

Lecture 12: Neural networks

STATS 202: Data Mining and Analysis

Linh Tran

stat202@gmail.com



Department of Statistics
Stanford University

August 7, 2023



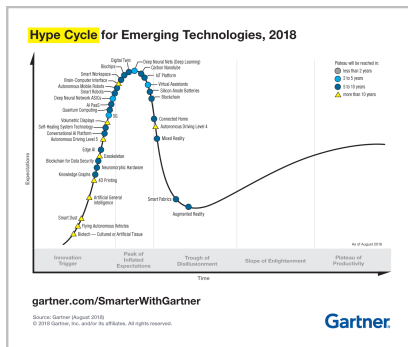
- ▶ Homework 4 due this Friday.
- ▶ Final project predictions due this Sunday night.
- ▶ Final project writeup is due next Wednesday.
 - ▶ Optional (will take max of Final project & Final exam).
- ▶ Final exam is a week from this Saturday
 - ▶ Time: August 19 @ 7:00 P.M. - 10:00 P.M.
 - ▶ Location: Skilling Auditorium
 - ▶ Practice exam to be released Friday
 - ▶ Accommodation requests should be made now
- ▶ No formal lectures next week
- ▶ Review this Friday



- ▶ Introduction
- ▶ Logistic regression
- ▶ Back propagation
- ▶ Function approximation
- ▶ Feature extraction
- ▶ Model generalization
- ▶ Advanced topics

Currently, the most popular algorithm amongst ML practitioners.

- ▶ Many times, used within the context of Artificial Intelligence.
- ▶ Simply a general function estimation algorithm.
- ▶ Though is often hyped up the media.



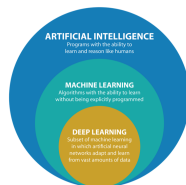
Gartner's hype cycle for 2018.



Lots of buzz words, but what do they mean?

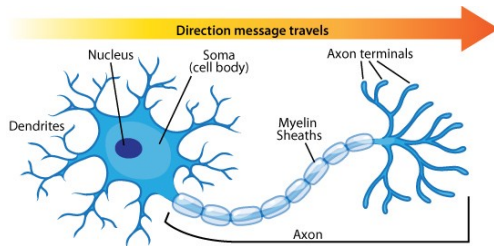
Definitions:

- ▶ *AI*: human-like machines or programs.
- ▶ *ML*: Algorithms that learn from data.
- ▶ *DL*: A type of ML algorithm, using neural networks (typically with many layers).



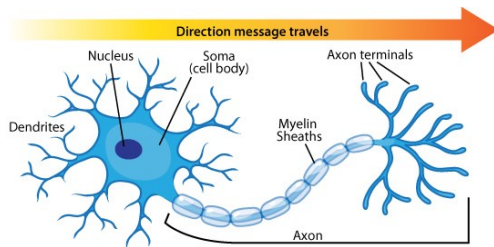


But: what exactly are neural networks?





But: what exactly are neural networks?

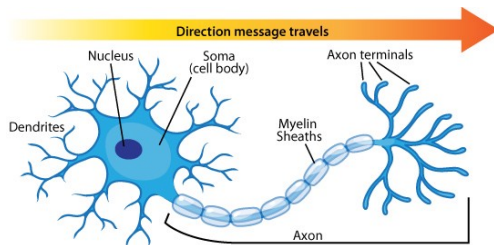


Some potential answers

- ▶ A *universal function approximator*.



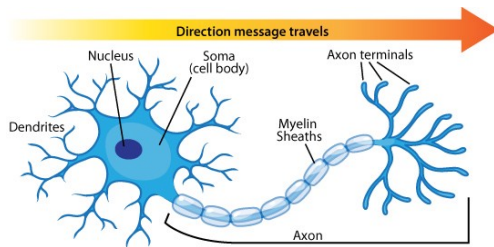
But: what exactly are neural networks?



Some potential answers

- ▶ A *universal function approximator*.
- ▶ A *feature extractor*.

But: what exactly are neural networks?



Some potential answers

- ▶ A *universal function approximator*.
- ▶ A *feature extractor*.
- ▶ A *model generalizer*.



Recall: logistic regression is a linear model with a logit link function, i.e.

$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p) : \sigma(z) = \frac{1}{1 + \exp(-z)} \quad (1)$$

Let's rephrase this by:

1. Using b to denote β_0 (aka the bias)
2. Using \mathbf{W} to denote $(\beta_0, \dots, \beta_p)$ (aka the weights)
3. Using matrix notation

$$\mathbb{P}(Y = 1|\mathbf{X}) = \underbrace{\sigma}_{\text{non-linearity}}(\mathbf{X}\mathbf{W} + b) \quad (2)$$



When the function is non-linear, our prior option was to do feature transformations, e.g.

- ▶ Expand predictor set (e.g. non-linear transformations, interactions, etc.).
- ▶ Define a kernel (e.g. find a function $f(\cdot, \cdot)$ that is positive definite).



When the function is non-linear, our prior option was to do feature transformations, e.g.

- ▶ Expand predictor set (e.g. non-linear transformations, interactions, etc.).
- ▶ Define a kernel (e.g. find a function $f(\cdot, \cdot)$ that is positive definite).

Another option: build the non-linearity into the model specification, e.g.

$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\sigma(\mathbf{X}\mathbf{W}_1 + b_1)\mathbf{W}_2 + b_2) \quad (3)$$



When the function is non-linear, our prior option was to do feature transformations, e.g.

- ▶ Expand predictor set (e.g. non-linear transformations, interactions, etc.).
- ▶ Define a kernel (e.g. find a function $f(\cdot, \cdot)$ that is positive definite).

Another option: build the non-linearity into the model specification, e.g.

$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\sigma(\mathbf{X}\mathbf{W}_1 + b_1)\mathbf{W}_2 + b_2) \quad (3)$$

This is a neural network (with 1 hidden layer)!



For logistic regression:

$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\mathbf{X} \underbrace{\mathbf{W}}_{p \times 1} + \underbrace{b}_{1 \times 1}) \quad (4)$$



For logistic regression:

$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\underbrace{\mathbf{X} \mathbf{W}}_{p \times 1} + \underbrace{b}_{1 \times 1}) \quad (4)$$

For neural networks:

$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\underbrace{\sigma(\mathbf{X} \mathbf{W}_1 + b_1)}_{p \times M} \underbrace{\mathbf{W}_2}_{M \times 1} + \underbrace{b_2}_{1 \times 1}) \quad (5)$$

M specifies how many hidden nodes we have

- ▶ Called '*hidden*' since it's not directly observed by us.
- ▶ Also referred to as '*embeddings*'.



$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\underbrace{\sigma(\mathbf{X} \mathbf{W}_1 + b_1)}_{p \times M} \underbrace{\mathbf{W}_2 + b_2}_{M \times 1}) \quad (6)$$

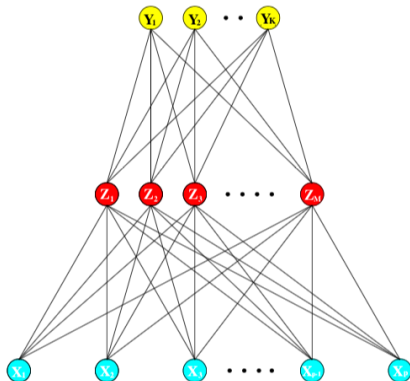


FIGURE 11.2. Schematic of a single hidden layer, feed-forward neural network.



We can iteratively apply our non-linear operations, e.g.

$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\cdots \sigma(\sigma(\mathbf{X}\mathbf{W}_1 + b_1)\mathbf{W}_2 + b_2) \cdots \mathbf{W}_B + b_B) \quad (7)$$

Where B is the number of iterations (i.e. *hidden layers*).



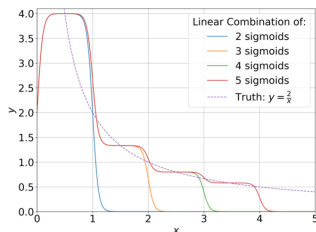
We can iteratively apply our non-linear operations, e.g.

$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\cdots \sigma(\sigma(\mathbf{X}\mathbf{W}_1 + b_1)\mathbf{W}_2 + b_2) \cdots \mathbf{W}_B + b_B) \quad (7)$$

Where B is the number of iterations (i.e. *hidden layers*).

Each activation (e.g. sigmoid) can approximate a local change

► B sigmoids \implies approximate at $\approx B$ points



n.b. Each layer needs the number of hidden nodes specified.



Our examples have been for binary outcomes so far.

Question: What about multinomial outcomes

- ▶ e.g. Which of digits 0 through 9 is this photo?



Our examples have been for binary outcomes so far.

Question: What about multinomial outcomes

- ▶ e.g. Which of digits 0 through 9 is this photo?

Recall: for logistic regression, we're modeling

$$\log \left[\frac{\mathbb{P}(Y = 1|\mathbf{X})}{1 - \mathbb{P}(Y = 1|\mathbf{X})} \right] = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p \quad (8)$$

$$= \mathbf{XW} \quad (9)$$

where \mathbf{X} is our $n \times p$ design matrix and \mathbf{W} is our $p \times 1$ parameter vector.



For multinomial regression, let $Y \in \{1, \dots, K\}$. We can model

$$\log \left[\frac{\mathbb{P}(Y = 1|\mathbf{X})}{\mathbb{P}(Y = K|\mathbf{X})} \right] = \mathbf{X}\mathbf{W}_1 \quad (10)$$

$$\log \left[\frac{\mathbb{P}(Y = 2|\mathbf{X})}{\mathbb{P}(Y = K|\mathbf{X})} \right] = \mathbf{X}\mathbf{W}_2 \quad (11)$$

$$\dots = \dots \quad (12)$$

$$\log \left[\frac{\mathbb{P}(Y = K - 1|\mathbf{X})}{\mathbb{P}(Y = K|\mathbf{X})} \right] = \mathbf{X}\mathbf{W}_{K-1} \quad (13)$$

where each \mathbf{W}_k is a $p \times 1$ parameter vector.

Exponentiating both sides and solving for $\mathbb{P}(Y = K|\mathbf{X})$ (using the fact that the probabilities have to sum to 1) gives us

$$\mathbb{P}(Y = K|\mathbf{X}) = \frac{1}{1 + \sum_{k=1}^{K-1} e^{\mathbf{X}\mathbf{W}_k}} \quad (14)$$



Equivalently, can represent the multinomial logistic model as

$$\log \mathbb{P}(Y = 1|\mathbf{X}) = \mathbf{XW}_1 - \log(Z) \quad (15)$$

$$\log \mathbb{P}(Y = 2|\mathbf{X}) = \mathbf{XW}_2 - \log(Z) \quad (16)$$

$$\dots = \dots \quad (17)$$

$$\log \mathbb{P}(Y = K|\mathbf{X}) = \mathbf{XW}_K - \log(Z) \quad (18)$$

resulting in the following probabilities

$$\mathbb{P}(Y = 1|\mathbf{X}) = \frac{\exp(\mathbf{XW}_1)}{\sum_{k=1}^K \exp(\mathbf{XW}_k)} \quad (19)$$

$$\mathbb{P}(Y = 2|\mathbf{X}) = \frac{\exp(\mathbf{XW}_2)}{\sum_{k=1}^K \exp(\mathbf{XW}_k)} \quad (20)$$

$$\dots = \dots \quad (21)$$

$$\mathbb{P}(Y = K|\mathbf{X}) = \frac{\exp(\mathbf{XW}_K)}{\sum_{k=1}^K \exp(\mathbf{XW}_k)} \quad (22)$$



This leads us to the softmax function, i.e.

$$\text{softmax}(\mathbf{XW}_1, \dots, \mathbf{XW}_K)_k = \frac{e^{\mathbf{XW}_k}}{\sum_{l=1}^K e^{\mathbf{XW}_l}} \quad (23)$$

Or, more succinctly, we have

$$\text{softmax}(\mathbf{XW}^K)_k = \frac{\exp((\mathbf{XW}^K)_{k.})}{\sum_{k=1}^K \exp((\mathbf{XW}^K)_{k.})} \quad (24)$$

where the $p \times K$ matrix \mathbf{W}^K is simply the (concatenated) matrix of $\mathbf{W}_1, \dots, \mathbf{W}_K$.

This is what multiclass neural networks are modeling!



Recall: In logistic regression we try to maximize the likelihood

- ▶ Equivalent to minimizing the cross-entropy

$$L(y_i, f(\mathbf{X}_i)) = -y_i \log(p_i) - (1 - y_i) \log(1 - p_i), \text{ where (25)}$$

$$p_i = \frac{1}{1 + \exp(-Z_i)} \quad (26)$$

$$Z_i = \mathbf{X}_i \mathbf{W} \quad (27)$$



Recall: In logistic regression we try to maximize the likelihood

- ▶ Equivalent to minimizing the cross-entropy

$$L(y_i, f(\mathbf{X}_i)) = -y_i \log(p_i) - (1 - y_i) \log(1 - p_i), \text{ where} \quad (25)$$

$$p_i = \frac{1}{1 + \exp(-Z_i)} \quad (26)$$

$$Z_i = \mathbf{X}_i \mathbf{W} \quad (27)$$

Can apply *the derivative chain rule* to get our gradient, i.e.

$$\frac{\partial L(y_i, f(\mathbf{X}_i))}{\partial \mathbf{W}} = \frac{\partial L(y_i, f(\mathbf{X}_i))}{\partial p_i} \times \frac{\partial p_i}{\partial Z_i} \times \frac{\partial Z_i}{\partial \mathbf{W}} \quad (28)$$



Recall: In logistic regression we try to maximize the likelihood

- ▶ Equivalent to minimizing the cross-entropy

$$L(y_i, f(\mathbf{X}_i)) = -y_i \log(p_i) - (1 - y_i) \log(1 - p_i), \text{ where} \quad (25)$$

$$p_i = \frac{1}{1 + \exp(-Z_i)} \quad (26)$$

$$Z_i = \mathbf{X}_i \mathbf{W} \quad (27)$$

Can apply *the derivative chain rule* to get our gradient, i.e.

$$\frac{\partial L(y_i, f(\mathbf{X}_i))}{\partial \mathbf{W}} = \frac{\partial L(y_i, f(\mathbf{X}_i))}{\partial p_i} \times \frac{\partial p_i}{\partial Z_i} \times \frac{\partial Z_i}{\partial \mathbf{W}} \quad (28)$$

Which gives us

$$\frac{\partial L(y_i, f(\mathbf{X}_i))}{\partial \mathbf{W}} = \mathbf{X}_i (y_i - p_i) \quad (29)$$



Neural networks are simply a generalization of the logistic regression case, e.g. for

$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\sigma(\mathbf{X}\mathbf{W}_1)\mathbf{W}_2) \quad (30)$$



Neural networks are simply a generalization of the logistic regression case, e.g. for

$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\sigma(\mathbf{X}\mathbf{W}_1)\mathbf{W}_2) \quad (30)$$

Our loss is

$$L(y_i, f(\mathbf{X}_i)) = -y_i \log(p_i) - (1 - y_i) \log(1 - p_i), \text{ where} \quad (31)$$

$$p_i = \frac{1}{1 + \exp(-Z_{2,i})} \quad (32)$$

$$Z_{2,i} = h_i \mathbf{W}_2 \quad (33)$$

$$h_i = \frac{1}{1 + \exp(-Z_{1,i})} \quad (34)$$

$$Z_{1,i} = \mathbf{X}\mathbf{W}_1 \quad (35)$$



Our loss is

$$L(y_i, f(\mathbf{X}_i)) = -y_i \log(p_i) - (1 - y_i) \log(1 - p_i), \text{ where} \quad (36)$$

$$p_i = \frac{1}{1 + \exp(-Z_{2,i})} \quad (37)$$

$$Z_{2,i} = h_i \mathbf{W}_2 \quad (38)$$

$$h_i = \frac{1}{1 + \exp(-Z_{1,i})} \quad (39)$$

$$Z_{1,i} = \mathbf{X} \mathbf{W}_1 \quad (40)$$

Applying *the derivative chain rule*:

$$\frac{\partial L(y_i, f(\mathbf{X}_i))}{\partial \mathbf{W}_2} = \frac{\partial L(y_i, f(\mathbf{X}_i))}{\partial p_i} \times \frac{\partial p_i}{\partial Z_{2,i}} \times \frac{\partial Z_{2,i}}{\partial \mathbf{W}_2}$$

$$\frac{\partial L(y_i, f(\mathbf{X}))}{\partial \mathbf{W}_1} = \frac{\partial L(y_i, f(\mathbf{X}))}{\partial p_i} \times \frac{\partial p_i}{\partial Z_{2,i}} \times \frac{\partial Z_{2,i}}{\partial h_i} \times \frac{\partial h_i}{\partial Z_{1,i}} \times \frac{\partial Z_{1,i}}{\partial \mathbf{W}_1}$$



Our gradient is estimated using our data, i.e.
 $(y_i, \mathbf{X}_i) : i = 1, 2, \dots, n$.



Our gradient is estimated using our data, i.e.

$(y_i, \mathbf{X}_i) : i = 1, 2, \dots, n.$

We can estimate it using, e.g.

- ▶ **Stochastic gradient descent:** estimating our (full) gradient using just one observation.
- ▶ **Gradient descent:** estimating our (full) gradient using all observations.
- ▶ **Mini-batch gradient descent:** using a (random) subsample of our observations.

Each will trade off between variance for the gradient and memory size.



Our gradient is estimated using our data, i.e.

$(y_i, \mathbf{X}_i) : i = 1, 2, \dots, n$.

We can estimate it using, e.g.

- ▶ **Stochastic gradient descent**: estimating our (full) gradient using just one observation.
- ▶ **Gradient descent**: estimating our (full) gradient using all observations.
- ▶ **Mini-batch gradient descent**: using a (random) subsample of our observations.

Each will trade off between variance for the gradient and memory size.

When done iteratively, we'll typically specify a stopping point (e.g. by using a dev set).



The use of non-linearities results in multiple minima, tendency to overfit, and can be unstable.

Some considerations to make:

- ▶ Set initial weight values near zero.
- ▶ Over parameterize and regularize heavily.
- ▶ Standardize input features.
- ▶ Use a dev set and stop training earlier.
- ▶ Try out different weight randomizations and take the one with the lowest (validated) error.
 - ▶ Or average the predictions (or apply bagging).



Estimating neural network parameters simply requires '*propagating back*' errors.

- ▶ We're just applying (matrix) multiplications
 - ▶ GPU's can be very good for this
- ▶ Matrix multiplications can get pretty big (for large networks)
 - ▶ Commonly not worth it to use the Hessian
- ▶ Should be careful with large values going into sigmoid activations
 - ▶ Results in saturated gradients



Hornik's theorem

Whenever the activation function is continuous, bounded, and non-constant, then, for arbitrary compact subsets $X \subseteq \mathbb{R}^k$, standard multilayer feedforward networks can approximate any continuous function on X arbitrarily well with respect to uniform distance, provided that sufficiently many hidden units are available.



Hornik's theorem

Whenever the activation function is continuous, bounded, and non-constant, then, for arbitrary compact subsets $X \subseteq \mathbb{R}^k$, standard multilayer feedforward networks can approximate any continuous function on X arbitrarily well with respect to uniform distance, provided that sufficiently many hidden units are available.

In words: A 2-layer neural network with enough hidden nodes can closely approximate any continuous function $f(x)$.

References:

Cybenko (1989) "Approximations by superpositions of sigmoidal function"

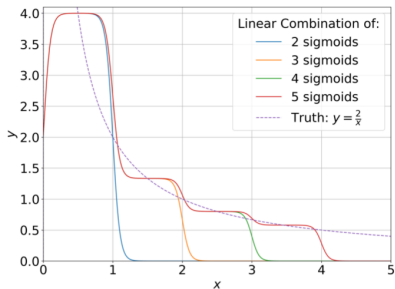
Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks"

Leshno and Schocken (1993) "Multilayer Feedforward Networks with Non-Polynomial Activation Functions Can Approximate Any Function"

Universal approximation theorem



Given enough hidden nodes, we can approximate any function.



Check out this *visual example* of this.



Some caveats to the theorem

- ▶ We're approximating the function within some bound, i.e.
 $|\hat{f}_n(x) - f(x)| < \epsilon.$
- ▶ Result is meant for *continuous* functions on *compact* subsets of \mathbb{R} .
- ▶ Nothing is guaranteed on the how quickly we can learn the function's parameters.
- ▶ Other function estimators also do a good job approximating!



For non-linear functions,

Logistic regression: expand our feature set via transformations

Neural network: define the model non-linearly



For non-linear functions,

Logistic regression: expand our feature set via transformations

- ▶ We have to specify the feature transformations

Neural network: define the model non-linearly

- ▶ The model learns the feature transformations



For non-linear functions,

Logistic regression: expand our feature set via transformations

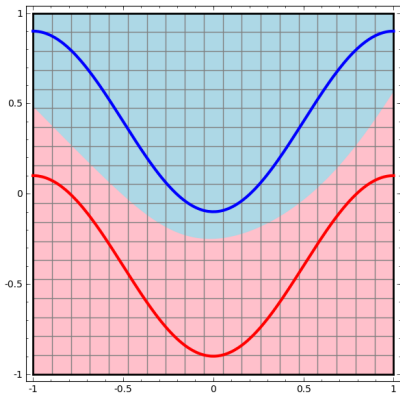
- ▶ We have to specify the feature transformations

Neural network: define the model non-linearly

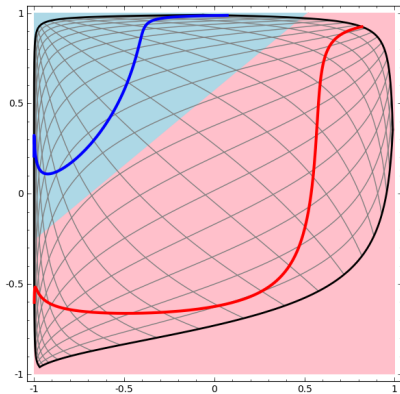
- ▶ The model learns the feature transformations
- ▶ *This helps us greatly when dealing with abstract or high dimensional problems (e.g. images & text)!*



How do the feature transformations get learned?



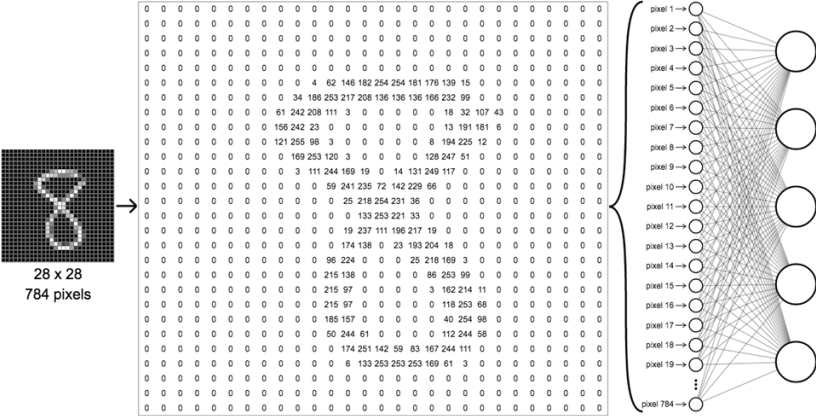
Original representation of curves



Hidden layer representation of curves

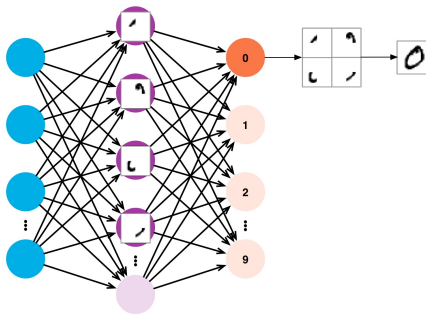
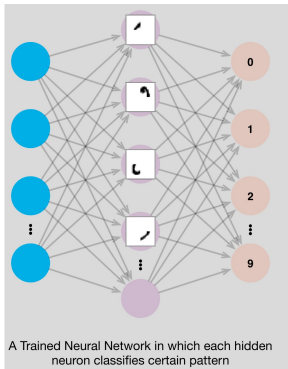
Well demonstrated by *Chris Olah's blog*.

Example: The MNIST data



A linear model (e.g. for multinomial logistic regression)

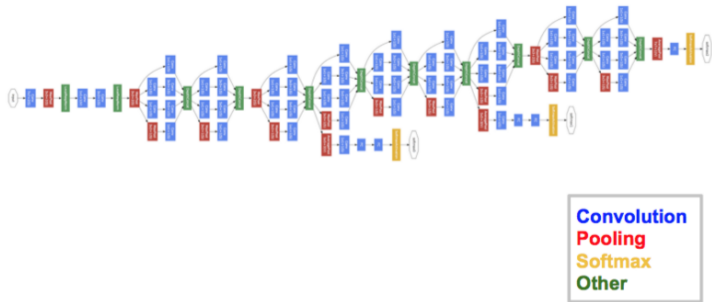
Example: The MNIST data



Neural network's learned (kernel) features.



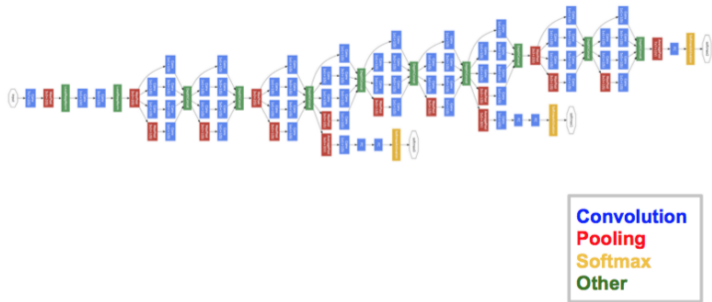
Knowing that we can simply back propagate errors via multiplication opens many doors for us, e.g.



Google's InceptionNet architecture



Knowing that we can simply back propagate errors via multiplication opens many doors for us, e.g.



Google's InceptionNet architecture

Problem: *large networks are vulnerable to vanishing/exploding gradients.*



Recall: Our gradient is simply a product of partial derivatives.

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



Neural Network

Example neural network and gradient.



Recall: Our gradient is simply a product of partial derivatives.

$$\frac{\partial C}{\partial b_1} = \sigma'(z_1) \times w_2 \times \sigma'(z_2) \times w_3 \times \sigma'(z_3) \times w_4 \times \sigma'(z_4) \times \frac{\partial C}{\partial a_4}$$



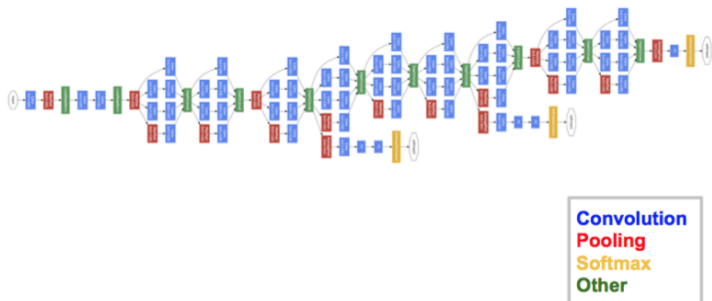
Neural Network

Example neural network and gradient.

Question: What is the derivative of the sigmoid function?



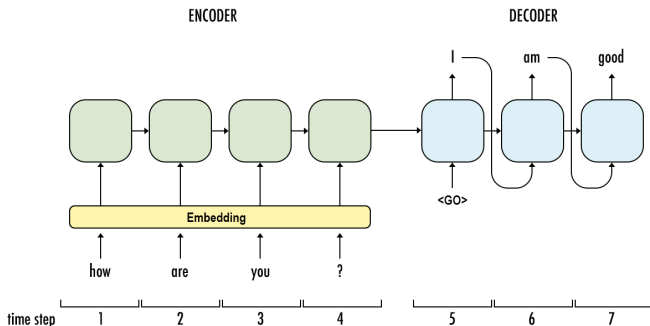
Google's InceptionNet address this via strategically placed loss functions.



Google's InceptionNet architecture



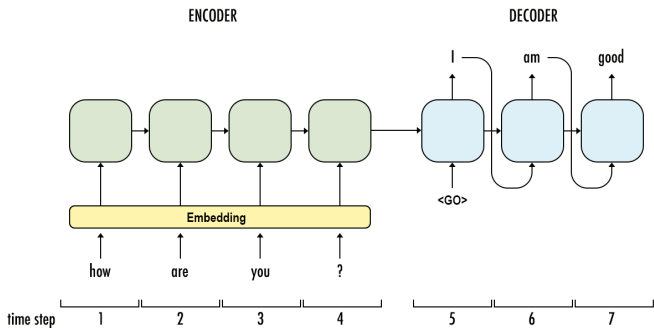
We can also train models end to end, e.g.



Encoder-decoder architecture



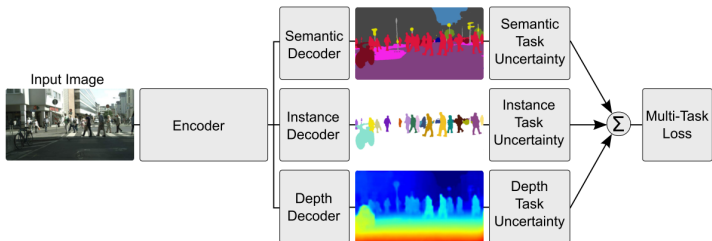
We can also train models end to end, e.g.



Encoder-decoder architecture

Or: in a modular fashion, e.g. pre-training.

Or over multiple tasks (i.e. multi-task learning), e.g.



Kendall et al. 2017's multi-task model

Many researchers will create unique architectures for specific problems, e.g. *Instacart*



The prediction problem

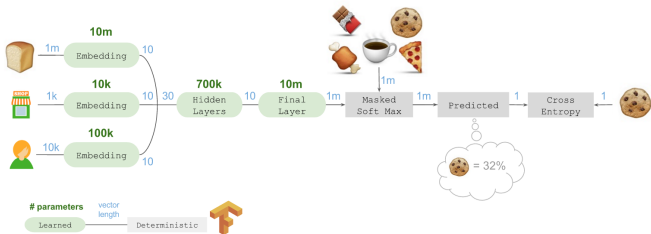


Model generalizations

Many researchers will create unique architectures for specific problems, e.g. *Instacart*



The prediction problem



The initial solution



Another example using the Netflix data.

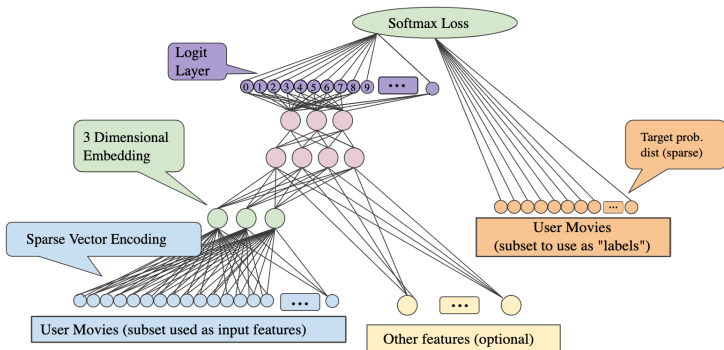
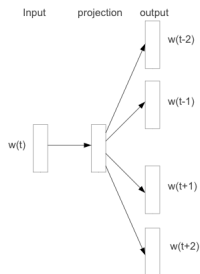


Figure 5. A sample DNN architecture for learning movie embeddings from collaborative filtering data.



The skip gram model



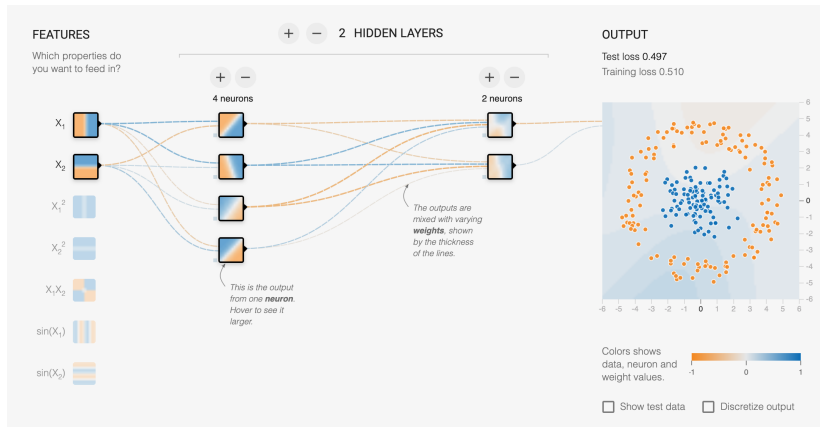
$$\frac{1}{T} \sum_{t=1}^T \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_t)$$

$$p(w_o | w_I) = \frac{\exp(v'_{w_o} \top v_{w_I})}{\sum_{w=1}^W \exp(v'_w \top v_{w_I})}$$

The tensorflow playground



An *interactive demo* that allows you to play with a neural network.





Neural net related core topics:

- ▶ Weight initializations
- ▶ Activation functions
- ▶ Optimization functions
- ▶ Loss functions
- ▶ Normalization
- ▶ Regularization / dropout
- ▶ Model architectures
- ▶ Hyperparameter optimization
- ▶ Bayesian neural networks
- ▶ Computation graphs
- ▶ Software / platforms
- ▶ Encoding / adding outside knowledge
- ▶ Hardware accelerators



Neural net applied topics:

- ▶ Computer vision
- ▶ Natural language processing
- ▶ Signal processing
- ▶ Generative models
- ▶ Unsupervised learning
- ▶ Reinforcement learning
- ▶ One/Zero shot learning
- ▶ Transfer learning
- ▶ Auto-ML
- ▶ Memory Augmented Neural Networks



- [1] ESL. Chapter 11
- [2] Pancha N, Zhai A, Leskovec J, Rosenberg C. PinnerFormer: Sequence Modeling for User Representation at Pinterest. arXiv 2022.