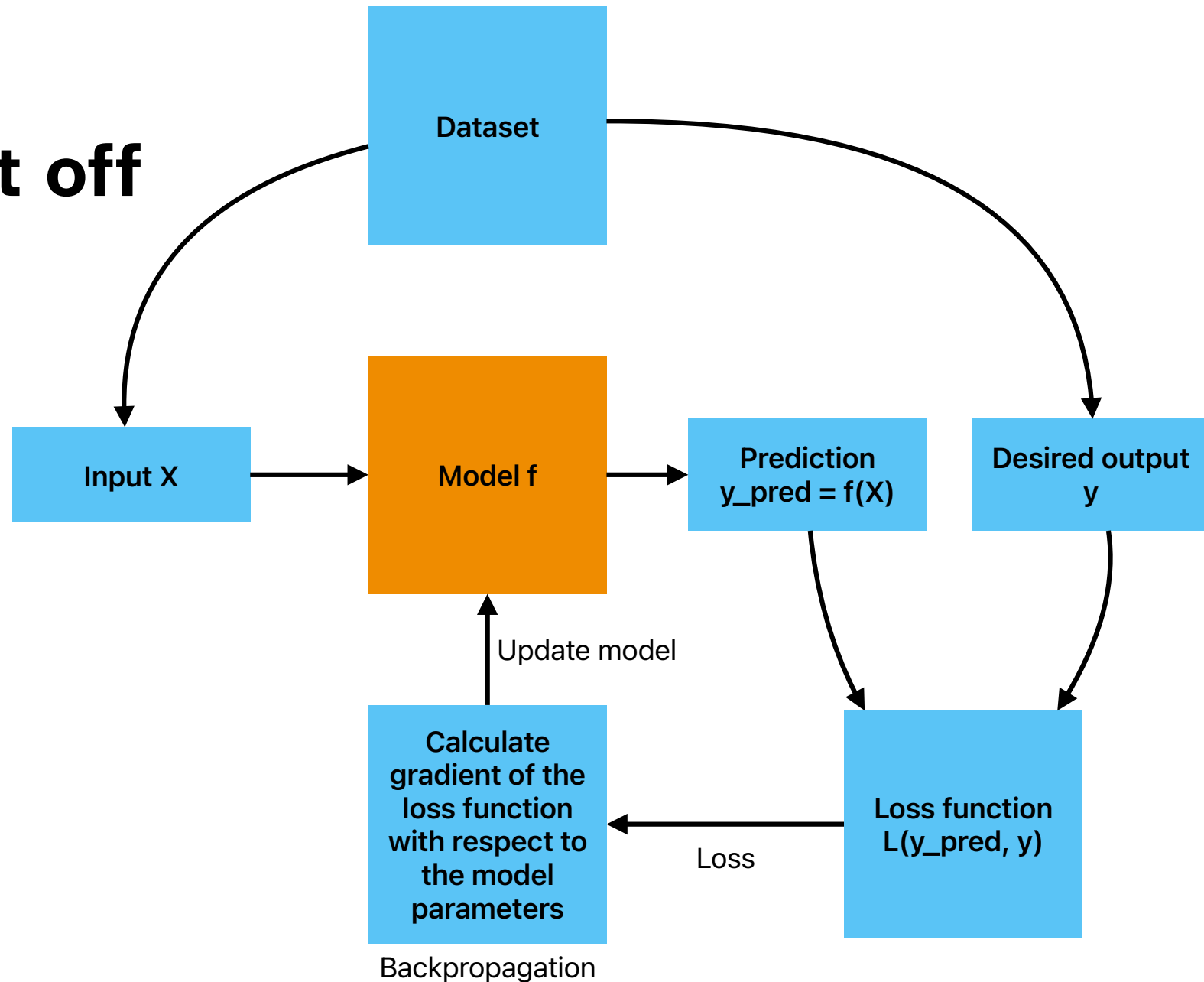


Advanced Machine Learning

Deep Feedforward Networks

Lukas Galke Poech • Spring 2026

Where we left off

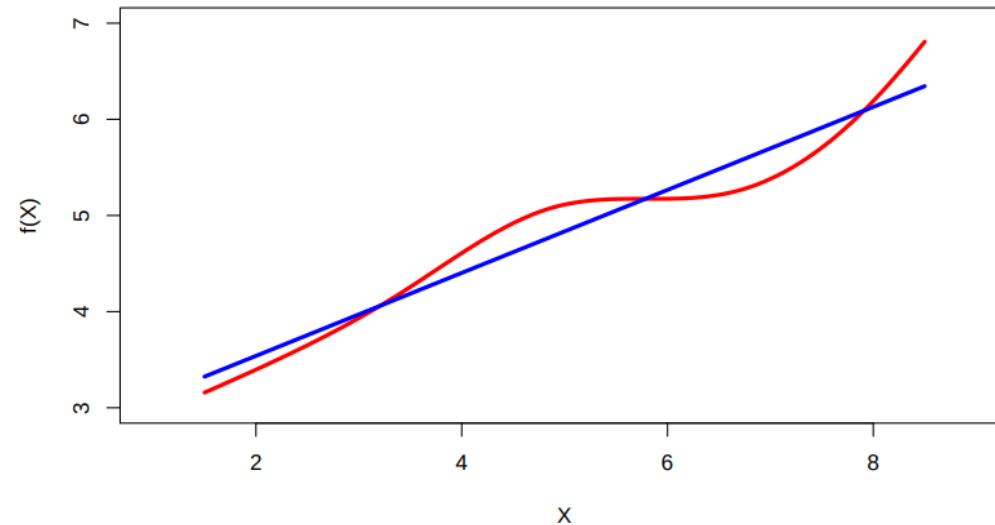


Keywords

- * **Single-Layer Networks**
- * Deep Neural Networks
- * Universal Approximation Theorem
- * Activation functions
- * Output units

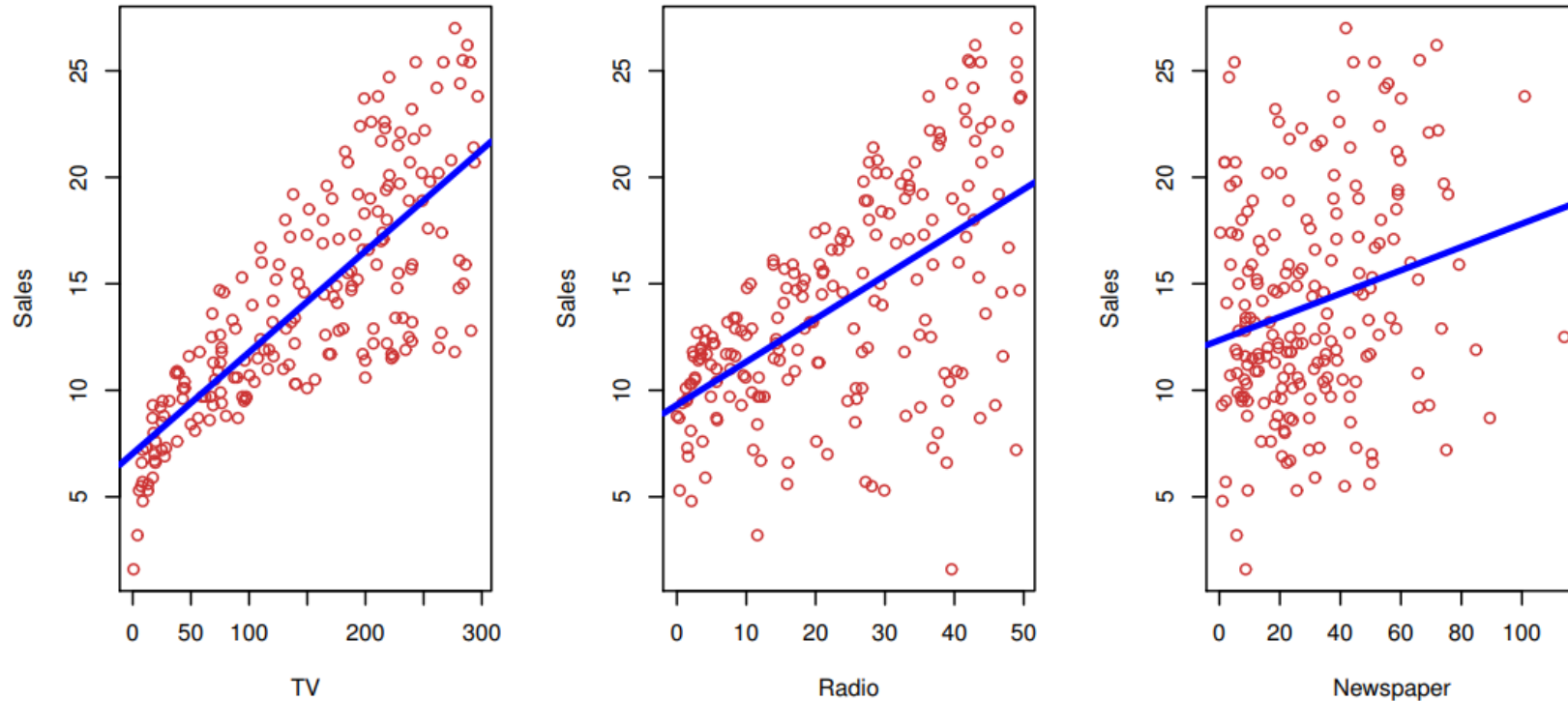
Linear Regression

- Linear regression is a simple approach to supervised learning. It assumes that the dependence of Y on X_1, X_2, \dots, X_p is linear.
- True regression functions are never linear!



- Although it may seem overly simplistic, linear regression is extremely useful both conceptually and practically.

Example: Advertising data



- The Advertising data set consists of the sales of that product in 200 different markets, along with advertising budgets for the product in each of those markets for three different media: TV, radio, and newspaper.

Example: Advertising data

Questions we might ask:

- Is there a relationship between advertising budget and sales?
- How strong is the relationship between advertising budget and sales?
- Which media contribute to sales?
- How accurately can we predict future sales?
- Is the relationship linear?
- Is there synergy among the advertising media?

Simple Linear Regression.

We assume a model

$$Y = \beta_0 + \beta_1 X + \epsilon$$

where β_0 and β_1 are two unknown constants that represent the intercept and slope, also known as coefficients or parameters, and ϵ is the error term.

Given some estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ model coefficients, we **predict** future sales using

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

Estimation of the Parameters by Least Squares

- Let $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ be the prediction of y for the i th observation.
- Then the **residual** is defined as

$$e_i = y_i - \hat{y}_i$$

- We define the **residual sum of squares** as

$$\begin{aligned} \text{RSS} &= e_1^2 + e_2^2 + \dots + e_n^2 = \\ &= \left(y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1\right)^2 + \left(y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2\right)^2 + \dots + \left(y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n\right)^2 \end{aligned}$$

- Our task is to find the $\hat{\beta}_0$ and $\hat{\beta}_1$ minimizing the RSS.

Estimation of the Parameters by Least Squares

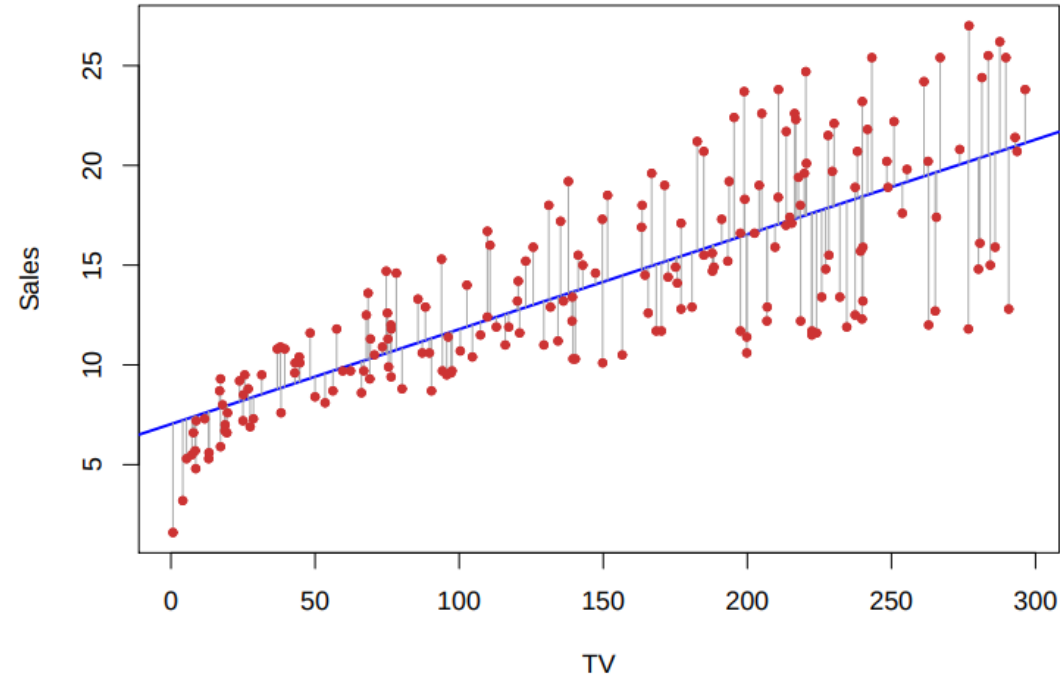
➤ With simple derivation of the RSS formula, we can show that

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

➤ with the sample means $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ and $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Example: Advertising data



- The least squares fit for the regression of sales onto TV.
- In this case a linear fit captures the essence of the relationship, although it is somewhat deficient in the left of the plot.

Multiple Linear Regression

- We can include several features or predictors into our model

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \epsilon$$

- We interpret β_j as the average effect on Y of a one unit increase in X_j , holding all other predictors fixed. In the advertising example, the model becomes

$$\text{sales} = \beta_0 + \beta_1 \cdot \text{TV} + \beta_2 \cdot \text{radio} + \beta_3 \cdot \text{newspaper} + \epsilon$$

In Matrix Form

We can formulate the entire linear regression also in Matrixform

$$y = Xb + e$$

$y = (y_1, \dots, y_n) \in \mathbb{R}^n$ is the $n \times 1$ response vector

$X = [1_n, X'] \in \mathbb{R}^{n \times (p+1)}$ is the $n \times (p+1)$ design matrix

1_n is an $n \times 1$ vector of ones

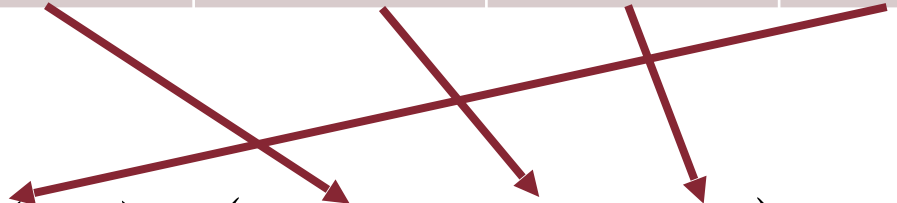
$X' = [X_1, X_2, \dots, X_p] \in \mathbb{R}^{n \times p}$ the $n \times p$ predictor Matrix

$b = (\beta_0, \beta_1, \dots, \beta_p) \in \mathbb{R}^{p+1}$ is regression coefficient vector

$e = (e_1, e_2, \dots, e_n) \in \mathbb{R}^n$ is the $n \times 1$ error vector

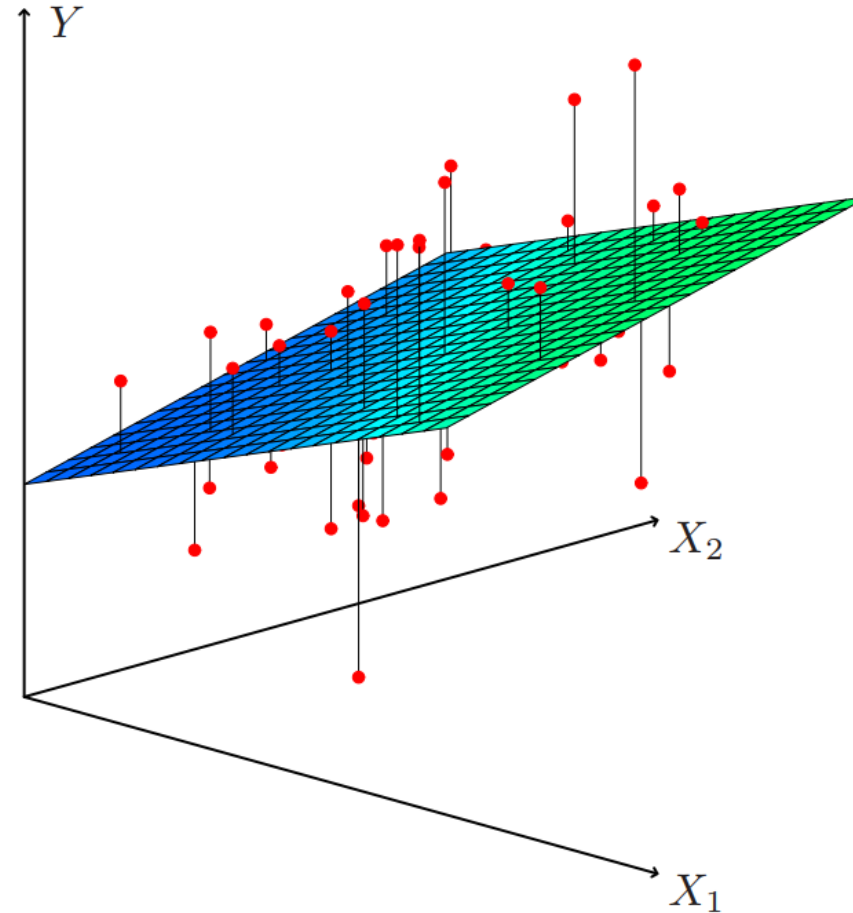
Example

#	TV	Radio	News	Sales
1	230.1	37.8	69.2	22.1
2	44.5	39.3	45.1	10.4
3	17.2	45.9	69.3	9.3
4	151.5	41.3	58.5	18.5
5	180.8	10.8	58.4	12.9



$$\begin{pmatrix} 22.1 \\ 10.4 \\ 9.3 \\ 18.5 \\ 12.9 \end{pmatrix} = \begin{pmatrix} 1 & 230.1 & 37.8 & 69.2 \\ 1 & 44.5 & 39.3 & 45.1 \\ 1 & 17.2 & 45.9 & 69.3 \\ 1 & 151.5 & 41.3 & 58.5 \\ 1 & 180.8 & 10.8 & 58.4 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \\ e_5 \end{pmatrix}$$

Linear Regression



Generalizations

There exist many different extensions and generalizations to the simple linear regression model:

- * **Classification problems:** logistic regression, support vector machines
- * **Non-linearity:** kernel smoothing, splines and generalized additive models; nearest neighbor methods.
- * **Interactions:** Tree-based methods, bagging, random forests and boosting (these also capture non-linearities)
- * **Regularized fitting:** Ridge regression and lasso

Logistic Regression

One type of generalized linear model

- * **Logistic regression** models the log-odds of an event as a linear combination of input variables
- * How to get there? Apply sigmoid to the output to ensure it is between 0 and 1, estimating the probability of an event (rather than a quantity)

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

Regression $p(X) = \beta_0 + \beta_1 X.$

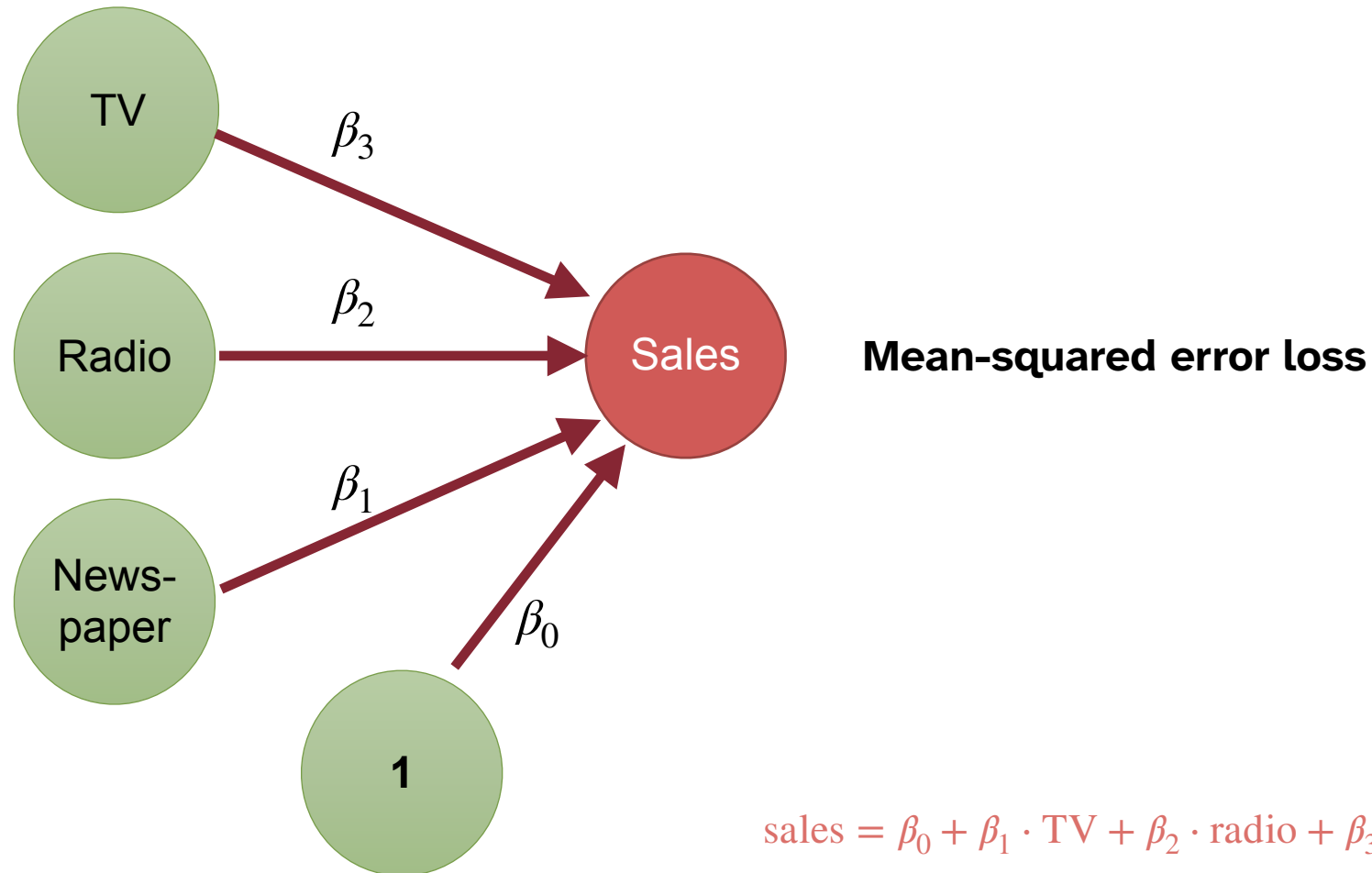
Apply sigmoid to RHS $p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}}.$

Rewrite to "Odds" $\frac{p(X)}{1 - p(X)} = e^{\beta_0 + \beta_1 X}.$

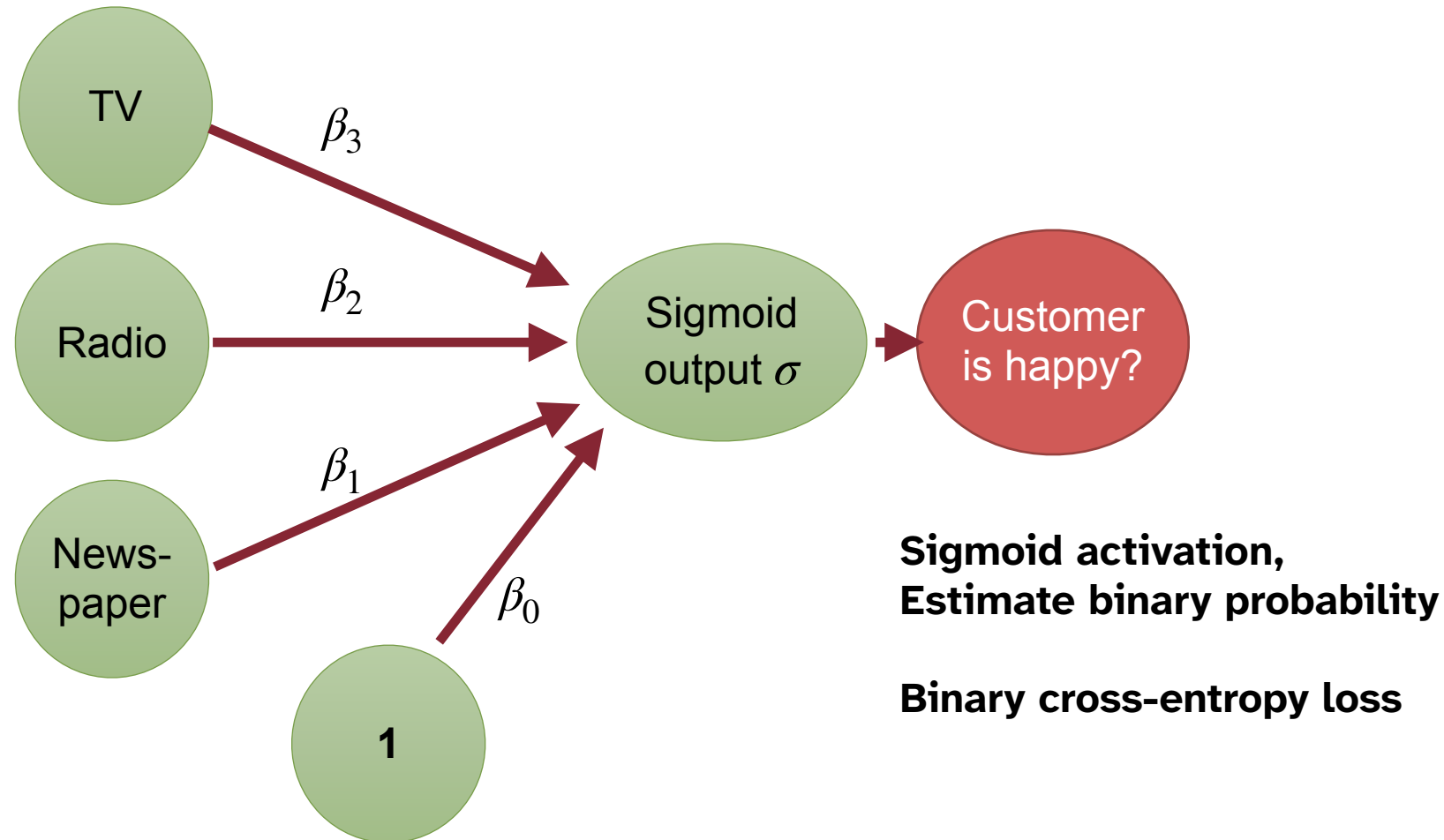
Log odds or "Logits" $\log \left(\frac{p(X)}{1 - p(X)} \right) = \beta_0 + \beta_1 X.$

This terminology has been recycled for the raw outputs of a neural network

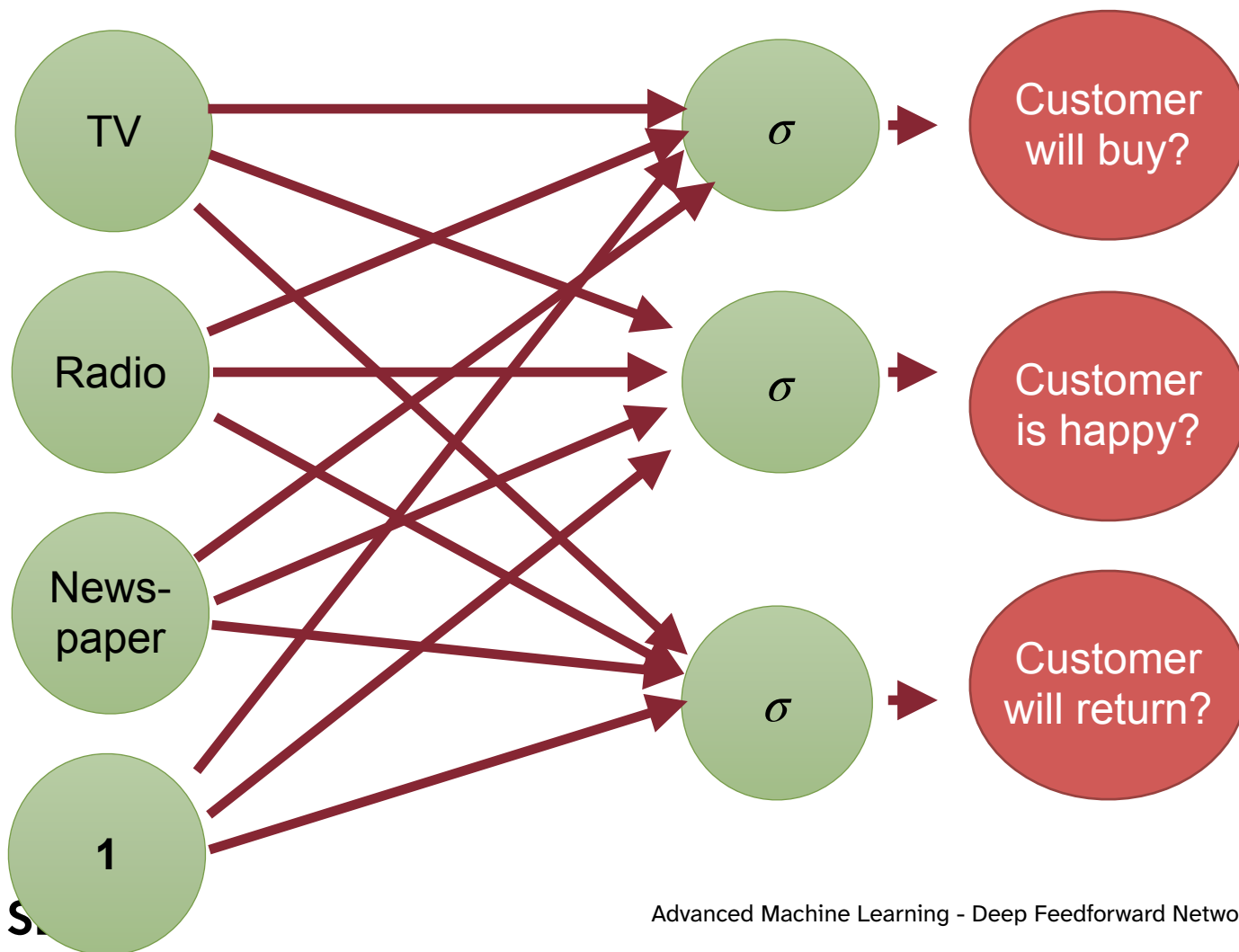
Single-layer neural network for linear regression



Single-layer neural network for logistic regression



Single-layer NN doing multiple logistic regressions



Here, all the logistic regressions are independent from each other

Multi-layered networks can find shared factors of variation, similar to the interactions seen before

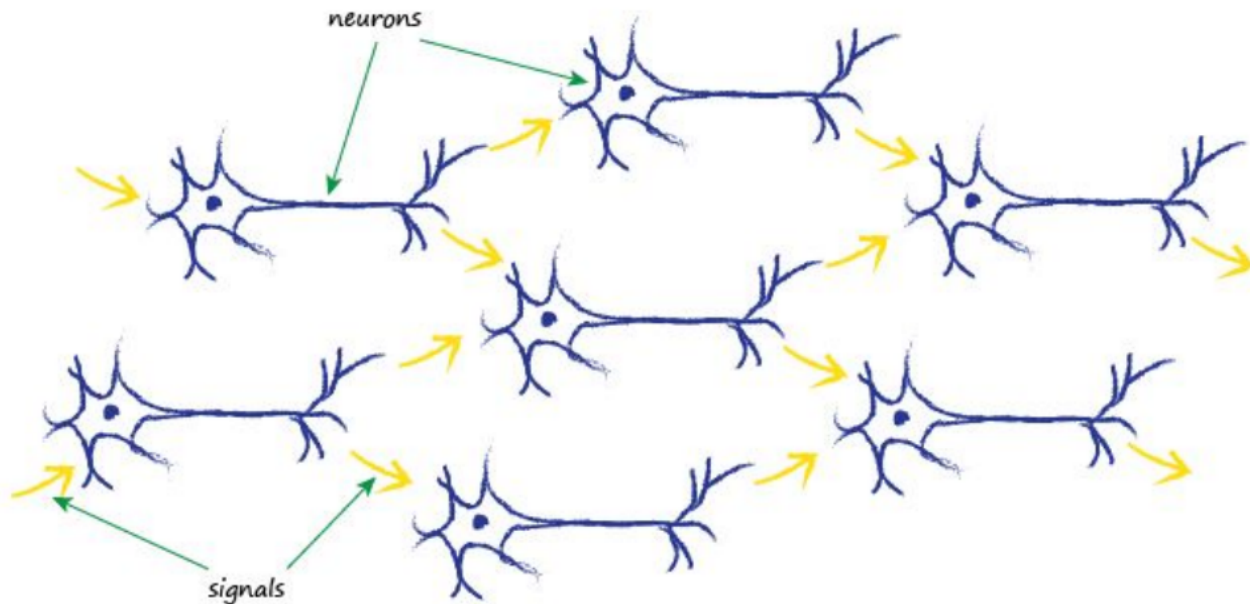
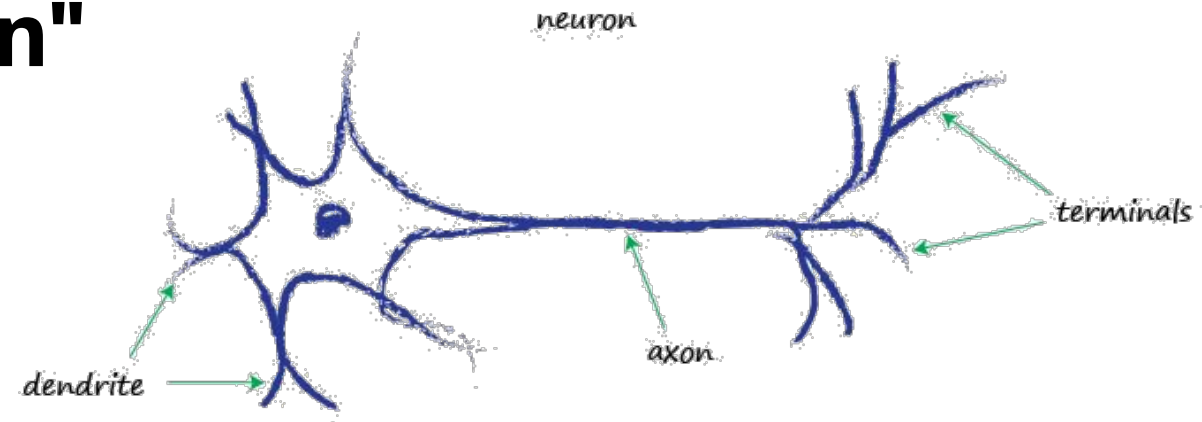
More scalable than modeling explicit interaction for every pair of variables

But less interpretable

Keywords

- * Single-Layer Networks
- * **Deep Neural Networks**
- * Universal Approximation Theorem
- * Activation functions
- * Output units

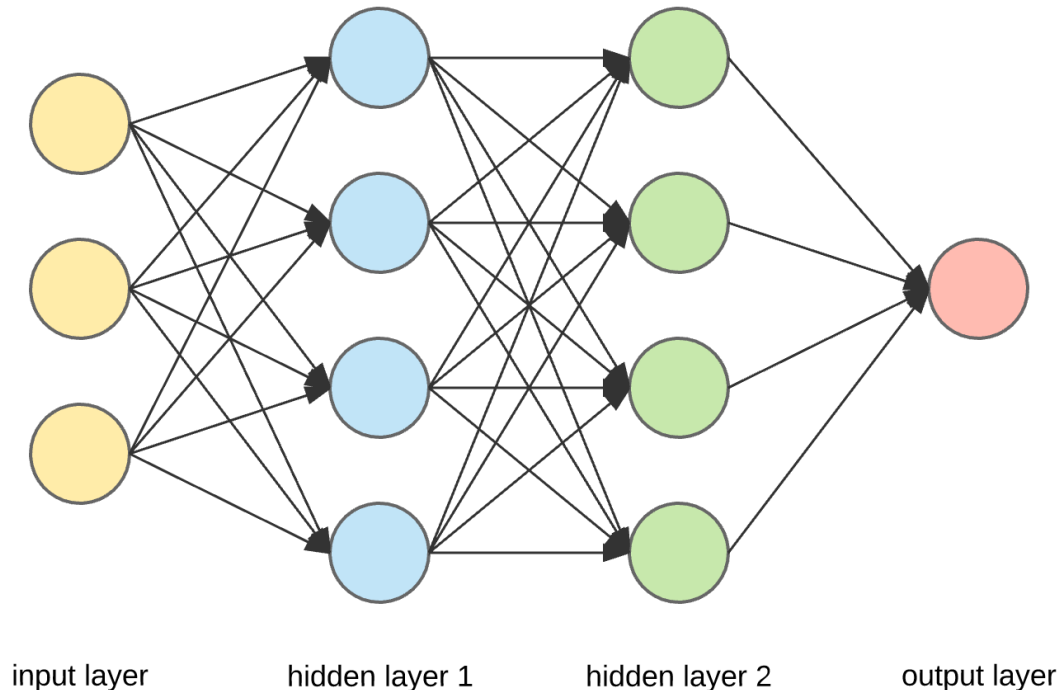
From Neuron to "Brain"



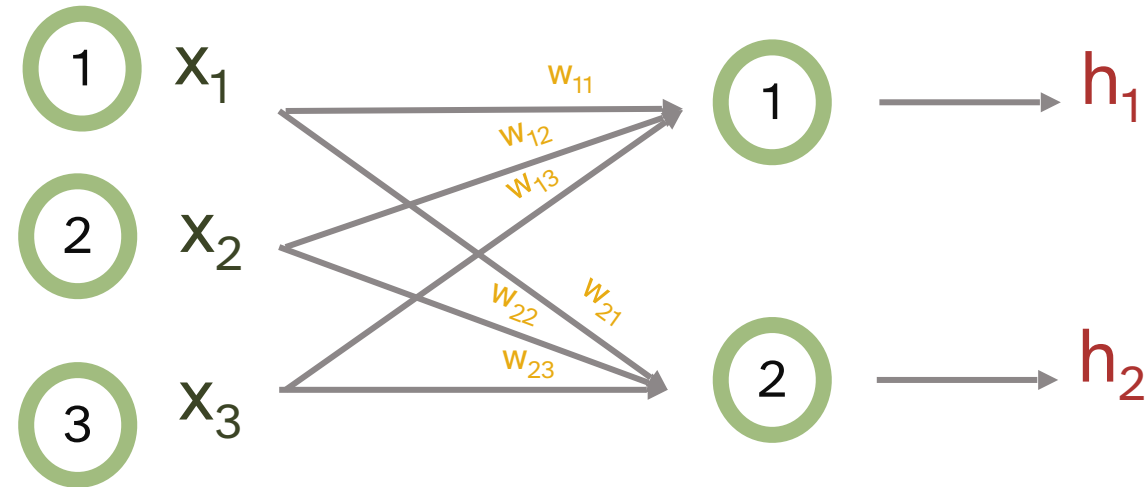
- * We now can connect neurons
- * The output of one neuron becomes the input of other neurons

In a more structured way

- * We normally have more than one node
- * Multiple Nodes are arranged in layers
- * Each layer receives the generated output from the previous layer



What does such a Network do?

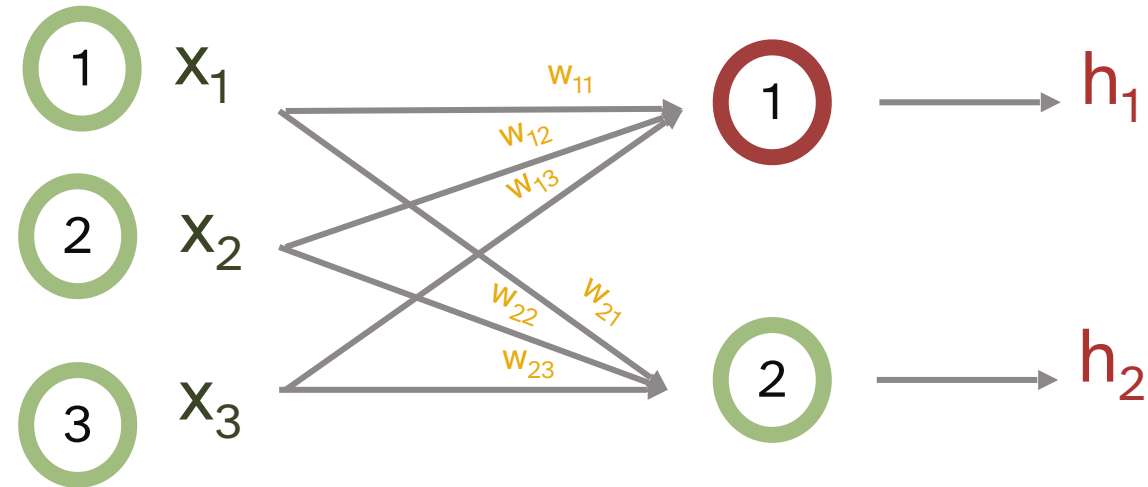


* Remember from before:

* We have a weighted sum

* Then we have an "activation function" (which we ignore for now)

What does such a Network do?



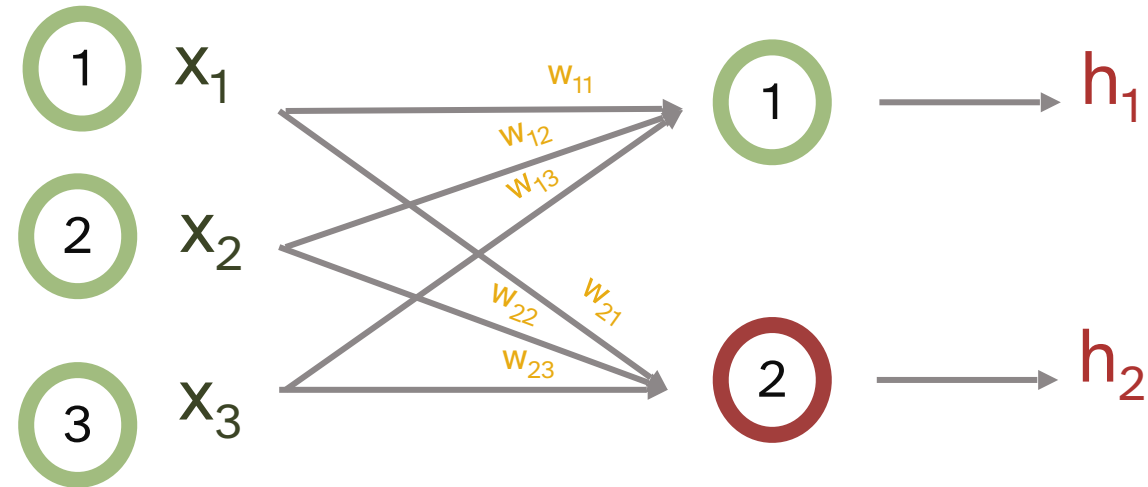
* Remember from before:

* We have a weighted sum

* Then we have an "activation function" (which we ignore for now)

* For **Node 1**, we get: $h_1 = w_{11} \cdot x_1 + w_{12} \cdot x_2 + w_{13} \cdot x_3$

What does such a Network do?



* Remember from before:

* We have a weighted sum

* Then we have an "activation function" (which we ignore for now)

* For **Node 2**, we get: $h_2 = w_{21} \cdot x_1 + w_{22} \cdot x_2 + w_{23} \cdot x_3$

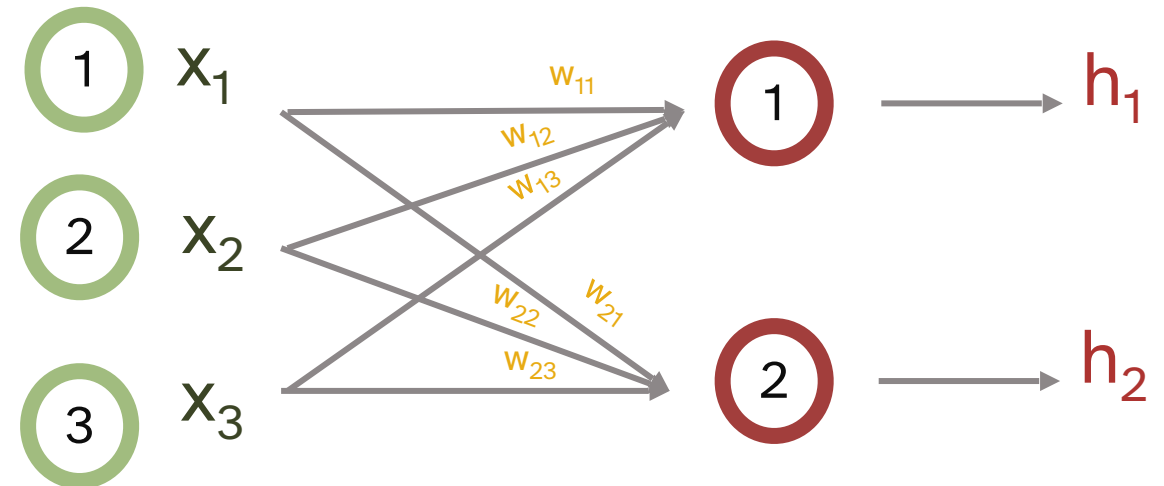
What does such a Network do?

Together:

$$h_1 = w_{11} \cdot x_1 + w_{12} \cdot x_2 + w_{13} \cdot x_3$$

$$h_2 = w_{21} \cdot x_1 + w_{22} \cdot x_2 + w_{23} \cdot x_3$$

Looks familiar, doesn't it?

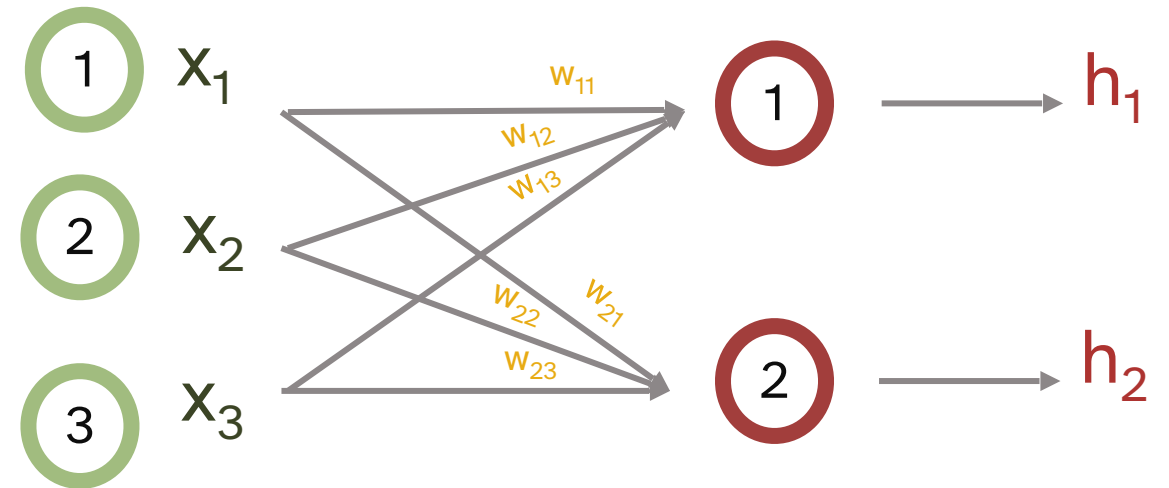


What does such a Network do?

Rewrite:

$$\begin{pmatrix} h_1 \\ h_2 \end{pmatrix} = \begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

$$h = Wx$$

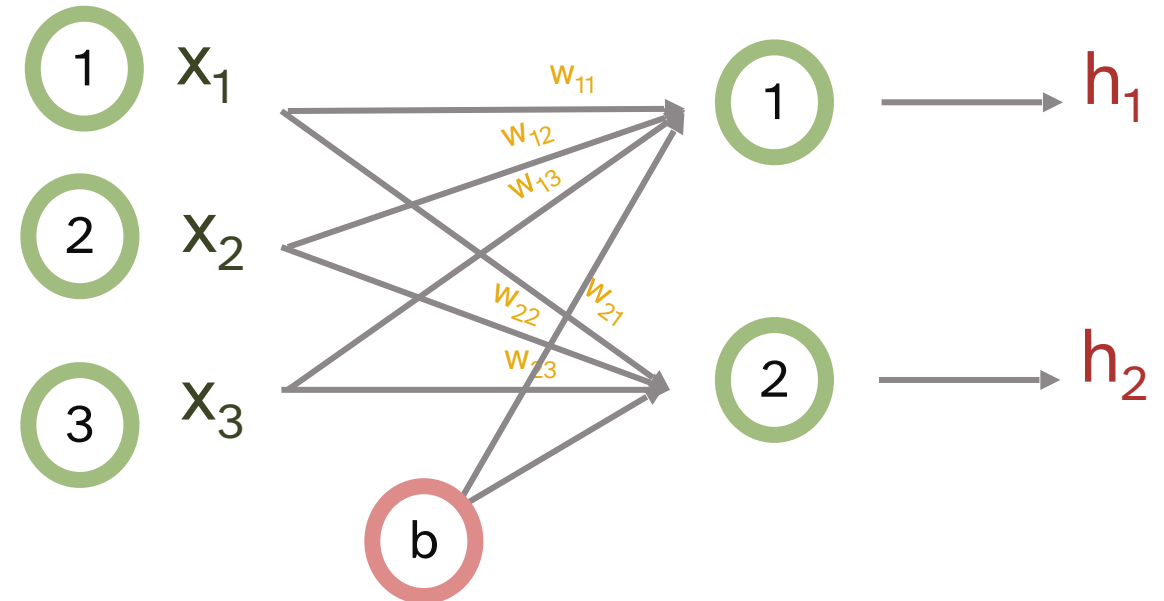


What does such a Network do?

Adding a **Bias**:

$$\begin{pmatrix} h_1 \\ h_2 \end{pmatrix} = \begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}$$

$$h = Wx + b$$

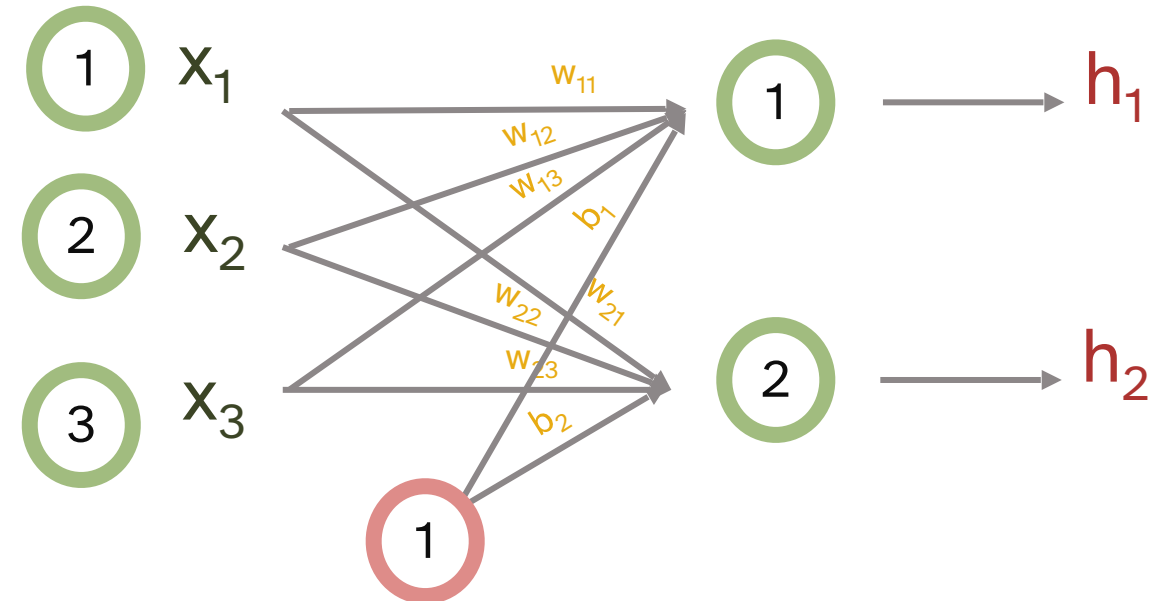


What does such a Network do?

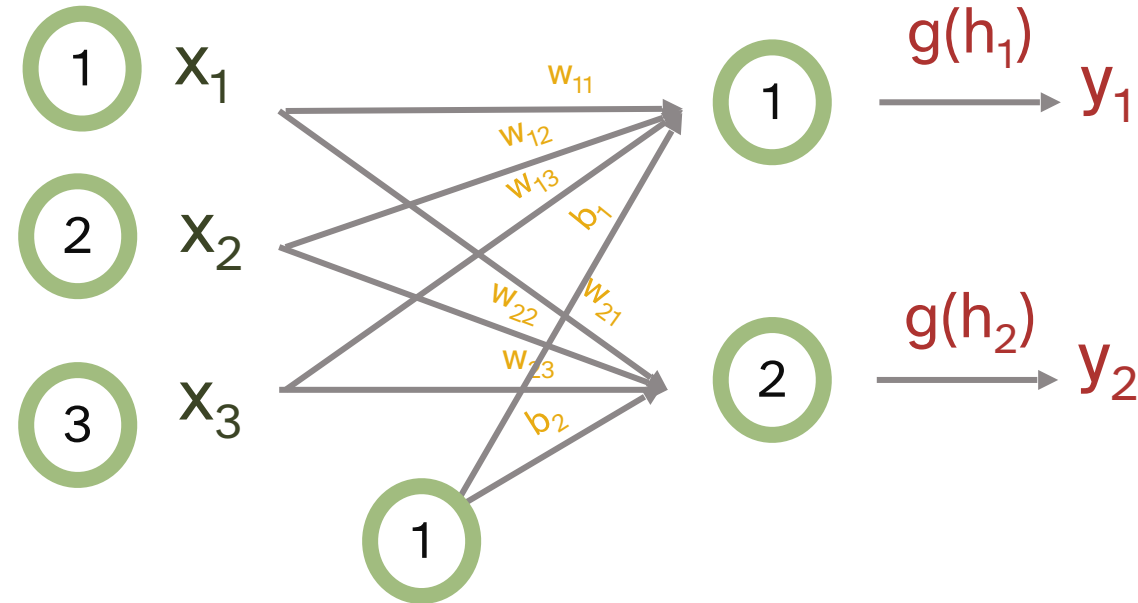
Adding a **Bias** (rewritten):

$$\begin{pmatrix} h_1 \\ h_2 \end{pmatrix} = \begin{pmatrix} w_{11} & w_{12} & w_{13} & b_1 \\ w_{21} & w_{22} & w_{23} & b_2 \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ 1 \end{pmatrix}$$

$$h = Wx$$



What does such a Network do?

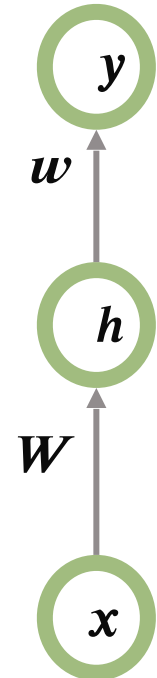


Adding the **activation function**:

$$y = g(Wx + b)$$

More Compact Representation

- * We can summarize these networks to their essential parts: We receive the vector \mathbf{x} and perform an affine transformation according to the weight matrix
- * $\mathbf{h}' = \mathbf{W}\mathbf{x} + \mathbf{b}$
- * Afterward, the internal representation is component wise feed through an activation function $h_i = g(h'_i)$ to generate the input vector \mathbf{h} for the next layer
- * Please note the slight shift in notation. From now on:
 - * \mathbf{x} denotes the input
 - * \mathbf{W} denotes weight matrices, \mathbf{w} denotes weight vectors
 - * \mathbf{h} denotes hidden layer "outputs"
 - * \mathbf{y} is the final output of the network



Composition of Complex Functions

Now let's look at a deeper Network:

E.g., functions $f^{(1)}$, $f^{(2)}$, $f^{(3)}$ connected in a chain to form

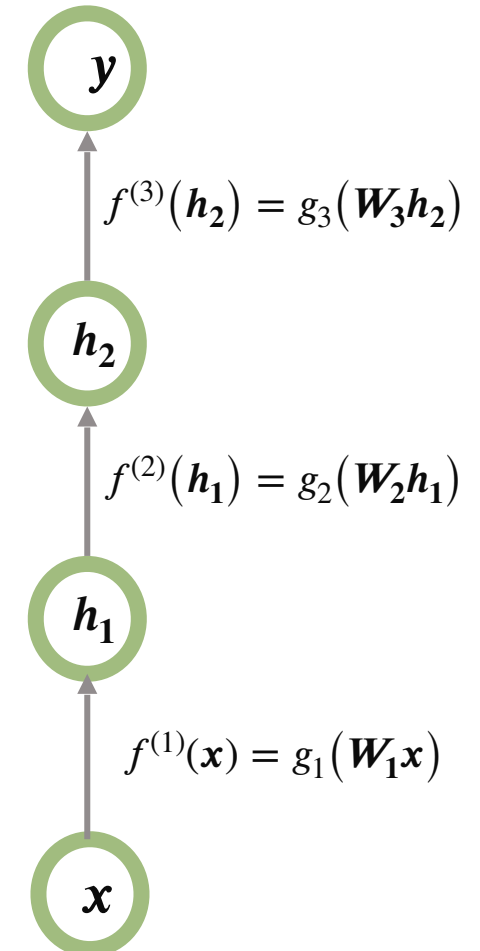
$$f(\mathbf{x}) = f^{(3)}\left(f^{(2)}\left(f^{(1)}(\mathbf{x})\right)\right)$$

$f^{(1)}$ is the first layer, $f^{(2)}$ the second, etc.

The length of the chain is the depth of the model

The name deep learning arises from this terminology

Final layer of a feedforward network, (here $f^{(3)}$) is the output layer



What are Hidden Layers

- * Behavior of hidden layers is not directly specified by the data
- * Learning algorithm must decide how to use those layers to produce a final value that is close to y
- * Training data does not say what individual layers should do
 - * This is one of the main tasks: Define the appropriate structure of the network and the hidden layers
- * Since the desired output for these layers is not shown, they are called hidden layers

Extending Linear Models

- * To represent non-linear functions of x
- * Apply linear model to transformed input $\phi(x)$ with non-linear ϕ
Equivalently kernel trick of SVM obtains nonlinearity
- * Function is nonlinear wrt x but linear wrt $\phi(x)$
We can think of ϕ as providing a set of features resulting in a new representation for x
- * Get ϕ
 - * Generic functions: There exists a set of useful kernel functions
 - * Manually engineer ϕ : Laborious and not transformable
 - * **Deep Learning**: Learn ϕ

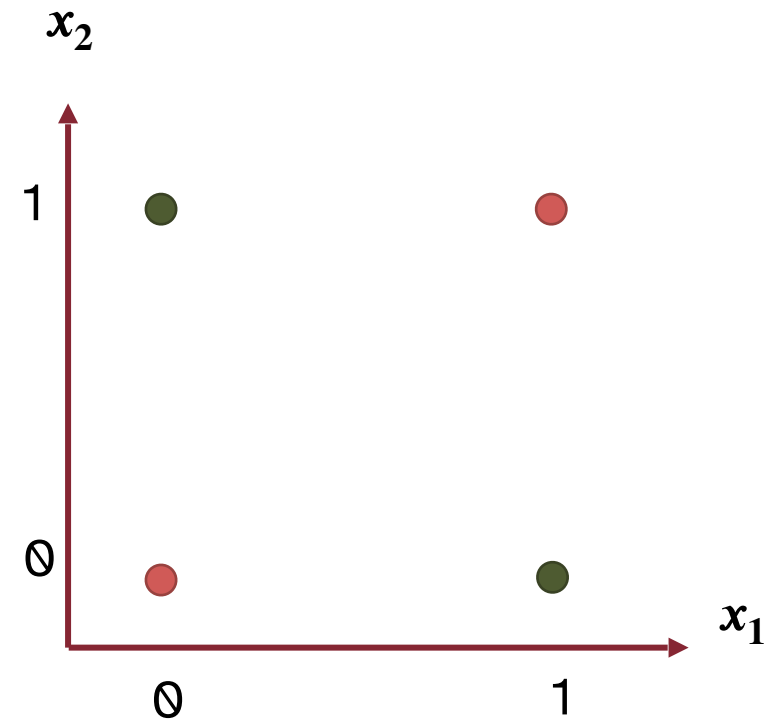
On the importance of Activation Functions

* The XOR Problem:

* XOR is a logical function

* It takes two inputs, x_1 and x_2

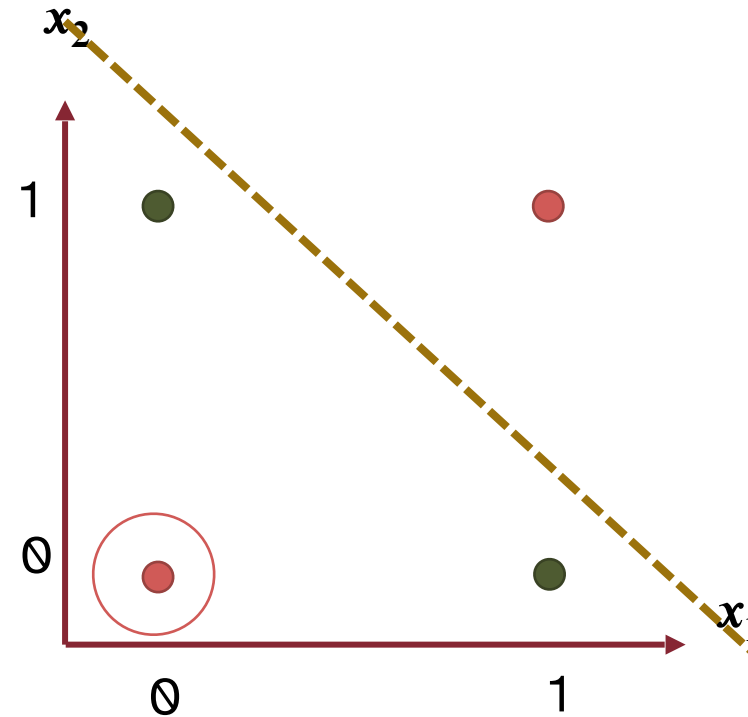
Input 1	Input 2	Output
0	0	0
0	1	1
1	0	1
1	1	0



On the importance of Activation Functions

- * The XOR Problem:
 - * XOR is a logical function
 - * It takes two inputs, x_1 and x_2
 - * **The problem is not linearly separable**

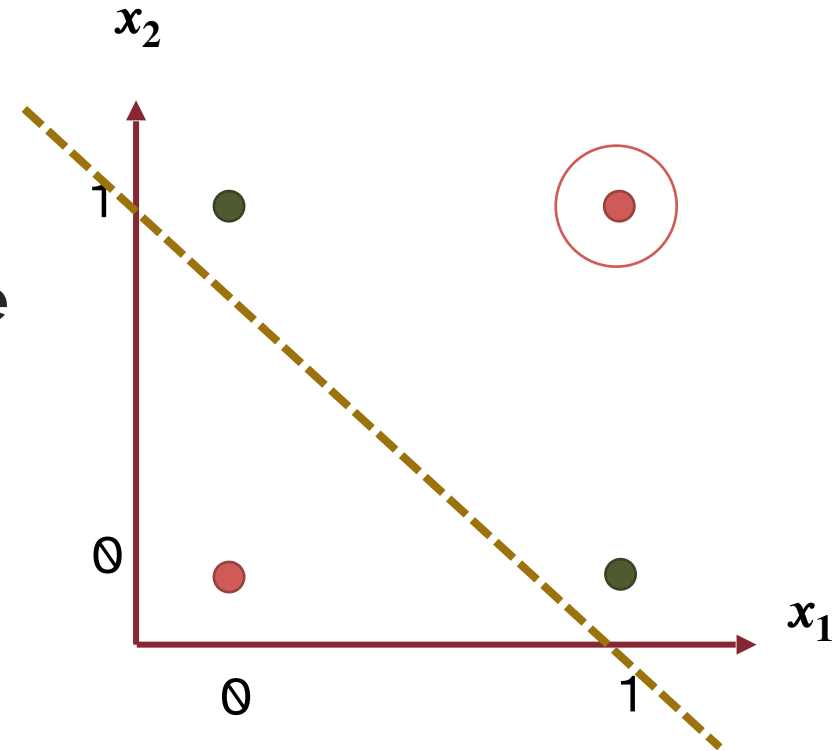
Input 1	Input 2	Output
0	0	0
0	1	1
1	0	1
1	1	0



On the importance of Activation Functions

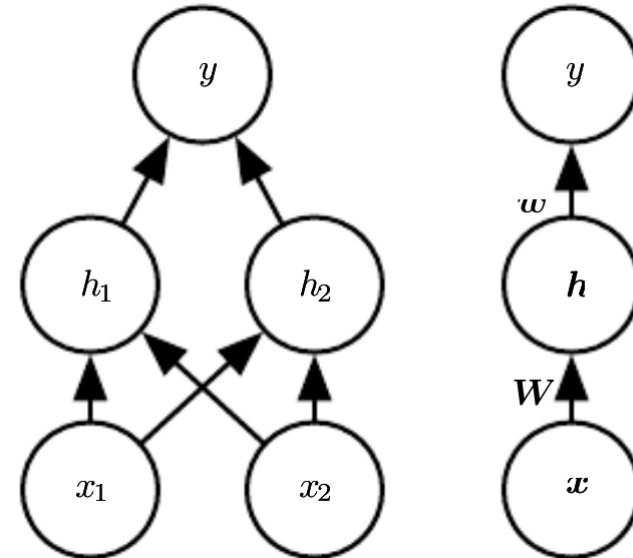
- * The XOR Problem:
 - * XOR is a logical function
 - * It takes two inputs, x_1 and x_2
 - * **The problem is not linearly separable**

Input 1	Input 2	Output
0	0	0
0	1	1
1	0	1
1	1	0



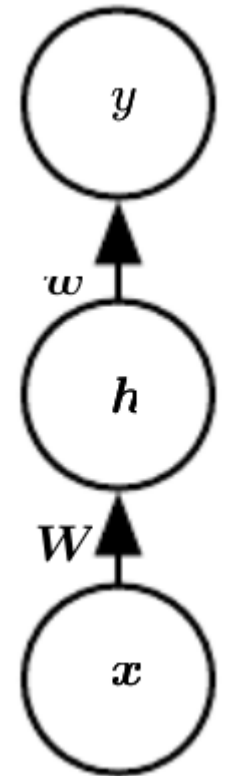
The XOR Problem

- * We are trying to approximate the XOR function
 - * Our training points: $X = \left\{ [0,0]^T, [1,0]^T, [0,1]^T, [1,1]^T \right\}$
 - * As we have seen, linear separation is not possible!
- * **Now lets use a FFN**
 - * One hidden layer, containing two units
 - * Both representations are equivalent
 - * Matrix \mathbf{W} describes mapping \mathbf{x} to \mathbf{h}
 - * Vector \mathbf{w} describes mapping \mathbf{h} to y
 - * (biases omitted)



What is computed

- * Layer 1 (hidden layer): vector of hidden units h computed by function $h = f^{(1)}(x; W, c)$ with bias c
- * Layer 2 (output layer) computes: $y = f^{(2)}(h; w, b)$
 - * w are the weights (it is basically a linear regression)
 - * Output is linear regression applied to h rather than to x
- * Complete model is $f(x; W, c, w, b) = f^{(2)}(f^{(1)}(x; W, c), w, b)$



Activation Function

✱ If we choose both $f^{(1)}$ and $f^{(2)}$ to be linear, the total function will still be linear; which is insufficient (as already seen)

✱ Therefore, we need an **Activation Function**:

✱ Remember, our function looks as follows:

$$y = \mathbf{x}^T \mathbf{w} + b$$

✱ Activation function g is typically chosen to be applied element-wise:

$$h_i = g(\mathbf{x}^T \mathbf{W}_{:i} + c_i)$$

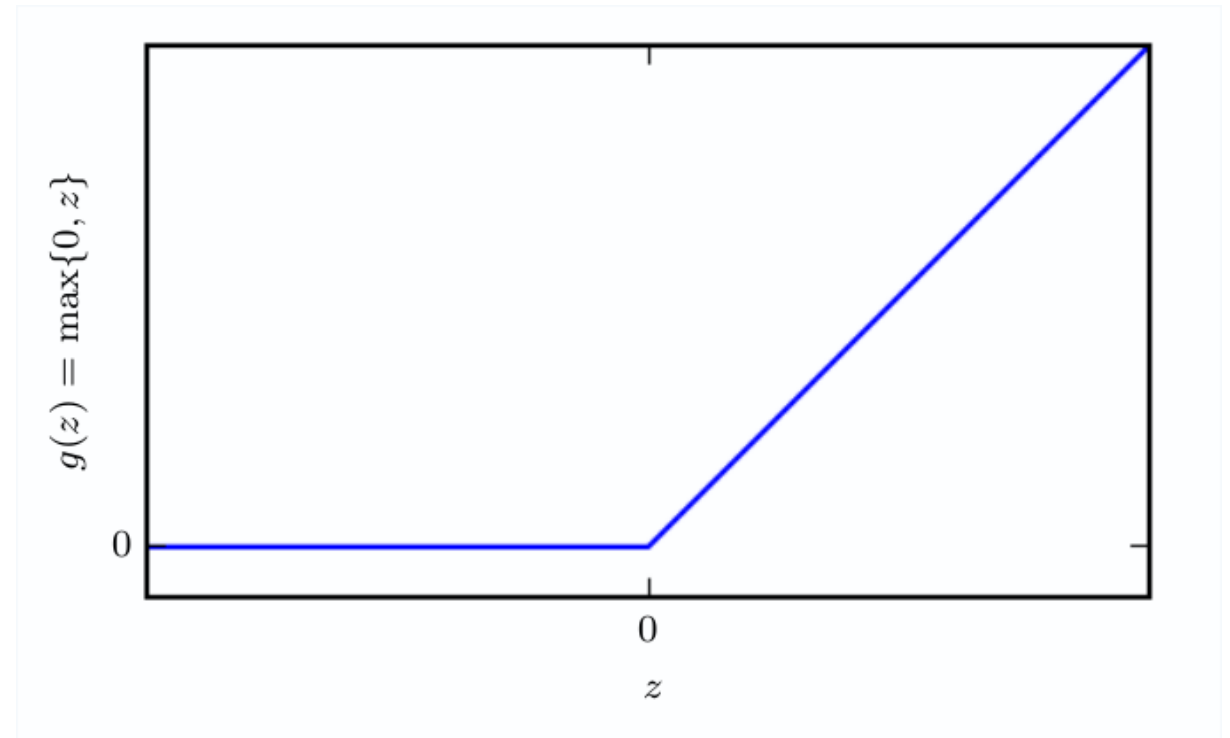
Default Activation Function: ReLU

➤ Activation function:

$$g(z) = \max\{0, z\}$$

* This already yields to a non-linear transformation

- * But still piecewise linear
- * Preserves couple of nice properties making gradient-based learning easy
- * Generalizes well



Back to our XOR Problem

➤ The complete equation is now:

$$f(\mathbf{x}; \mathbf{W}, c, \mathbf{w}, b) = \mathbf{w}^T \max\{0, \mathbf{W}^T \mathbf{x} + c\} + b$$

Lets use the following weights:

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, b = 0 \text{ and } \mathbf{X} = \begin{bmatrix} 0101 \\ 0011 \end{bmatrix}^T$$

Back to our XOR Problem

➤ The complete equation is now:

$$f(\mathbf{x}; \mathbf{W}, c, \mathbf{w}, b) = \mathbf{w}^T \max\{0, \mathbf{W}^T \mathbf{x} + c\} + b$$

Lets use the following weights:

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, b = 0 \text{ and } \mathbf{X} = \begin{bmatrix} 0101 \\ 0011 \end{bmatrix}^T$$

$$\mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

Back to our XOR Problem

➤ The complete equation is now:

$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^T \max\{0, \mathbf{W}^T \mathbf{x} + \mathbf{c}\} + b$$

Lets use the following weights:

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, b = 0 \text{ and } \mathbf{X} = \begin{bmatrix} 0101 \\ 0011 \end{bmatrix}^T$$

$$\mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

$$\mathbf{XW} + \mathbf{c} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

Back to our XOR Problem

➤ The complete equation is now:

$$f(\mathbf{x}; \mathbf{W}, \mathbf{c}, \mathbf{w}, b) = \mathbf{w}^T \max\{0, \mathbf{W}^T \mathbf{x} + \mathbf{c}\} + b$$

Lets use the following weights:

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, b = 0 \text{ and } \mathbf{X} = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}^T$$

$$\mathbf{XW} = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

$$\mathbf{XW} + \mathbf{c} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

$$\max\{0, \mathbf{XW} + \mathbf{c}\} =$$

$$\begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

Back to our XOR Problem

➤ The complete equation is now:

$$f(x; W, c, w, b) = w^T \max\{0, W^T x + c\} + b$$

Lets use the following weights:

$$W = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, b = 0 \text{ and } X = \begin{bmatrix} 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}^T$$

$$XW = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

$$XW + c = \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

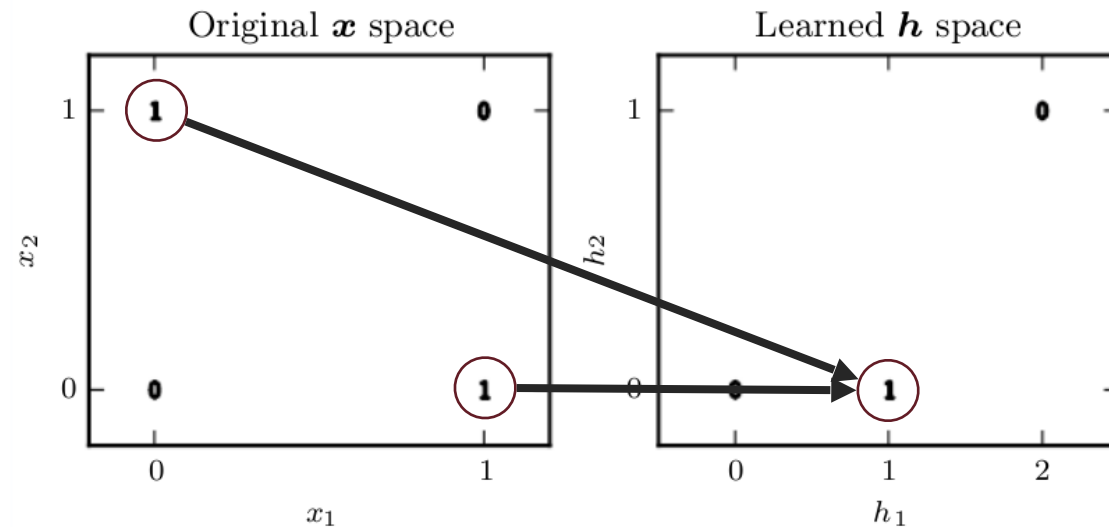
$$\max\{0, XW + c\} =$$

$$\begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$$

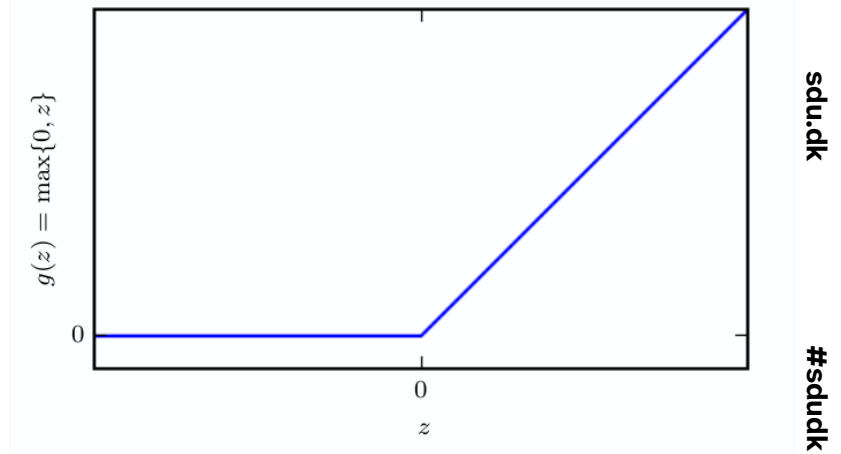
$$\begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix} w = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

“Learned XOR”

- Two points that must have output 1 have been collapsed into one
- Can be separated in a linear model:
 - For fixed h_2 output has to increase in h_1
- Of course, we “cheated”: We have “guessed” the parameters correctly which lead to the desired result
- In reality: Millions of parameters which we have to actually learn



Interim Summary



- * Linear models cannot solve the X-OR (exclusive-or) problem
- * Composition of linear models are still Linear, due to associativity.

In other words, we could rewrite $\mathbf{W}^{(2)}\mathbf{W}^{(1)}\mathbf{x}$ as we could set

$$\mathbf{W}' := \mathbf{W}^{(1)} \cdot \mathbf{W}^{(2)},$$

s.th. $\mathbf{W}'\mathbf{x} = \mathbf{W}^{(2)}\mathbf{W}^{(1)}\mathbf{x}$ are equivalent.

- * Therefore, we need to add **nonlinear activation functions** in-between
- * A simple but effective nonlinearity is the rectified linear unit (**ReLU**)

$$\text{ReLU}(\mathbf{x}) = \max(0, \mathbf{x}), \text{ i.e., } \mathbf{x} \text{ if } \mathbf{x} > 0, \text{ else } 0$$

Keywords

- * Single-Layer Networks
- * Deep Neural Networks
- * **Universal Approximation Theorem**
- * Activation functions
- * Output units

On the Power of Neural Networks

A linear model, mapping from features to outputs via matrix multiplication, can by definition represent only linear functions. But:

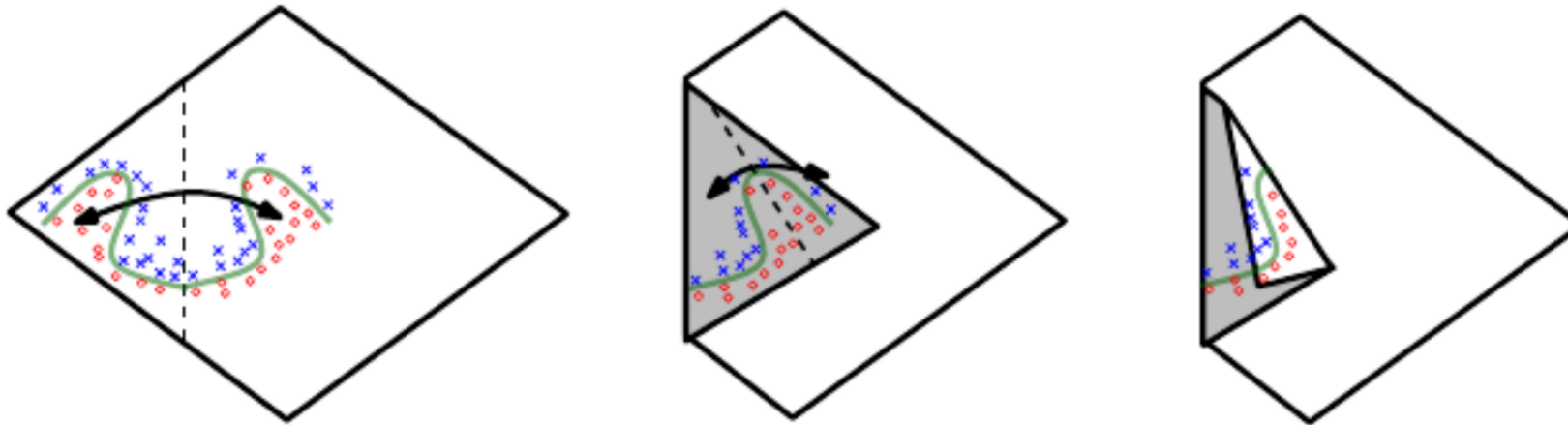
The **universal approximation theorem** states that a feedforward neural network with at least one hidden layer and a nonlinear activation function can approximate any compact, continuous function to an arbitrary degree of accuracy (depending on the hidden layer's width).

Implication of Theorem

- Whatever function we are trying to learn, a sufficiently large MLP will be able to represent it
- However we are not guaranteed that the training algorithm will learn this function
- No Free Lunch: There is no universal procedure for examining a training set of samples and choosing a function that will generalize to points not in training set

Advantage of deeper networks

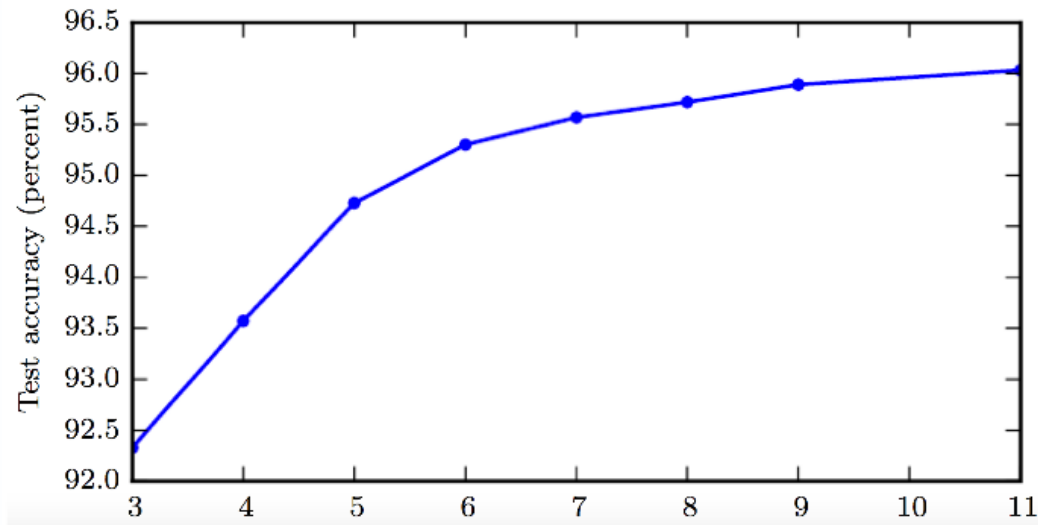
- Absolute value rectification creates mirror images of function computed on top of some hidden unit, wrt the input of that hidden unit.
- Each hidden unit specifies where to fold the input space to create mirror responses.
- By composing these folding operations, we obtain an exponentially large no. of piecewise linear regions which can capture all kinds of repeating patterns



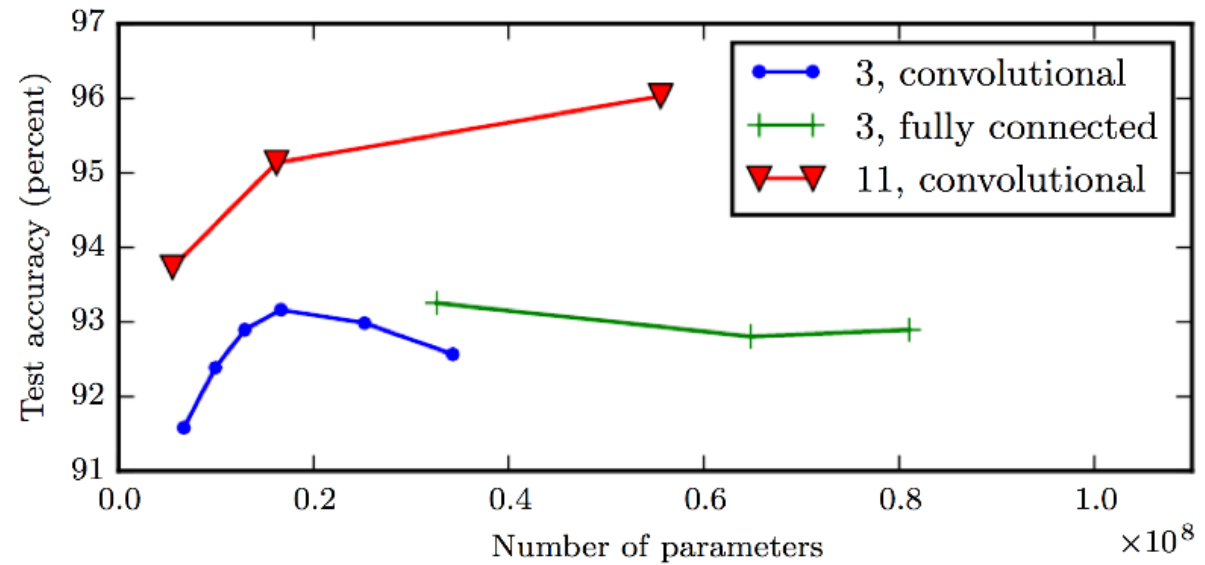
Intuition of Depth

- Any time we choose a ML algorithm we are implicitly stating a set of beliefs about what kind of functions that algorithm should learn
- Choosing a deep model encodes a belief that the function should be a composition several simpler functions
- The learning problem consists of discovering a set of underlying factors of variation that can in turn be described in terms of other, simpler underlying factors of variation.
- Empirically: Deeper Networks perform better 😊

Deeper Networks



Test accuracy consistently increases with depth



Increasing parameters without increasing depth is not as effective

Keywords

- * Single-Layer Networks
- * Deep Neural Networks
- * Universal Approximation Theorem
- * **Activation functions**
- * Output units

General Construction of a Hidden Unit

- * Accepts a vector of inputs x and computes an affine transformation $h' = Wx + b$
- * Computes an element-wise non-linear function $g(h')$
- * Most hidden units are distinguished from each other by the choice of activation function $g(h')$
 - * We look at: ReLU, Sigmoid and tanh, and other hidden units

Choice of Hidden Unit

- * We now look at how to choose the type of hidden unit in the hidden layers of the model
- * Design of hidden units is an active research area that does not have many definitive guiding theoretical principles
- * ReLU is an excellent default choice
- * But there are many other types of hidden units available
- * When to use which kind?
 - * Hard to predict in advance which will work best for a problem at hand
 - * Design process is trial and error

Logistic Sigmoid

- * Prior to introduction of ReLU, most neural networks used logistic sigmoid activation

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

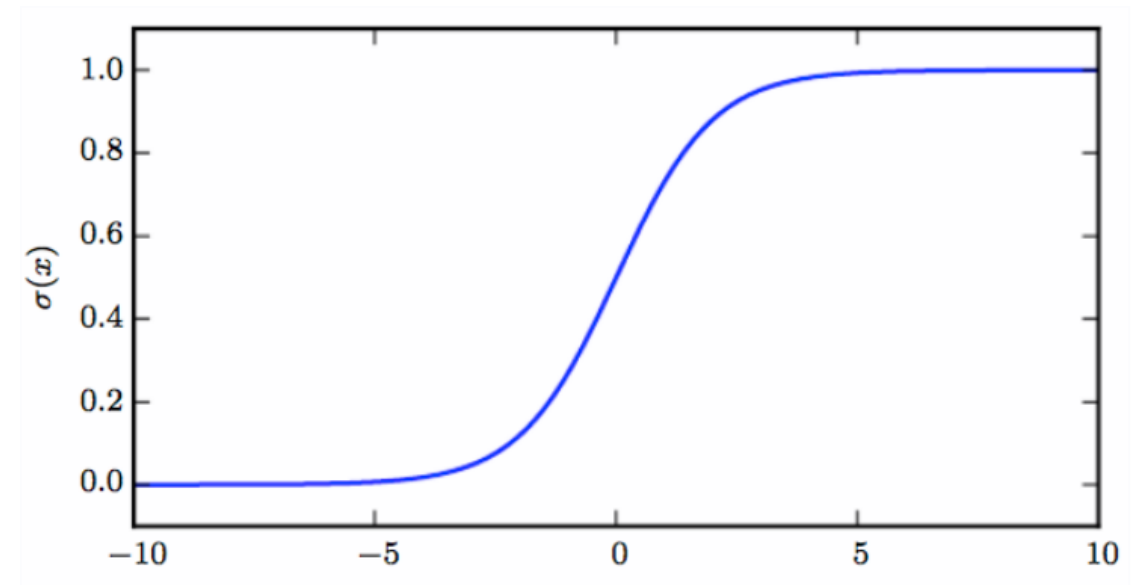
- * Or the hyperbolic tangent

$$g(z) = \tanh(z)$$

- * These activation functions are closely related because

$$\tanh(z) = 2\sigma(2z) - 1$$

- * Sigmoid units are used to predict probability that a binary variable is 1



Other Activation Functions

* Cosine

$$h = \cos(\mathbf{W}\mathbf{x} + \mathbf{b})$$

* Feedforward networks on MNIST obtained error rate of less than 1%

* Radial Basis

$$h_i = \exp\left(-\frac{1}{\sigma^2} \|\mathbf{W}_{:,i} - \mathbf{x}\|^2\right)$$

* Becomes more active as \mathbf{x} approaches a template $\mathbf{W}_{:,i}$

* Softplus

$$g(a) = \zeta(a) = \log(1 + e^a)$$

* Smooth version of the rectifier

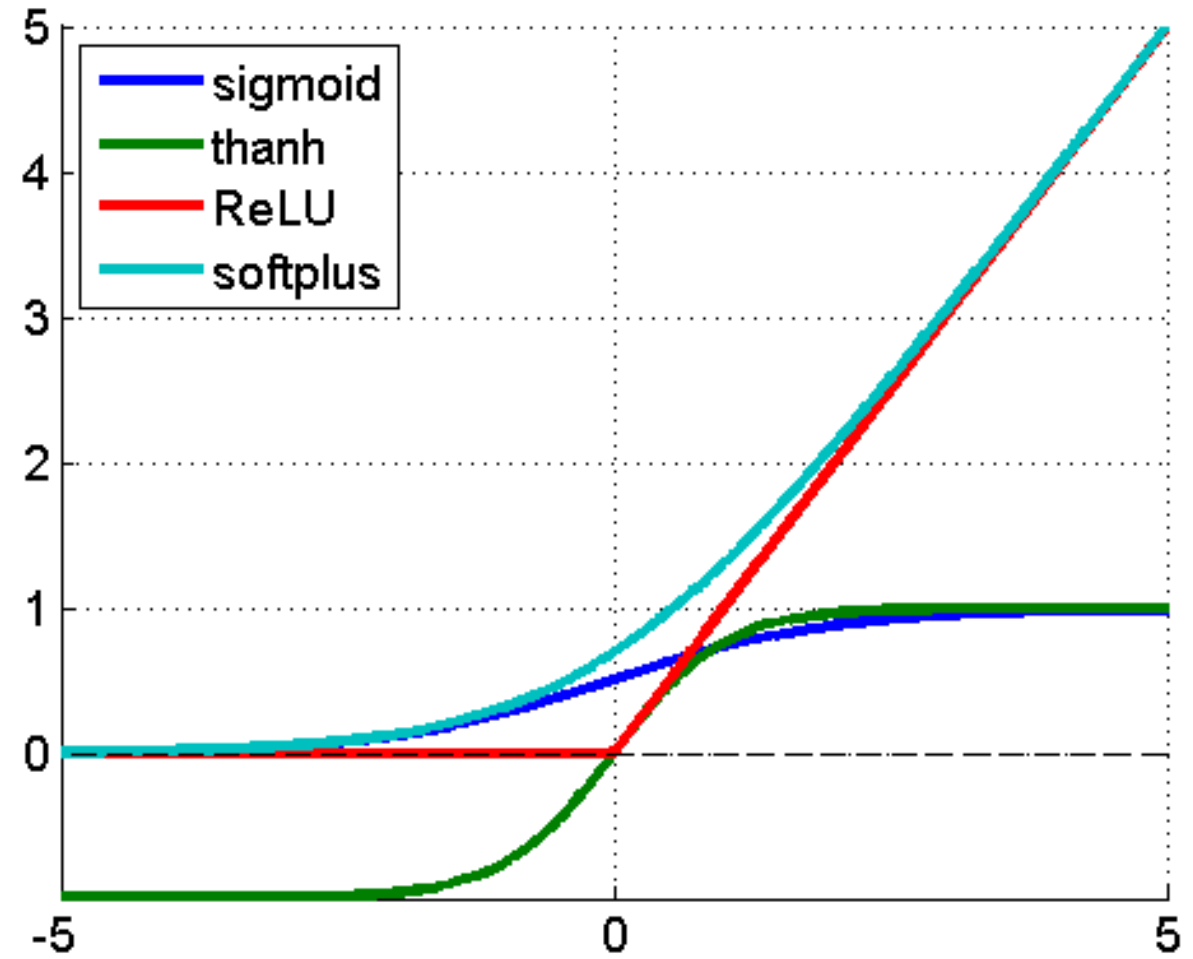
* Hard tanh

$$g(a) = \max(-1, \min(1, a))$$

* Shaped similar to tanh and the rectifier but it is bounded

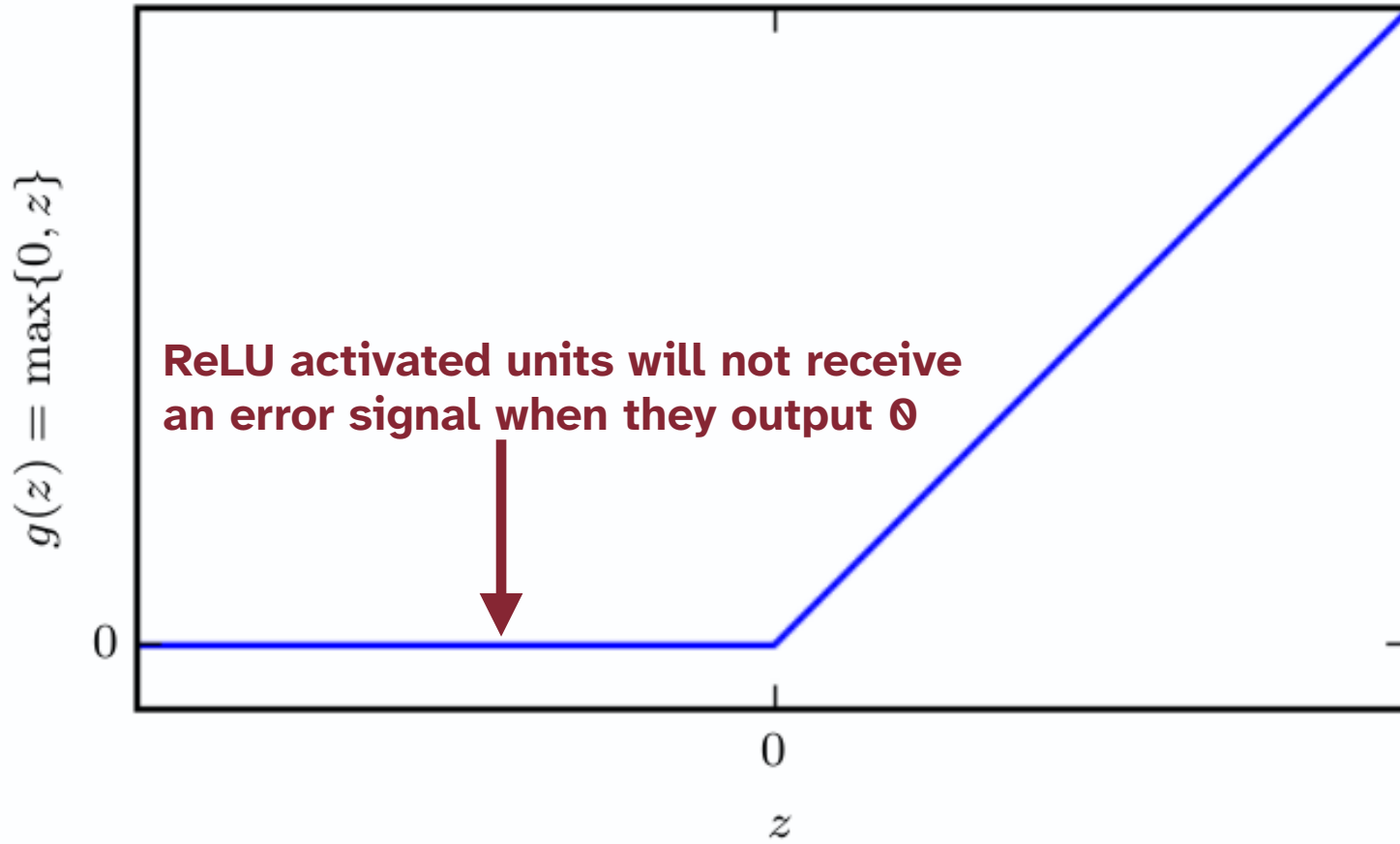
Sigmoid Saturation

- * Sigmoidals saturate
- * ReLU and softplus increase for input > 0
- * As networks get very deep, this linear increase of the ReLU in the positive domain is important, such that the error signals don't get too small in the backward pass. Sigmoid units would learn **much** slower.



But is ReLU the end of the story?

$$h = g(W^T x + b)$$
$$g(x) = \max\{0, x\}$$



Generalizations of ReLU

- * Three methods based on using a non-zero slope α_i when $z_i < 0$:

$$h_i = g(\mathbf{z}, \boldsymbol{\alpha})_i = \max(0, z_i) + \alpha_i \min(0, z_i)$$

or equivalently: $g(\mathbf{z}, \boldsymbol{\alpha})_i = z_i$ **if** $z_i > 0$ **else** $\alpha_i \cdot z_i$

- * Absolute-value rectification:

- * fixes $\alpha_i = -1$ to obtain $g(\mathbf{z}) = |\mathbf{z}|$

- * Leaky ReLU:

- * fixes α_i to a small value like 0.01

- * Parametric ReLU or PReLU

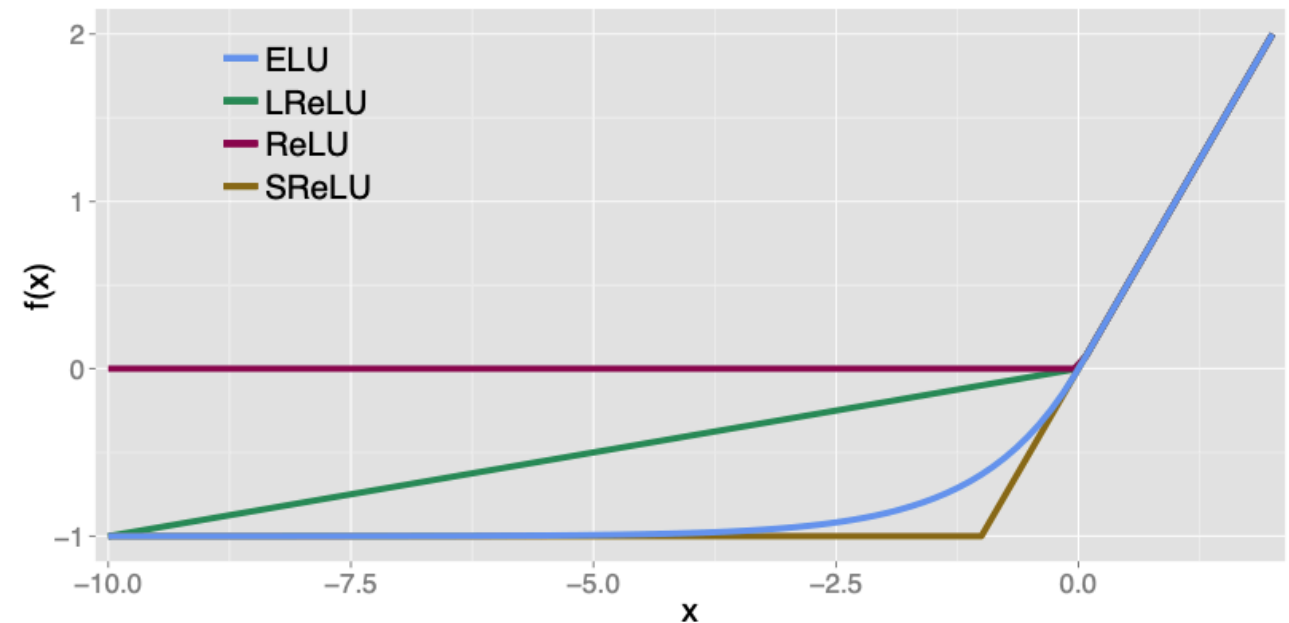
- * treats α_i as a learnable parameter

Dead ReLU's and how to do better: ELU

- * **ReLU is a good default choice**
- * **But:** ReLU's don't receive any gradient (error signal) when they output 0
- * It can happen that they output 0 for all possible inputs – then you have **dead ReLU's**
- * Variants that mitigate this problem:
 - * **LeakyReLU:**
 $h_i = z_i$ if $z_i > 0$ else αz_i
(with small positive α)
 - * **Exponential Linear Unit (ELU),**
better than LeakyReLU

The *exponential linear unit* (ELU) with $0 < \alpha$ is

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$



From arXiv:1511.07289v5

Dead ReLU's and how to do even better: GELU

- Popular in Transformers is the Gaussian Error Linear Unit (GELU)
- Exceeds performance of ReLU and ELU

Idea: GELU scales activations by how much larger they are than other activations, assuming $N(0,1)$ distribution

↓

$$\text{GELU}(x) = xP(X \leq x) = x\Phi(x) = x \cdot \frac{1}{2} \left[1 + \text{erf}(x/\sqrt{2}) \right]$$

We can approximate the GELU with

$$0.5x(1 + \tanh[\sqrt{2/\pi}(x + 0.044715x^3)])$$

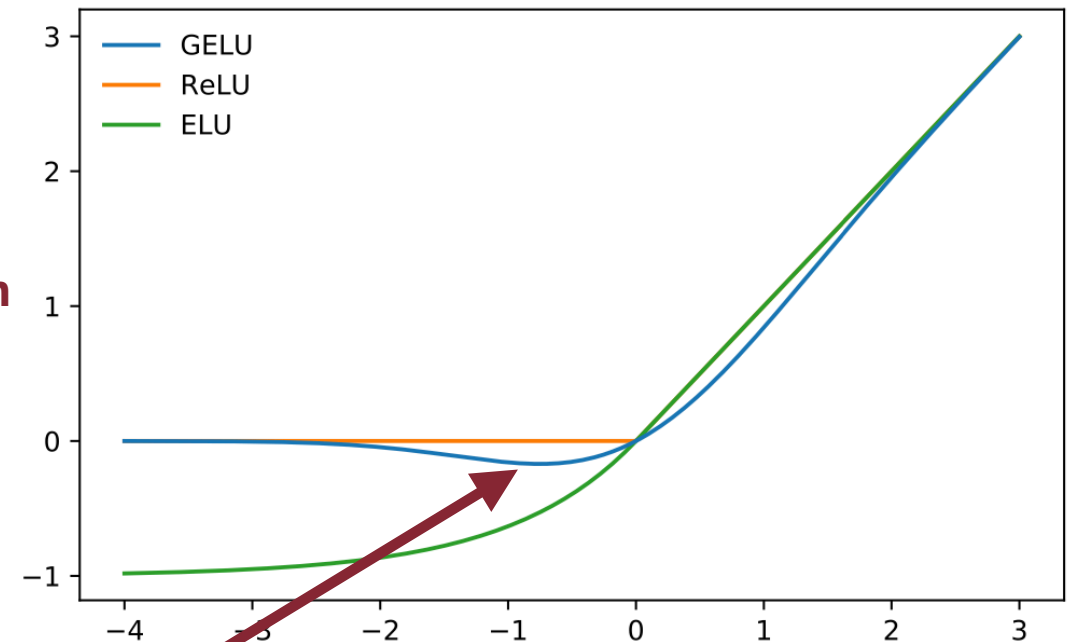
or

↑

$$x\sigma(1.702x),$$

**Pattern also seen in other recent activation functions:
x multiplied by logistic sigmoid of x**

From arXiv:1606.08415v5



Leads to a smooth “bump” just before zero

Summary of activation functions

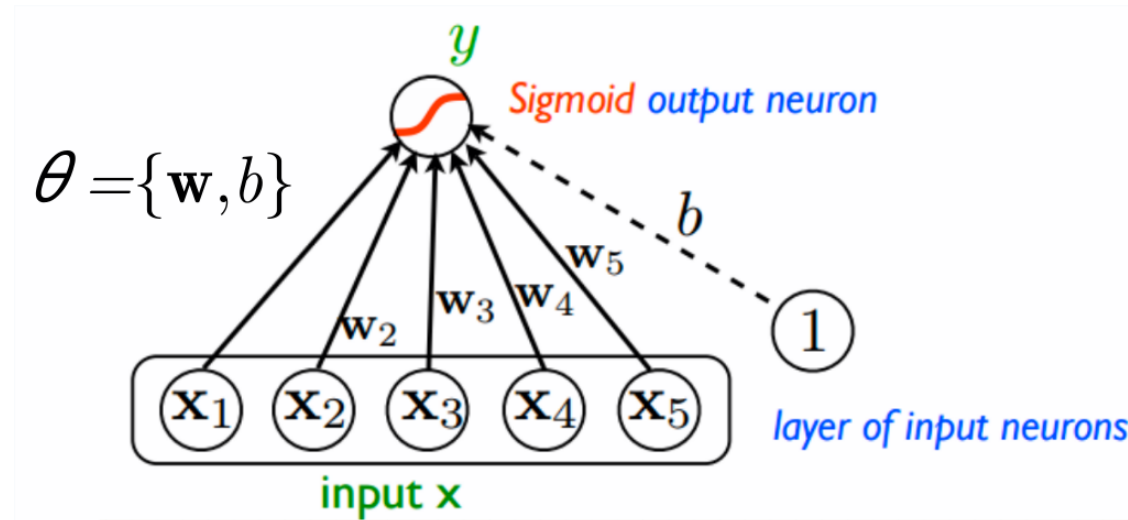
- * Sigmoidal functions were used traditionally
- * **But** they saturate and make learning difficult – **Sigmoid Saturation**
- * **Therefore** ReLU's, which don't saturate on the positive side
- * **But** ReLU's don't receive an update when they output 0 – **Dead ReLU's**
- * **Therefore** do allow small, downscaled negative values (e.g., **LeakyReLU**)
- * **But** we don't want them to become too negative, else they will have difficulty to re-activate
- * **Therefore** have them saturate on the negative range (e.g., **ELU**)
- * **And** often used recently, introducing a small bump close to zero (e.g., **GELU**)

Keywords

- * Single-Layer Networks
- * Deep Neural Networks
- * Universal Approximation Theorem
- * Activation functions
- * **Output units**

Output Units

- * The choice of cost function is **tightly** coupled with the choice of output unit (later...)
- * Most of the time we use cross-entropy between data distribution and model distribution (later ...)
- * Choice of how to represent the output then determines the form of the cross-entropy function



Role of an Output Unit

* Any output unit is in principle also usable as a hidden unit

* A feedforward network provides a hidden set of features

$$h = f(x; \theta)$$

* Role of output layer is to provide some additional transformation from the features to the task that network must perform

* Common Output Units

* **Linear units**: no non-linearity (used for Regressions)

* **Sigmoid units** each individual output is between 0 and 1 (used for many-out-of-K classification)

* **Softmax units** each individual output is between 0 and 1 **and** all outputs together sum up to 1
(e.g., used for one-out-of-K classification)

Linear Units for Regressions

- * Linear unit: simple output based on affine transformation with no nonlinearity

- * Given features h , a layer of linear output units produces a vector

$$\hat{y} = W^T h + b$$

- * Linear units are often used to produce mean \hat{y} of a conditional Gaussian distribution

- * $P(y|x) = N(y; \hat{y}, \sigma^2)$

We will see later why this is interesting

Most commonly used for regression:
Mean squared error loss (MSELoss)

Sigmoid Units for Binary Classification

- * Task of predicting value of binary variable y
 - * Classification problem with two classes
- * Maximum likelihood approach is to define a Bernoulli distribution over y conditioned on x
- * Neural net needs to predict $p(y = 1 \mid x) \in [0,1]$:
$$P(y = 1 \mid x) = \max \left\{ 0, \min \{ 1, \mathbf{w}^T \mathbf{h} + \mathbf{b} \} \right\}$$
 - * We would define a valid conditional distribution, but cannot train it effectively with gradient descent
 - * A gradient of 0: learning algorithm cannot be guided

Binary Cross-Entropy

Sigmoid Units

➤ Sigmoid always gives a gradient

$$\hat{y} = \sigma(\mathbf{w}^T \mathbf{h} + b)$$

with

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

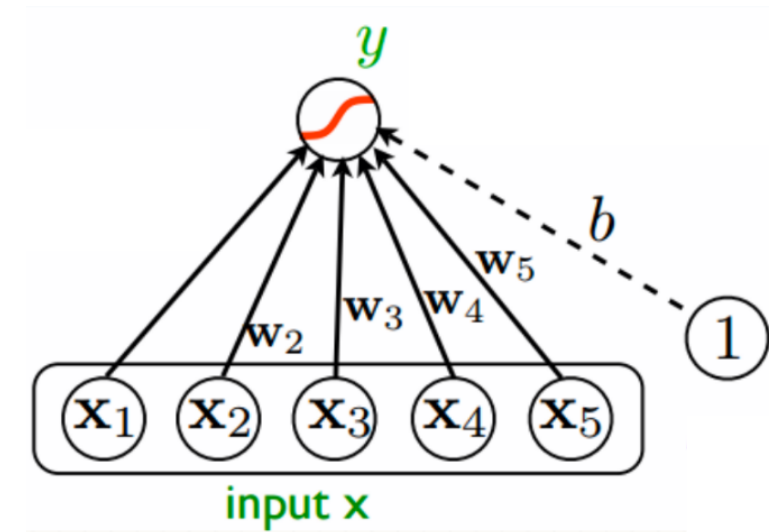
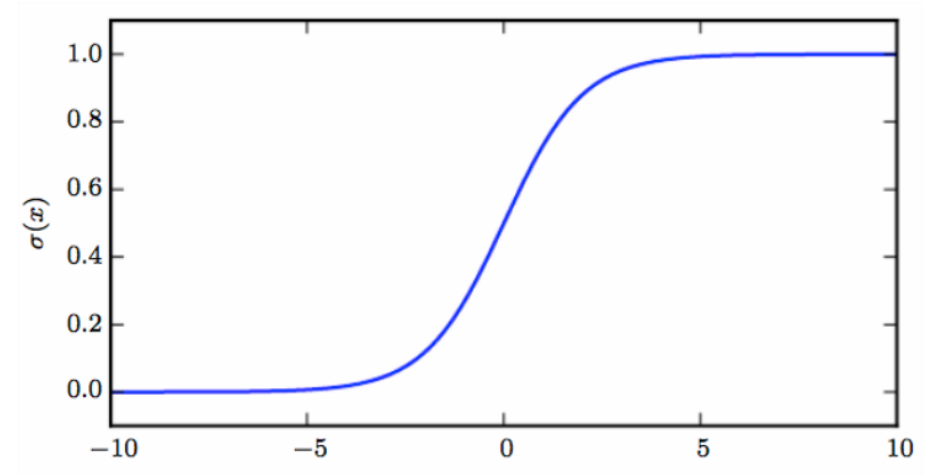
➤ Sigmoid Unit has two components:

➤ A linear layer to compute

$$z = \mathbf{w}^T \mathbf{h} + b$$

➤ A sigmoid activation function to convert z into a probability

➤ You will understand these slides later ☺.



Softmax Units for General Classification

- Any time we want a probability distribution over a discrete variable with n values we may use the softmax function
- Softmax most often used for output of classifiers to represent distribution over n classes
 - Also, inside the model itself when we wish to choose between one of n options
- For the Binary case we have produced a single number
$$\hat{y} = P(y = 1 | x)$$
- Now, we have to produce a vector $\hat{\mathbf{y}}$
$$\hat{y}_i = P(y = i | x)$$

Cross-entropy Loss

➤ You will understand these slides later ☺.

Softmax Definition

- The **softmax function** (aka normalized exponential), is defined as follows:

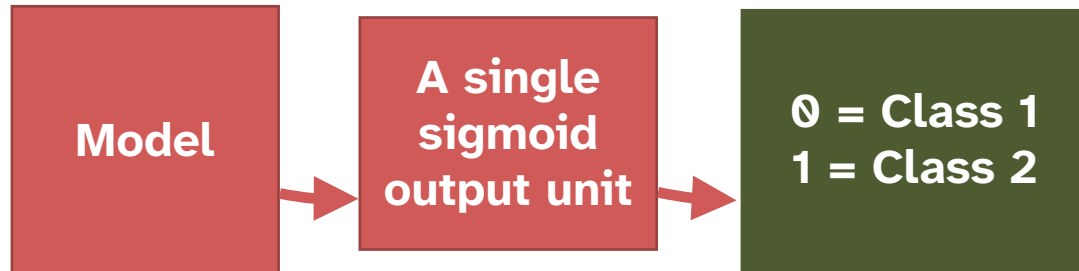
$$\text{softmax}(a_1, \dots, a_K) = \left(\frac{\exp a_1}{\sum_j^K \exp a_j}, \frac{\exp a_2}{\sum_j^K \exp a_j}, \dots, \frac{\exp a_K}{\sum_j^K \exp a_j} \right)$$

- We will learn later more about this function
- But it should be visible that it creates a valid distribution
all values are between 0 and 1 and all values sum up to 1

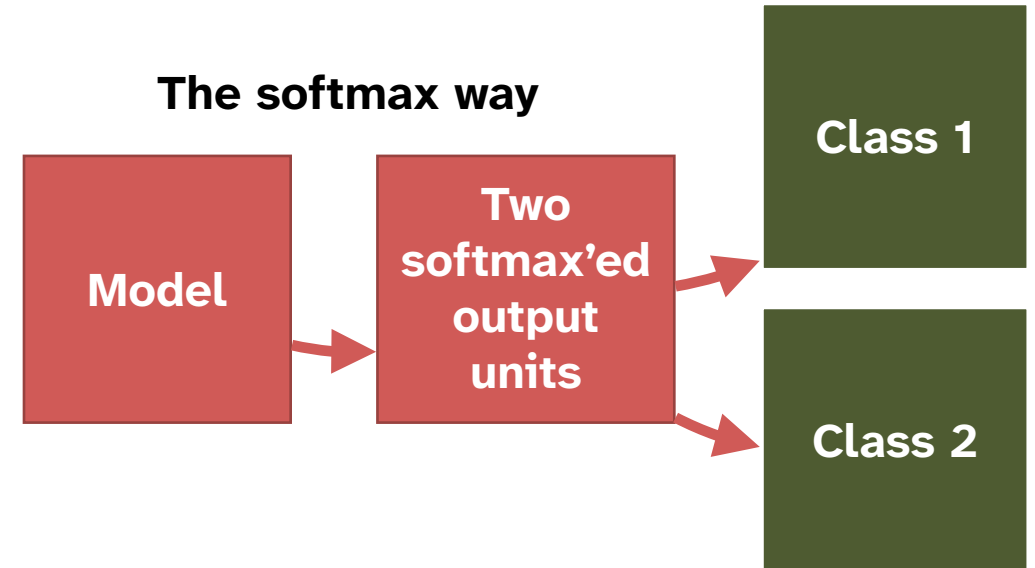
➤ You will understand these slides later ☺.

Softmax, even for binary classification

The sigmoid way



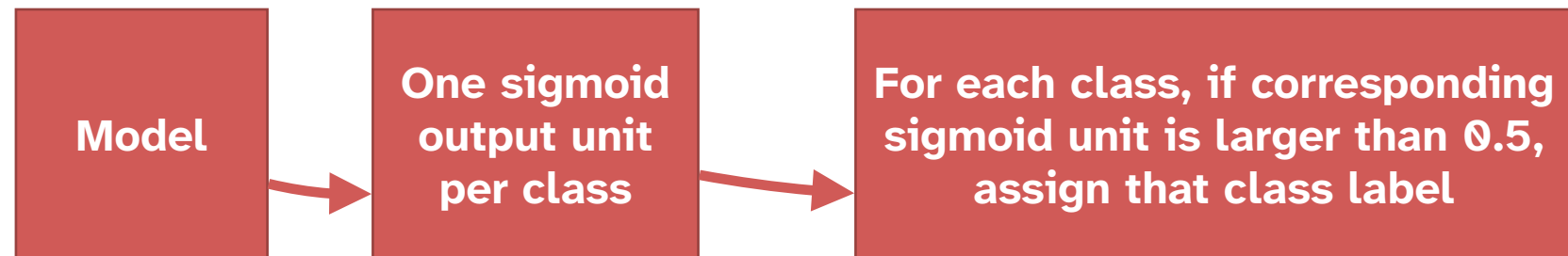
The softmax way



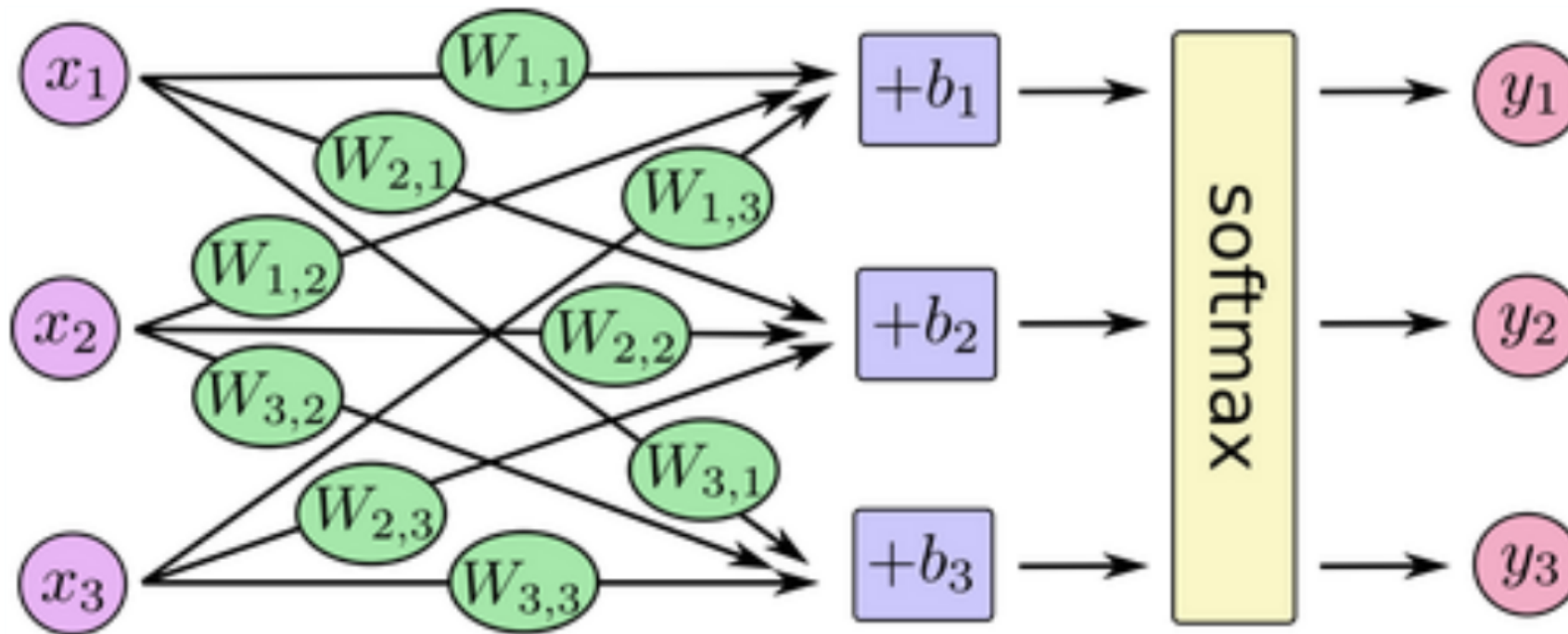
Softmax way often works better in practice!

But when do I need sigmoid output units then?

- **Prime example: Multi-label classification**
- Consider the case that you have **M** possible classes.
- And multiple class labels can be assigned to a single example.
- For example, you have a news article, and you need to classify it into topics {Sports, Football, Basketball, Europe, Denmark, Germany, Politics, Economics, Science}
 - **multiple can be correct**
- “A **scientific** investigation of the **football** match between **Denmark** and **Germany**”



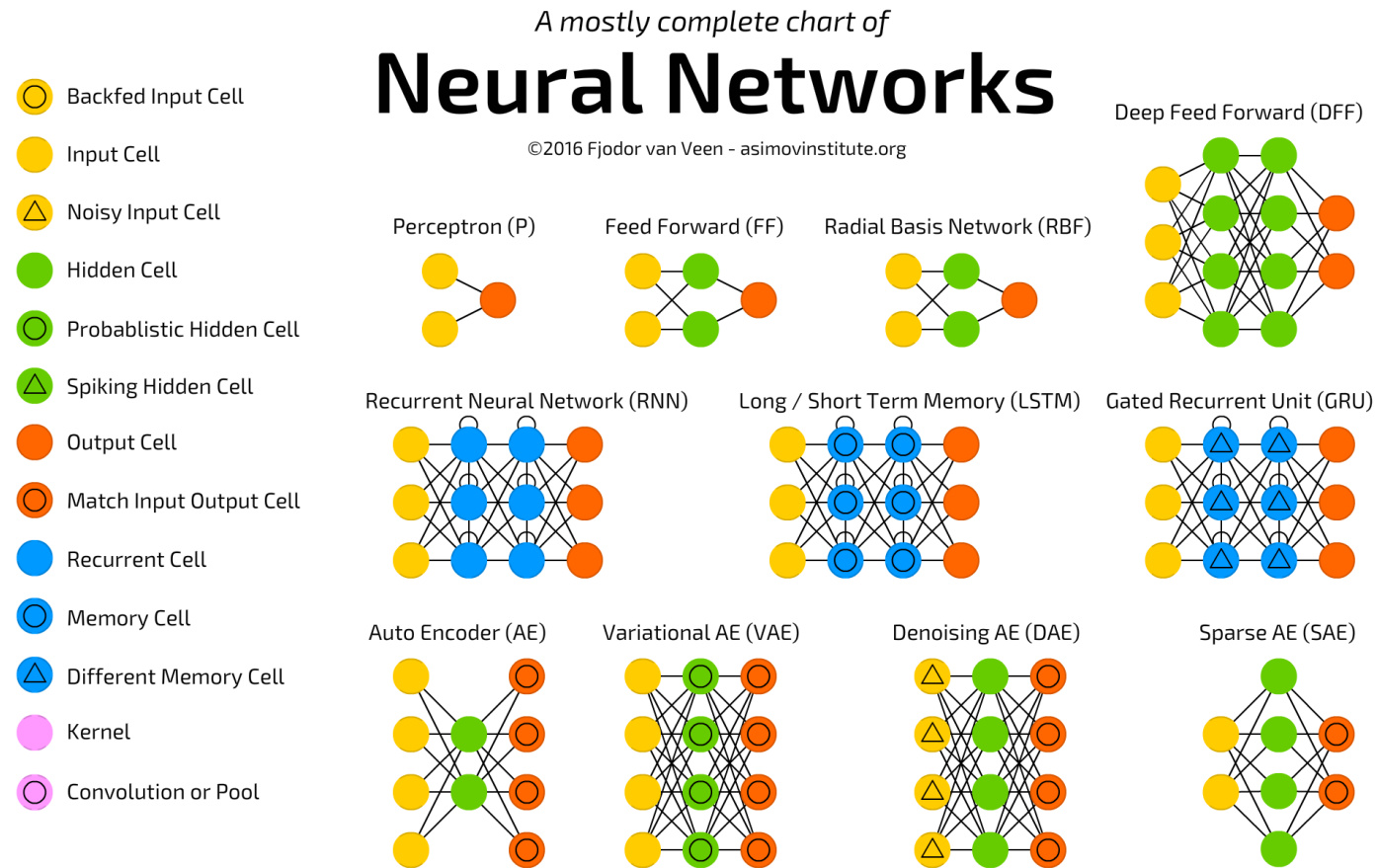
Softmax Regression



Keywords

- * Single-Layer Networks
- * Deep Neural Networks
- * Universal Approximation Theorem
- * Activation functions
- * Output units

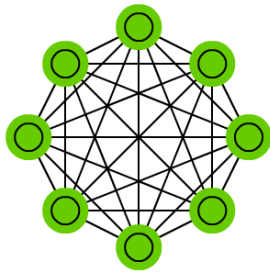
A mostly complete chart of Neural Networks



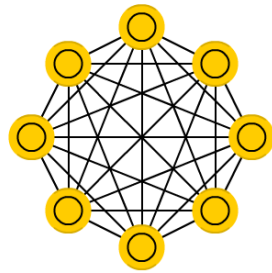
➤ <https://towardsdatascience.com/the-mostly-complete-chart-of-neural-networks-explained-3fb6f2367464>

A mostly complete chart of Neural Networks

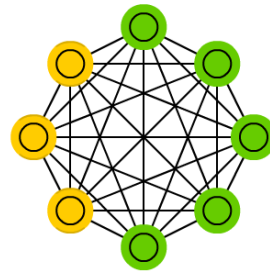
Markov Chain (MC)



Hopfield Network (HN)



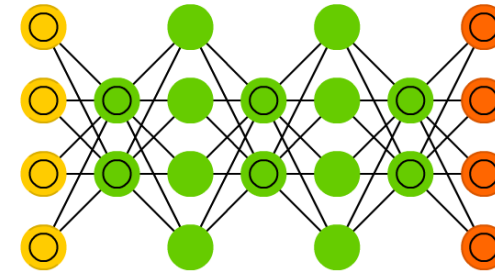
Boltzmann Machine (BM)



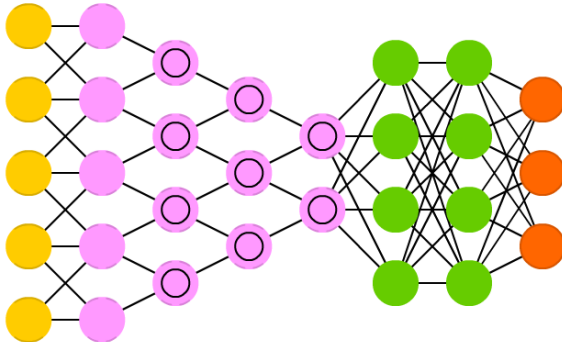
Restricted BM (RBM)



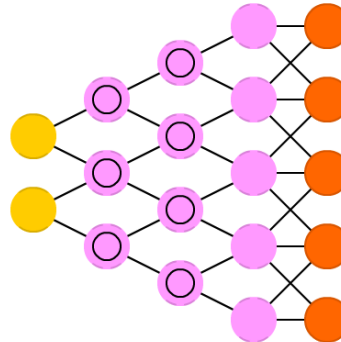
Deep Belief Network (DBN)



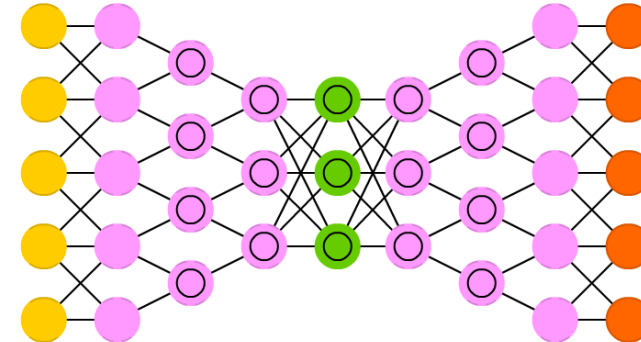
Deep Convolutional Network (DCN)



Deconvolutional Network (DN)

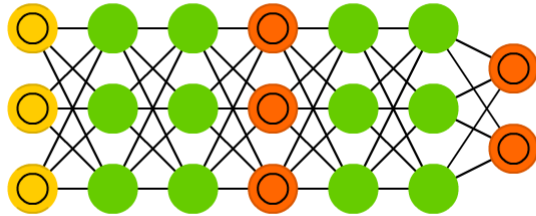


Deep Convolutional Inverse Graphics Network (DCIGN)

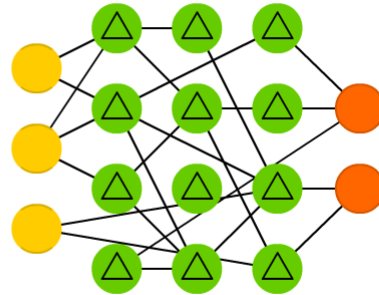


A mostly complete chart of Neural Networks

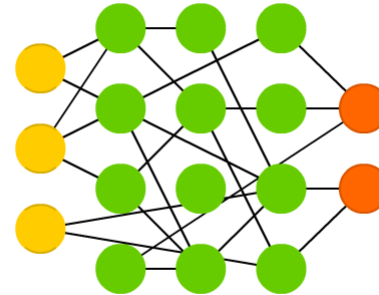
Generative Adversarial Network (GAN)



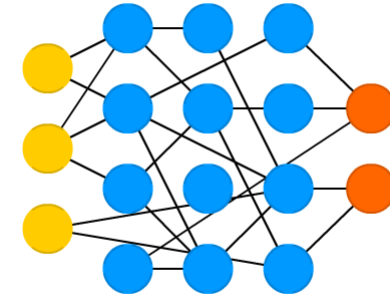
Liquid State Machine (LSM)



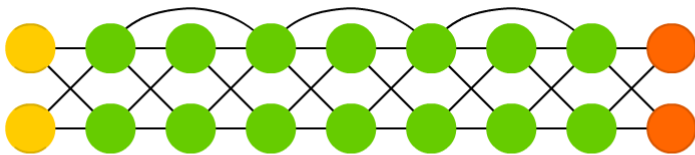
Extreme Learning Machine (ELM)



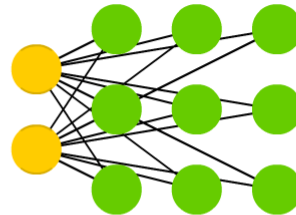
Echo State Network (ESN)



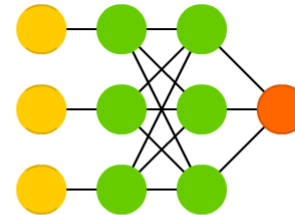
Deep Residual Network (DRN)



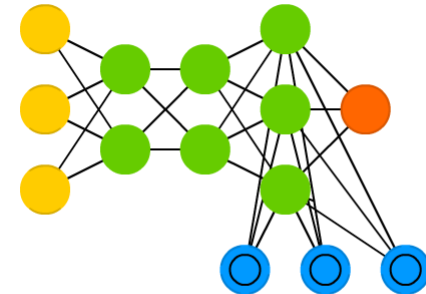
Kohonen Network (KN)



Support Vector Machine (SVM)



Neural Turing Machine (NTM)



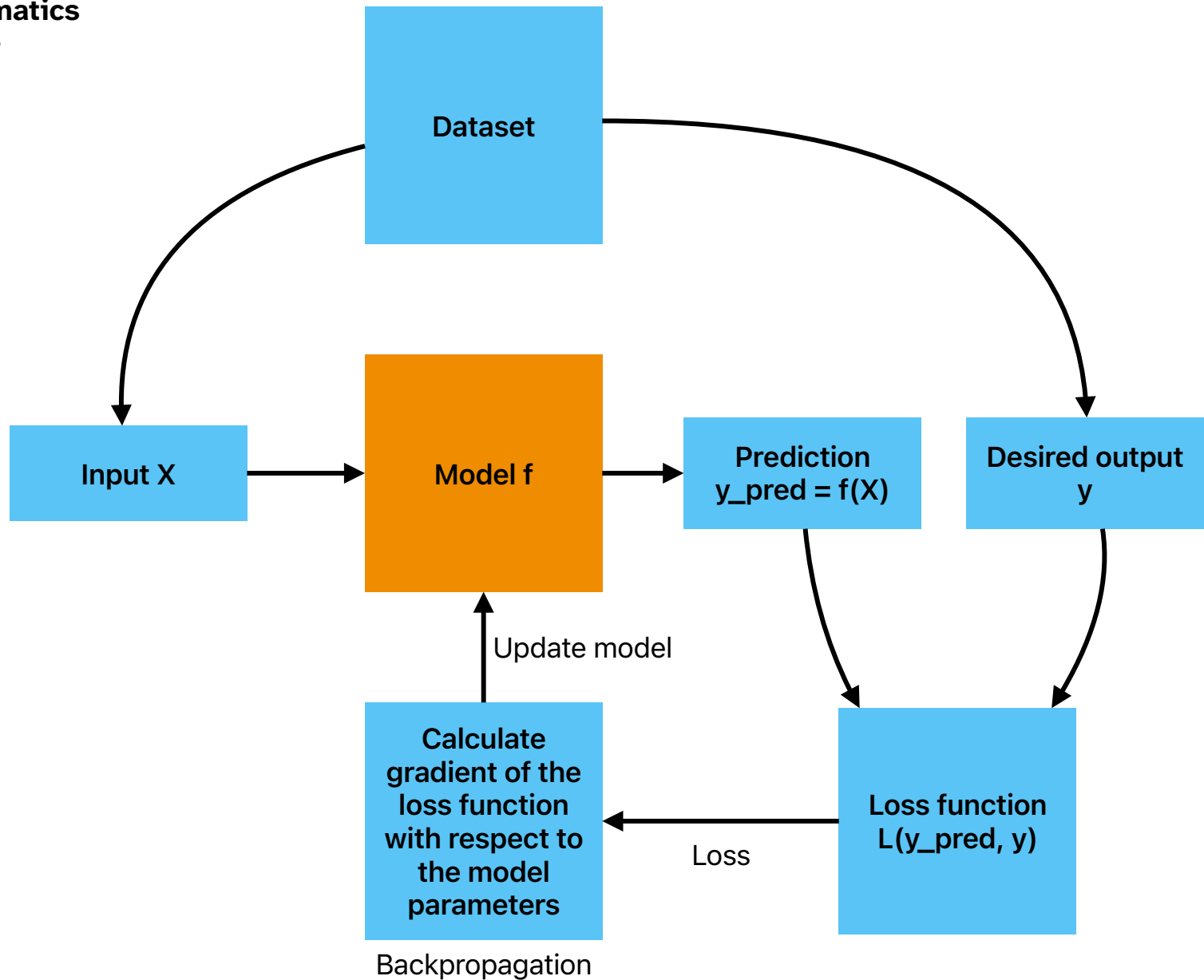
Main Architectural Considerations

1. **Choice of depth of network**
2. **Choice of width of each layer**

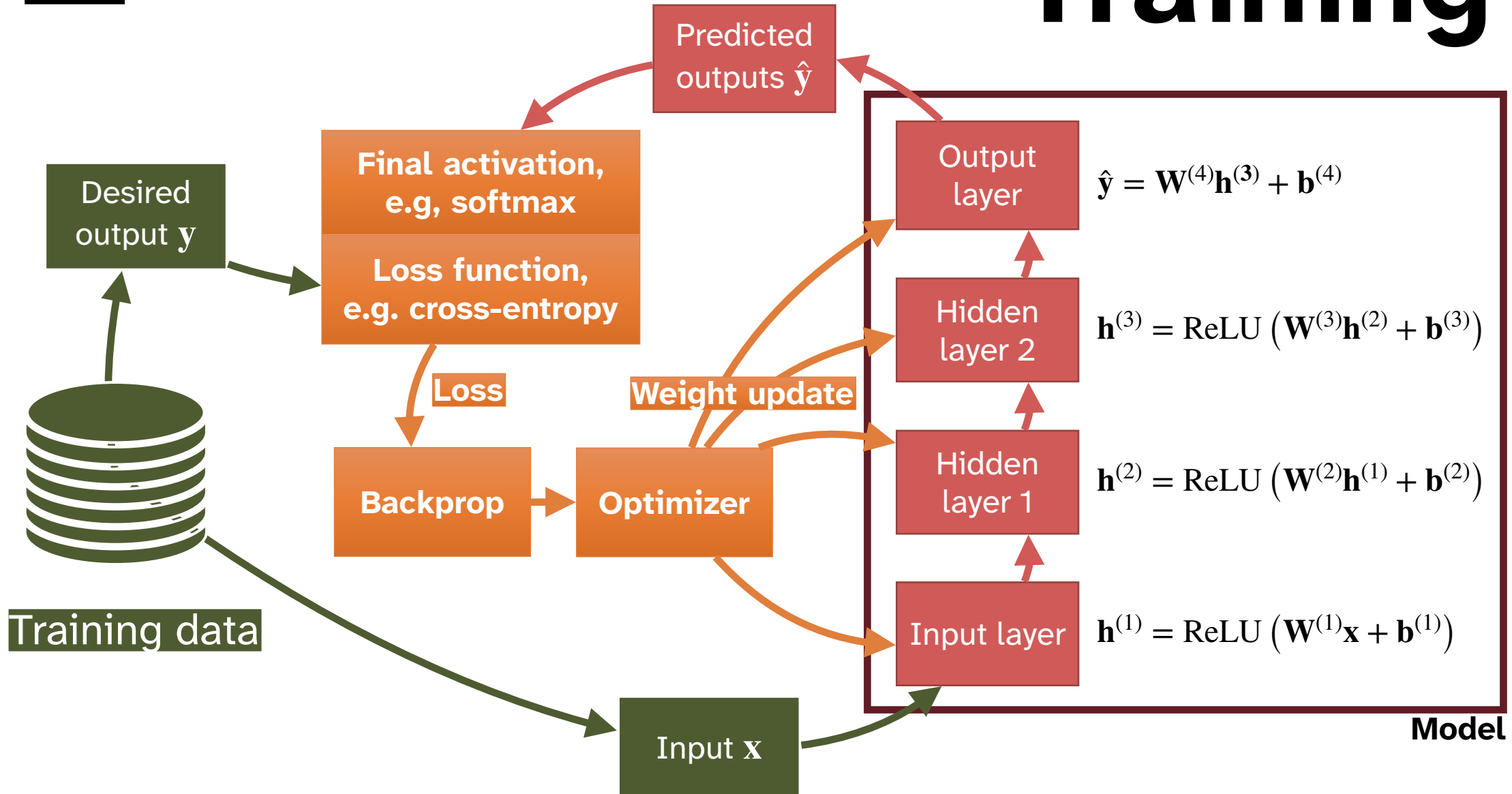
- * **Deeper networks have**

- * Far fewer units in each layer
- * Far fewer parameters
- * Often generalize well to the test set
- * But are often more difficult to optimize

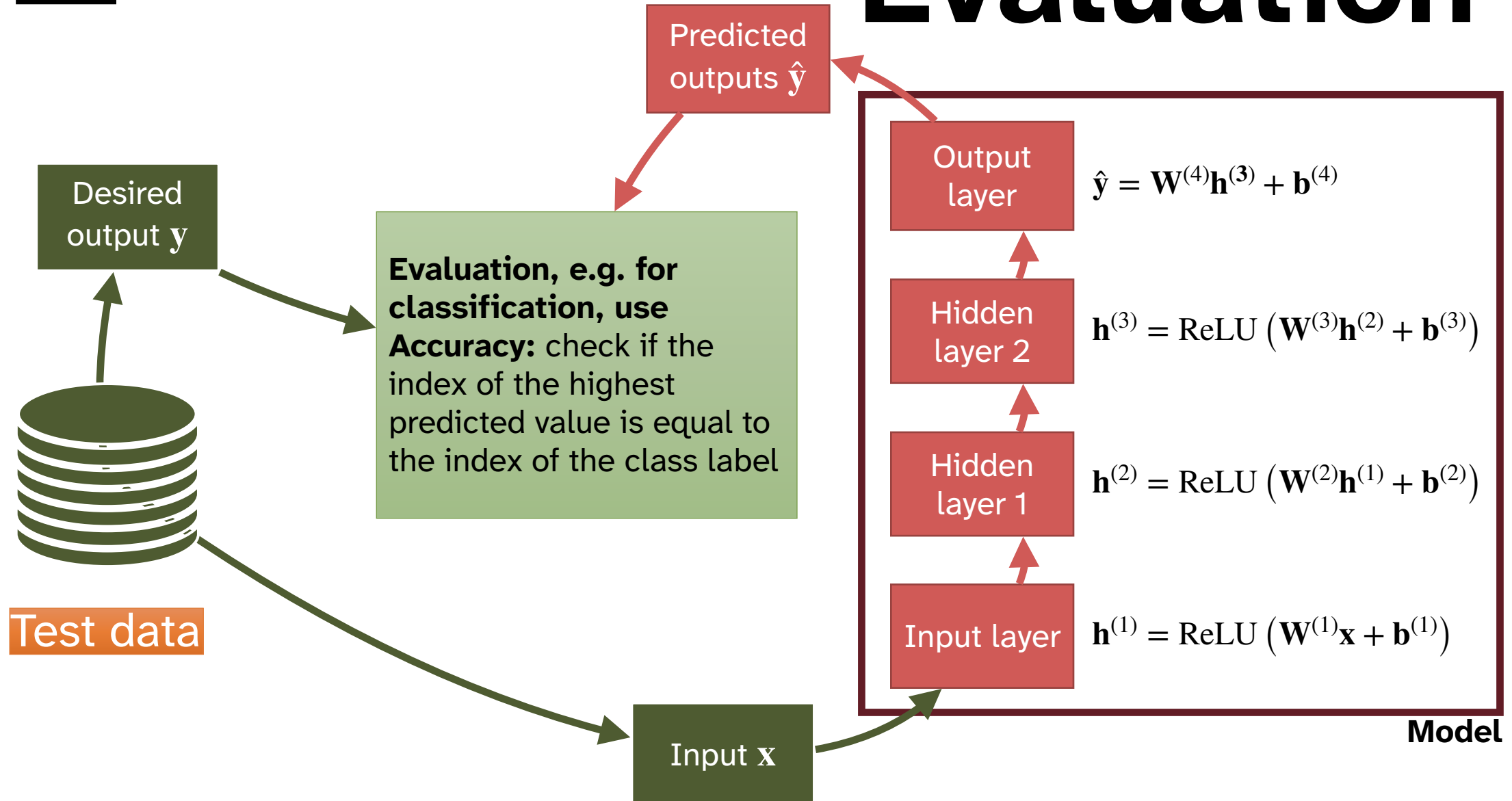
- * **Ideal network architecture must be found via experimentation guided by validation set error**



Training



Evaluation



Training loop (pseudo-PyTorch)

```
# [...] model & training_set defined
loss_function = nn.CrossEntropyLoss() # define the loss
function, Softmax activation included
optimizer = optim.SGD(model.parameters(), lr=0.001) #
Connects your optimizer with parameters
# Inside training loop
For each x with label y in training_set:
    model.zero_grad() # Resets gradient to zero - important!
    y_hat = model(x) # Feed data into the model
    loss = loss_function(y_hat, y) # Calculate the loss
    loss.backward() # Back-propagate the error
    optimizer.step() # Update model weights
```

Coding example

* example_xor.py

Take-home

- * Linear models – even multi-layered – are not enough to solve a simple XOR problem.
- * **Universal Approximation**; already with a few hidden layers and nonlinearities, a neural network model can **approximate any compact function to an arbitrary degree of accuracy**
- * **Loss function** and **output units** are tightly coupled and typically determined by the task
- * **Hidden units**: Rectified Linear Units are a good choice, yet there are some variants that tackle its shortcomings (dead ReLU's)
- * **Depth** and **Width** of a neural network are important for its performance
- * **The basic training loop – will be picked up again in the exercises**

Keywords

- * Single-Layer Networks
- * Deep Neural Networks
- * Universal Approximation Theorem
- * Activation functions
- * Output units

How's my teaching?



<https://www.admonymous.co/lukasgalke>