Networks

- Forward propagation

$$g : \mathcal{X} \rightarrow \mathcal{Y}$$

- Forward propagation results in error

$$E(\mathbf{W}) = \frac{1}{2} \sum_i^N (y_i - g(\mathbf{x}_i, \mathbf{W}))^2 + \text{reg}(\mathbf{W})$$

- Propagate error back to update \mathbf{W}

Backpropagation

- "Invented" many times but usually attributed to Rummelhart 1986
- Gradient based update of weights

$$\frac{\partial}{\partial \mathbf{W}} E$$

- Iteratively updated weights

$$\mathbf{W}^{(t)} = \mathbf{W}^{(t-1)} + \eta \frac{\partial}{\partial \mathbf{W}} E,$$

where η is the learning rate

Backpropagation

- Perceptron (no hidden layer)

$$E(\mathbf{W}) = \frac{1}{2} \sum_i \left(y_i - f \left(\sum_j w_{jk} x_j \right) \right)^2$$

- Calculate update on first layer

$$\frac{\partial}{\partial w_{ji}} E(\mathbf{W}) = \frac{\partial}{\partial w_{ik}} \frac{1}{2} \sum_i \left(y_i - f \left(\sum_j w_{jk} x_j \right) \right)^2$$

Backpropagation

- Perceptron (no hidden layer)

$$E(\mathbf{W}) = \frac{1}{2} \sum_i \left(y_i - f \left(\sum_j w_{jk} x_j \right) \right)^2$$

- Calculate update on first layer

$$\begin{aligned} \frac{\partial}{\partial w_{ji}} E(\mathbf{W}) &= \frac{\partial}{\partial w_{ik}} \frac{1}{2} \sum_i \left(y_i - f \left(\sum_j w_{ji} x_j \right) \right)^2 \\ &= \frac{1}{2} \sum_i 2(y_i - \underbrace{g(x_i)}_{\text{output of network}}) \frac{\partial}{\partial w_{ji}} \left(y_i - \sum_j w_{ji} x_j \right) \end{aligned}$$

Backpropagation

- Perceptron (no hidden layer)

$$E(\mathbf{W}) = \frac{1}{2} \sum_i \left(y_i - f \left(\sum_j w_{jk} x_j \right) \right)^2$$

- Calculate update on first layer

$$\begin{aligned} \frac{\partial}{\partial w_{ji}} E(\mathbf{W}) &= \frac{\partial}{\partial w_{ik}} \frac{1}{2} \sum_i \left(y_i - f \left(\sum_j w_{jk} x_j \right) \right)^2 \\ &= \frac{1}{2} \sum_i 2(y_i - \underbrace{g(x_i)}_{\text{output of network}}) \frac{\partial}{\partial w_{ji}} \left(y_i - \sum_j w_{ji} x_j \right) \end{aligned}$$

$$= \sum_i (y_i - g(x_i))(-x_i)$$

Backpropagation: hidden layers

- Chain rule of derivatives

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial t} \frac{\partial t}{\partial x}$$

- Derivatives with respect to hidden layers

$$\frac{\partial}{\partial w_{ji}} E(\mathbf{W}) = \frac{\partial E(\mathbf{W})}{\partial h_i} \frac{\partial h_i}{\partial w_{ji}}$$

Backpropagation: hidden layers

- Chain rule of derivatives

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial t} \frac{\partial t}{\partial x}$$

- Derivatives with respect to hidden layers

$$\frac{\partial}{\partial w_{ji}} E(\mathbf{W}) = \frac{\partial E(\mathbf{W})}{\partial h_i} \frac{\partial h_i}{\partial w_{ji}} = \frac{\partial E(\mathbf{W})}{\partial h_i} \frac{\partial \sum_l w_{li} a_l}{\partial w_{ji}}$$

Backpropagation: hidden layers

- Chain rule of derivatives

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial t} \frac{\partial t}{\partial x}$$

- Derivatives with respect to hidden layers

$$\begin{aligned}\frac{\partial}{\partial w_{ji}} E(\mathbf{W}) &= \frac{\partial E(\mathbf{W})}{\partial h_i} \frac{\partial h_i}{\partial w_{ji}} = \frac{\partial E(\mathbf{W})}{\partial h_i} \frac{\partial \sum_l w_{li} a_l}{\partial w_{ji}} \\ &= \frac{\partial E(\mathbf{W})}{\partial h_i} \sum_l a_l \frac{\partial w_{li}}{\partial w_{ji}}\end{aligned}$$

Backpropagation: hidden layers

- Chain rule of derivatives

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial t} \frac{\partial t}{\partial x}$$

- Derivatives with respect to hidden layers

$$\begin{aligned}\frac{\partial}{\partial w_{ji}} E(\mathbf{W}) &= \frac{\partial E(\mathbf{W})}{\partial h_i} \frac{\partial h_i}{\partial w_{ji}} = \frac{\partial E(\mathbf{W})}{\partial h_i} \frac{\partial \sum_l w_{li} a_l}{\partial w_{ji}} \\ &= \frac{\partial E(\mathbf{W})}{\partial h_i} \sum_l a_l \frac{\partial w_{li}}{\partial w_{ji}} \\ &= \frac{\partial E(\mathbf{W})}{\partial h_i} a_j\end{aligned}$$

Backpropagation

- Now one term remains

$$\frac{\partial}{\partial h_i^{\text{output}}} E(\mathbf{W}) = \frac{\partial E(\mathbf{W})}{\partial f(h_i^{\text{output}})} \frac{\partial f(h_i^{\text{output}})}{\partial h_i^{\text{output}}} = \left\{ f\left(\sum_j w_{ji} a_j^{\text{hidden}}\right) \right\}$$

- The second term is just the derivative of the activation function (more on that later)
- The first term,

$$\frac{\partial E(\mathbf{W})}{\partial f(h_i^{\text{output}})} = \frac{\partial}{\partial f(h_i^{\text{output}})} \left(\frac{1}{2} \sum_i (f(h_i^{\text{output}}) - y_i))^2 \right)$$

Backpropagation

- Now one term remains

$$\frac{\partial}{\partial h_i^{\text{output}}} E(\mathbf{W}) = \frac{\partial E(\mathbf{W})}{\partial f(h_i^{\text{output}})} \frac{\partial f(h_i^{\text{output}})}{\partial h_i^{\text{output}}} = \left\{ f\left(\sum_j w_{ji} a_j^{\text{hidden}}\right) \right\}$$

- The second term is just the derivative of the activation function (more on that later)
- The first term,

$$\begin{aligned} \frac{\partial E(\mathbf{W})}{\partial f(h_i^{\text{output}})} &= \frac{\partial}{\partial f(h_i^{\text{output}})} \left(\frac{1}{2} \sum_i (f(h_i^{\text{output}}) - y_i)^2 \right) \\ &= (f(h_i^{\text{output}}) - y_i) \end{aligned}$$

Backpropagation

- Nearly there

$$\sum_i \frac{\partial E(\mathbf{W})}{\partial h_i^{\text{output}}} \frac{\partial h_i^{\text{output}}}{\partial h_j^{\text{hidden}}}$$

- Hidden units

$$h_i^{\text{output}} = f \left(\sum_l w_{li} h_l^{\text{hidden}} \right)$$

- Derivatives

$$\frac{\partial h_i^{\text{output}}}{\partial h_j^{\text{hidden}}} = \frac{\partial f \left(\sum_l w_{li} h_l^{\text{hidden}} \right)}{\partial h_j^{\text{hidden}}} = w_{ji} f'(a_i)$$

Backpropagation

- Putting it all together
- First layer

$$w_{ji}^{(t)} = w_{ji}^{(t-1)} - \eta \frac{\partial E(\mathbf{W})}{\partial w_{ji}} = w_{ji}^{(t-1)} - \eta (f(h_i^{\text{output}})) f'(h_i^{\text{output}}) a_j$$

- Hidden layer

$$\begin{aligned} w_{ji}^{(t)} &= w_{ji}^{(t-1)} - \eta \frac{\partial E(\mathbf{W})}{\partial w_{ji}} \\ &= w_{ji}^{(t-1)} - \eta f'(a_i) \sum_k ((f(h_i^{\text{output}})) - y) w_{ik} x_j \end{aligned}$$

Backpropagation

1. Randomly initialise weights
2. Forward pass (get output error)
3. Compute gradients for one step back
4. Iteratively push gradients back
5. Update each layer
6. Run till convergence (or well...)

Activation functions

- Sigmoid

$$f(x) = \frac{1}{1 + e^{-x}}$$

- Tanh

$$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$$

- ReLu

$$f(x) = \max(0, x) \approx \log(1 + e^x)$$



- Tensorflow, OpenMX, Theano, Torch, etc.
- Very useful when building composite functions

$$\log(\det(K))$$

- Compute log determinant of covariance matrix

$$\log(\det(K))$$

- Compute log determinant of covariance matrix
- Covariance matrix $K \succ 0 \rightarrow$ Cholesky

$$K = L^T L$$

$$\log(\det(K))$$

- Compute log determinant of covariance matrix
- Covariance matrix $K \succ 0 \rightarrow$ Cholesky

$$K = L^T L$$

- Simplification

$$\log(|K|) = \log(|L^T L|) = \log(|L^T| \cdot |L|) = 2 \cdot \log|L|$$

$$\log(\det(K))$$

- Compute log determinant of covariance matrix
- Covariance matrix $K \succ 0 \rightarrow$ Cholesky

$$K = L^T L$$

- Simplification

$$\log(|K|) = \log(|L^T L|) = \log(|L^T| \cdot |L|) = 2 \cdot \log|L|$$

- L triangular matrix $|A| = \prod \text{diag}(A)$

$$\log(|K|) = 2 \cdot \log\left(\prod \text{diag}(L)\right) = 2 \sum \log(\text{diag}(L))$$

Architectures

Architectures

- Backpropagation allows for updates of any differentiable composite function
- With auto-diff this means that you can **hack** together lots of structures
- If we know a-priori what we want to learn we can reduce the search space
- We have to specify everything in terms of hard constraints

Convolutional Neural Networks

- Designed by Yann Lecun late eighties, did well on some tasks but really took off when data grew
- Includes two special types of transformations
 - Convolutions: linear transformation
 - Pooling: non-linear transformation
- Highly designed structures

Convolutional Neural Networks

- Convolution

$$s(x) = \int x(t)w(x-t)dt$$

- A convolution takes a weighted average within a spatial region (eg. blurring)

$$S_{ij} = \sum_m \sum_n I(i-m, j-n)K(m, n)$$

- Makes a lot of sense for images, average nearby pixels

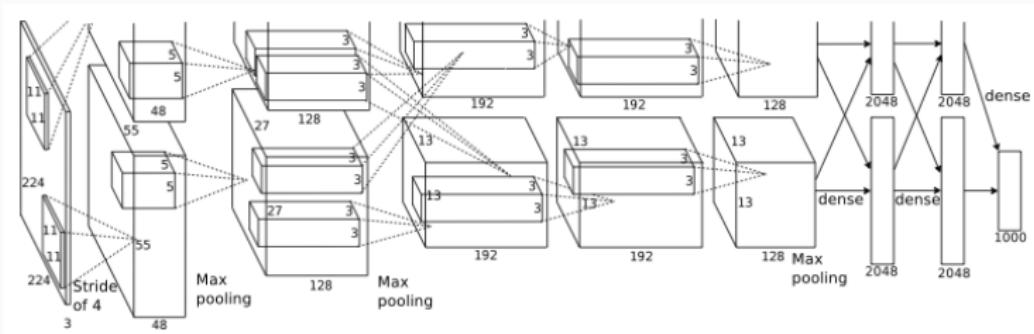
Convolutional Neural Networks

- Pooling

$$S_{ij} = \max\{x_{ij}\}_{i=1,j=1}^{i=N,j=M}$$

- Introduces invariances
 - translation
 - scaling
- *it is more important that something is present than where it is present*

Convolutional Neural Networks



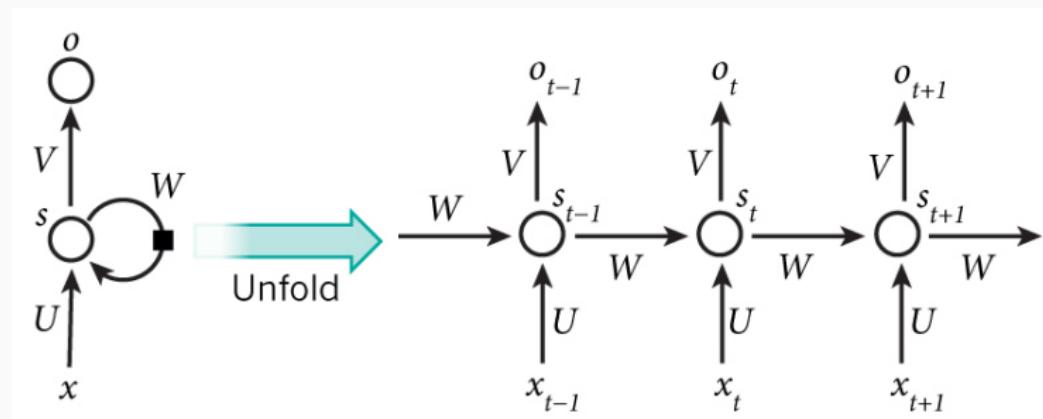
Recurrent Neural Networks

- Recurrent Neural Networks designed to model sequential data
- Dynamical system

$$s^{(t)} = f(s^{(t-1)}, x^{(t)})$$

- Idea:
 - unfold the computation into a single graph
 - becomes a standard neural network

Recurrent Neural Networks²



²<http://www.wildml.com/2015/09/>

recurrent-neural-networks-tutorial-part-1-introduction-to-rnns/

Architectures

- There exists many many different architectures
- Few are derived or justified on any scientific basis
- There are many post-justifications³
- They do work really well but don't listen to the pseudo-science
dare to have your own intuitions instead

³Chomsky, N. A., & Fodor, J. A. (1980). The inductivist fallacy.

The problem with lack of understanding

Decision vs. Understanding

- Neural networks, they way they are described are discriminative this means

Decision vs. Understanding

- Neural networks, they way they are described are discriminative this means
 - they do not model the data

Decision vs. Understanding

- Neural networks, they way they are described are discriminative this means
 - they do not model the data
 - they do not understand the data

Decision vs. Understanding

- Neural networks, they way they are described are discriminative this means
 - they do not model the data
 - they do not understand the data
 - there is no such thing as uncertainty in inputs

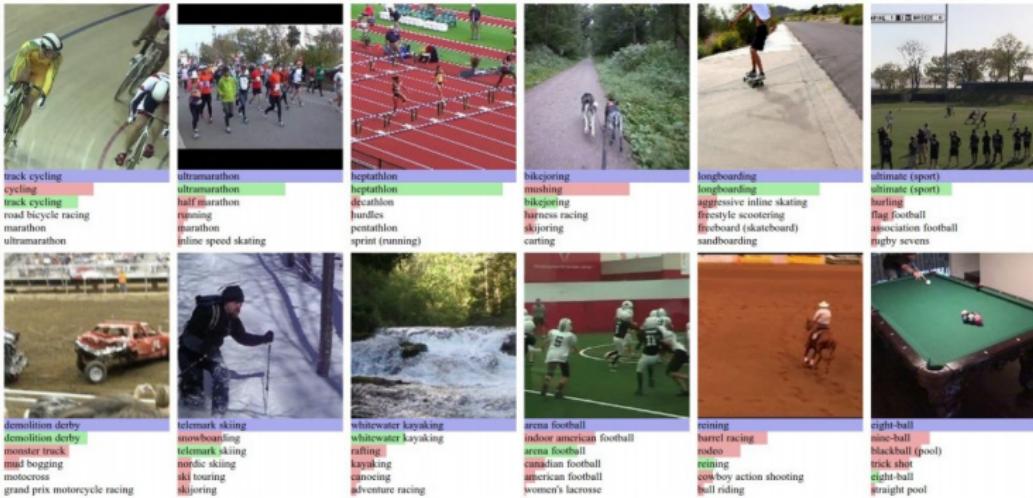
Decision vs. Understanding

- Neural networks, they way they are described are discriminative this means
 - they do not model the data
 - they do not understand the data
 - there is no such thing as uncertainty in inputs
- They are not models they are decision machines

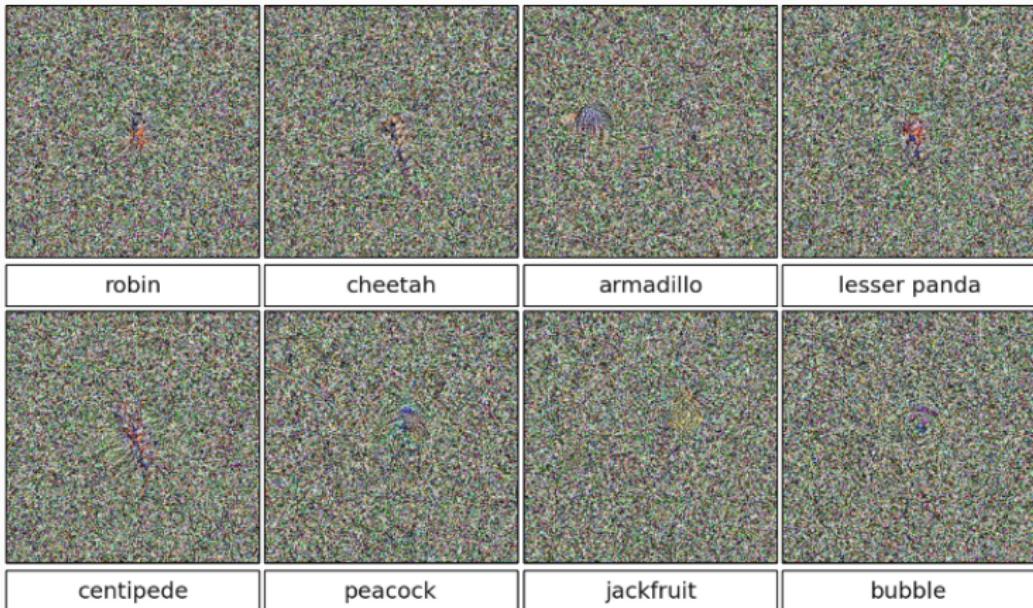
Decision vs. Understanding

- Neural networks, they way they are described are discriminative this means
 - they do not model the data
 - they do not understand the data
 - there is no such thing as uncertainty in inputs
- They are not models they are decision machines
- *Like only revising old exams*

ImageNet

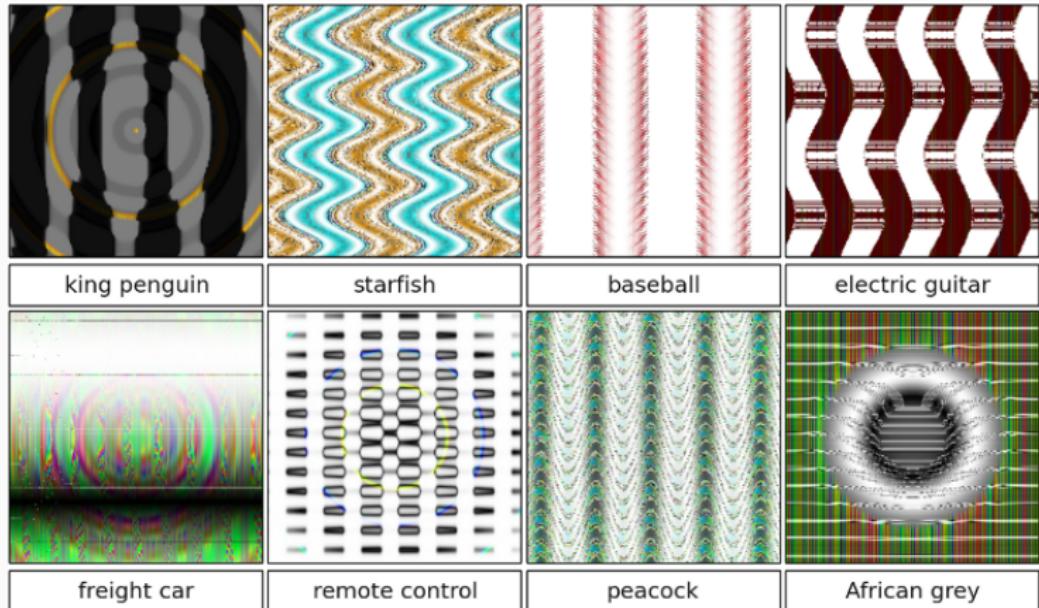


Fools⁴



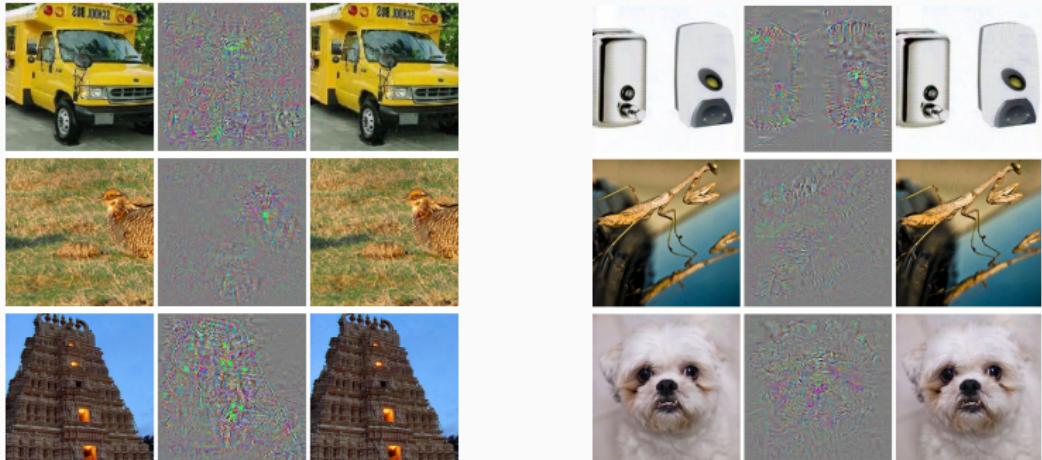
⁴http://www.evolvingai.org/files/DNNsEasilyFooled_cvpr15.pdf

Fools⁴



⁴http://www.evolvingai.org/files/DNNsEasilyFooled_cvpr15.pdf

Fools⁴



Left correctly classified, all right images are classified as "ostrich"

⁴http://www.evolvingai.org/files/DNNsEasilyFooled_cvpr15.pdf

The Black-arts

- Early stopping
- Layer-wise training
- Denoising
- Adveserial training
- Dropout
- Etc.

Decision machines

- Decision machines are useful
 - logistic regression, Support Vector Machines
 - nearest neighbours, etc.
- These methods do not warp the input space to the same extent
 - do not mess with proximity to the same extent
- You need enormous amounts of data to make sure you retain representation of space

Summary

Summary

- Neural networks is nothing special, forget the neuroscience, it has nothing to do with it

Summary

- Neural networks is nothing special, forget the neuroscience, it has nothing to do with it
- Composite functions **cannot** model more things

Summary

- Neural networks is nothing special, forget the neuroscience, it has nothing to do with it
- Composite functions **cannot** model more things
- However,

Summary

- Neural networks is nothing special, forget the neuroscience, it has nothing to do with it
- Composite functions **cannot** model more things
- However,
 - they can easily warp the input space to model **less** things

Summary

- Neural networks is nothing special, forget the neuroscience, it has nothing to do with it
- Composite functions **cannot** model more things
- However,
 - they can easily warp the input space to model **less** things
 - they are very simple and parallel

Summary

- Neural networks is nothing special, forget the neuroscience, it has nothing to do with it
- Composite functions **cannot** model more things
- However,
 - they can easily warp the input space to model **less** things
 - they are very simple and parallel
 - computationally trivial

Summary

- Neural networks is nothing special, forget the neuroscience, it has nothing to do with it
- Composite functions **cannot** model more things
- However,
 - they can easily warp the input space to model **less** things
 - they are very simple and parallel
 - computationally trivial
- However,

Summary

- Neural networks is nothing special, forget the neuroscience, it has nothing to do with it
- Composite functions **cannot** model more things
- However,
 - they can easily warp the input space to model **less** things
 - they are very simple and parallel
 - computationally trivial
- However,
 - lacks principle so hard to know what they do or how to do new things

Summary

- Neural networks is nothing special, forget the neuroscience, it has nothing to do with it
- Composite functions **cannot** model more things
- However,
 - they can easily warp the input space to model **less** things
 - they are very simple and parallel
 - computationally trivial
- However,
 - lacks principle so hard to know what they do or how to do new things
- With sufficient data we do not need to "learn" look-up is sufficient

eof

References
