Knowledge Discovery & Data Mining Neural Networks —

<u>vong.zhuang@gvsu.edu</u>

Yong Zhuang

Instructor: Yong Zhuang

Based on the original version by Professor Yizhou Sun

AlphaGo

- Artificial Neural Networks (ANN)
 - History
 - Architecture
 - Activation function
 - Loss function
 - Optimization
 - Regularization
 - Training
 - Stochastic gradient descent + chain rule
 - Backpropagation

Source: https://www.nytimes.com/2016/03/10/world/asia/google-alphago-lee-se-dol.html; https://www.bbc.com/news/technology-35785875





A little bit of History



Yong Zhuang

Source: https://i.pinimg.com/originals/6a/f0/3c/6af03c5026bb680ebe6d8db4bdbb8428.jpg





ImageNet Challenge



Yong Zhuang

Source: https://semiengineering.com/new-vision-technologies-for-real-world-applications/





AlphaGo

Master of Go Board Game Is Walloped by Google Computer Program



By Choe Sang-Hun and John Markoff

March 9, 2016

SEOUL, South Korea — Computer, one. Human, zero.

A <u>Google</u> computer program stunned one of the world's top players on Wednesday in a round of Go, which is believed to be the most complex board game ever created.

Source: https://www.bbc.com/news/technology-35785875









OpenAI announces ChatGPT successor GPT-4

14 March 2023

By Ben Derico and Zoe Kleinman, BBC News



OpenAI has released GPT-4, the latest version of its hugely popular artificial intelligence chatbot ChatGPT. Source: <u>https://www.bbc.com/news/technology-64959346</u>

Yong Zhuang

Knowledge Discovery & Data Mining

< Share





Artificial Neural Networks (ANN)

Consider humans

- Number of neurons ~10¹⁰
- Connections per neuron ~10⁴⁻⁵
- Neuron switching time ~.001 second
- Scene recognition time ~.1 second
- 100 inference steps doesn't seem like enough -> parallel computation

Artificial neural networks

- Many neuron-like threshold switching units
- Many weighted interconnections among units
- Highly parallel, distributed process
- Emphasis on tuning weights automatically



Knowledge Discovery & Data Mining



7

Important Concepts

- Architecture
- Activation function
- Loss function
- Optimization
- Regularization



Architecture

- Decide the network topology:
 - \circ # of units in the input layer,
 - \circ # of hidden layers (if > 1),
 - \circ # of units in each hidden layer,
 - the unit types,
 - connection between layers,
 - and # of units in the output layer
- to be learned.

• Architecture specifies the function that maps input to output, which contains parameters





Activation function

- An activation function $f(\cdot)$ in the output layer can control the nature of the output (e.g., probability value in [0, 1])
- Activation functions bring nonlinearity into hidden layers, which increases the complexity of the model.
- Good activation functions should be **differentiable** for optimization purpose

Neural Network Activation Functions: a small subset!



Loss Functions

- How good are the outputs compared with the labels (target)?
 - Empirical risk

• w: parameters in the model

 Loss function: difference between actual value and predicted value $l(y, \hat{y})$

Examples of Loss Functions

• Squared error $l(y, \hat{y}) = (y - \hat{y})^2$

• (Binary) cross entropy loss $\begin{aligned} l(y, \hat{y}) &= -y \log \hat{y} - (1 - y) \log(1 - \hat{y}) \\ y \in \{0, 1\}, \hat{y} \in [0, 1] \end{aligned}$

• Hinge loss $l(y, \hat{y}) = \max(0, 1 - y\hat{y})$ $y \in \{-1,1\}, \hat{y} \in (-\infty, +\infty)$



11

Optimization

- Given a training dataset, minimize the empirical risk • Find w, such that $\mathcal{L}(w) = \frac{1}{n} \sum_{i} l(y^{(i)}, \hat{y}^{(i)})$, where $\hat{y}^{(i)} = f(x^{(i)}, w)$ is minimized.
- Solution:
 - Stochastic gradient descent + chain rule = **backpropagation** 0





Regularization

- Avoid overfitting
- Explicit Regularization
 - L2/L1 regularization
 - **Dropout**
 - **Data Augmentation** 0
- Implicit Regularization
 - **Early stopping** Ο
 - **Batch Normalization** 0

SGD Ο

Source: https://medium.com/analytics-vidhya/a-simple-introduction-to-dropout-regularization-with-code-5279489dda1e



Knowledge Discovery & Data Mining



13

Single Unit: Perceptron

- node 2 is larger, neuron 1 has a larger influence on this neuron.
- Input layer: is made of artificial input neurons and takes the initial data into the system for further processing. • Weight: It represents the dimension or strength of the connection between units. If the weight from node 1 to
- Bias: It is the same as the intercept added in a linear equation. It is an additional parameter which task is to modify the output along with the weighted sum of the input to the other neuron.
- Net sum: It calculates the total sum.
- Activation Function: A neuron can be activated or not, is determined by an activation function. The activation function calculates a weighted sum and further adding bias with it to give the result.

X1

Inputs









Single Unit: Perceptron

- Architecture:
 - A single neuron
- Activation function
 - Training: identity function
 - Inference: sign function/step function
- Loss function
 - $o l(y, \hat{y}) = max(0, -y\hat{y})$
- Optimization
- η : learning rate



$\circ w \leftarrow w + \eta y^{(i)} x^{(i)}$, for a misclassified training data point $(x^{(i)}, y^{(i)})$, i.e., $y^{(i)} w^T x^{(i)} \leq 0$





Example: 1 for "Y" and -1 for "N"; $\eta = 0.9$

X0	X1	X2	True Label	Predicted Label	W (before update)	W (after update)
1	0	1	Y	Ν	(0.0, 0.0, 0.0)	(0.9, 0.0, 0.9)
1	1	1	Ν	Υ	(0.9, 0.0, 0.9)	(0.0, -0.9, 0.0)
1	0	0	Y	Ν	(0.0, -0.9, 0.0)	(0.9, -0.9, 0.0)
1	1	0	Y	Ν	(0.9, -0.9, 0.0)	(1.8, 0.0, 0.0)
1	0	1	Y	Υ	(1.8, 0.0, 0.0)	(1.8, 0.0, 0.0)
1	1	1	N	Υ	(1.8, 0.0, 0.0)	(0.9, -0.9, -0.9)
1	0	0	Y	Υ	(0.9, -0.9, -0.9)	(0.9, -0.9, -0.9)
1	1	0	Y	Ν	(0.9, -0.9, -0.9)	(1.8, 0.0, -0.9)
1	0	1	Y	Υ	(1.8, 0.0, -0.9)	(1.8, 0.0, -0.9)
1	1	1	Ν	Υ	(1.8, 0.0, -0.9)	(0.9, -0.9, -1.8)
1	0	0	Y	Υ	(0.9, -0.9, -1.8)	(0.9, -0.9, -1.8)
1	1	0	Y	Ν	(0.9, -0.9, -1.8)	(1.8, 0.0, -1.8)





Single Unit: Logistic Regression

- Architecture:
 - A single neuron
- Activation function
 - Sigmoid function
- Loss function
 - $\circ l(y, \hat{y}) = -y \log \hat{y} (1-y) \log(1-\log \hat{y})$
 - \circ **Note** \hat{y} is the predicted probability of taking class 1.
- Optimization
 - $\circ w \leftarrow w + \eta (y^{(i)} \sigma(w^T x^{(i)})) x^{(i)}$, for a training data point $(x^{(i)}, y^{(i)})$.
- η : learning rate

Knowledge Discovery & Data Mining



17

A Multi-Layer Feed-Forward Neural Network



Yong Zhuang

Knowledge Discovery & Data Mining

18

How A Multi-Layer Neural Network Works

- The inputs to the network correspond to the attributes measured for each training tuple Inputs are fed simultaneously into the units making up the input layer They are then weighted and fed simultaneously to a hidden layer
- The number of hidden layers is arbitrary
- The weighted outputs of the last hidden layer are input to units making up the output layer, which emits the network's prediction
- The network is feed-forward: None of the weights cycles back to an input unit or to an output unit of a previous layer
- From a math point of view, networks perform nonlinear regression: Given enough hidden units and enough training samples, they can closely approximate any continuous function













Learning by Backpropagation

- Iteratively process a set of training tuples & compare the network's prediction with the actual known target value
- For each training tuple, the weights are modified to minimize the loss function between the network's prediction and the actual target value, say mean squared error Stochastic gradient descent + chain rule
- Modifications are made in the "backwards" direction: from the output layer, through each hidden layer down to the first hidden layer, hence "backpropagation"







Recap: Chain Rule

differentiable functions f and g in terms of the derivatives of f and g.

If y = f(u) and u = g(x) are both differentiable functions, then

$\frac{c}{c}$	$\frac{ly}{lx} =$	$= \frac{d}{d}$	
$\frac{dy}{dx} = de$	rivative	of y v	~
$\frac{dy}{du} = de$	erivative	of y v	N
$\frac{du}{dx} = de$	erivative	of u \	٨

• The chain rule is a formula that expresses the derivative of the composition of two

 $y \, du$ u dx

vith respect to x

vith respect to u

vith respect to x



Example

• Loss function: $\mathcal{L} = \frac{1}{2} ||y - \hat{y}||^2$



Yong Zhuang



Gradient for Layer 3 (Last Layer)

- Stochastic gradient for $W_{ij}^{(3)}$ and $b_{ij}^{(3)}$
 - Recall:
 - $\mathcal{L} = \frac{1}{2} \left| \left| \mathbf{y} \mathbf{x}^{(3)} \right| \right|^2 = \frac{1}{2} \sum_i \left(y_i x_i^{(3)} \right)^2$ • $x_i^{(3)} = f^{(3)}(z_i^{(3)})$ • $z_i^{(3)} = \sum_j W_{ij}^{(3)} x_j^{(2)} + b_i^{(3)}$

$$\frac{\partial \mathcal{L}}{\partial W_{ij}^{(3)}} = \frac{\partial \mathcal{L}}{\partial x_i^{(3)}} \frac{\partial x_i^{(3)}}{\partial z_i^{(3)}} \frac{\partial z_i^{(3)}}{\partial W_{ij}^{(3)}} = -(y_i - x_i^{(3)}) f'^{(3)} \left(z_i^{(3)} \right) x_j^{(2)}$$

$$\frac{\partial \mathcal{L}}{\partial b_i^{(3)}} = \frac{\partial \mathcal{L}}{\partial x_i^{(3)}} \frac{\partial x_i^{(3)}}{\partial z_i^{(3)}} \frac{\partial z_i^{(3)}}{\partial b_i^{(3)}} = -(y_i - x_i^{(3)}) f'^{(3)} \left(z_i^{(3)} \right)$$





Gradient for Layer 2

- Stochastic gradient for $W_{jk}^{(2)}$
 - Recall:
 - $\mathcal{L} = \frac{1}{2} \left| \left| \mathbf{y} \mathbf{x}^{(3)} \right| \right|^2 = \frac{1}{2} \sum_i \left(y_i x_i^{(3)} \right)^2$ • $x_i^{(3)} = f^{(3)} \left(z_i^{(3)} \right); \ z_i^{(3)} = \sum_j W_{ij}^{(3)} x_j^{(2)} + b_i^{(3)}$ • $x_j^{(2)} = f^{(2)} \left(z_j^{(2)} \right); \ z_j^{(2)} = \sum_k W_{jk}^{(2)} x_k^{(1)} + b_j^{(2)}$

$$\cdot \frac{\partial \mathcal{L}}{\partial W_{jk}^{(2)}} = \sum_{i} \underbrace{\frac{\partial \mathcal{L}}{\partial x_{i}^{(3)}} \frac{\partial x_{i}^{(3)}}{\partial z_{i}^{(3)}} \frac{\partial z_{i}^{(3)}}{\partial x_{j}^{(2)}} \frac{\partial x_{j}^{(2)}}{\partial z_{j}^{(2)}} \frac{\partial z_{j}^{(2)}}{\partial W_{jk}^{(2)}} }{\delta W_{jk}^{(2)}}$$

$$= \sum_{i} -(y_{i} - x_{i}^{(3)})f'^{(3)}\left(z_{i}^{(3)}\right) W_{ij}^{(3)}f'^{(2)}\left(z_{j}^{(2)}\right) x_{k}^{(1)}$$

$$\frac{\delta_{i}^{(3)}}{\delta_{i}^{(2)}}$$





Gradient for Layer 1

• Stochastic gradient for $W_{ks}^{(1)}$

 $\frac{\partial \mathcal{L}}{\partial W_{ks}^{(1)}} = \sum_{j} \delta_{j}^{(2)} W_{jk}^{(2)} f'^{(1)} \left(z_{k}^{(1)} \right) x_{s}^{(0)}$





Backpropagation Steps to Learn Weights

- Initialize weights to small random numbers, associated with biases • **Repeat** until terminating condition meets
- For each training example
 - **Propagate the inputs forward** (by applying activation function)
 - For layer I = 1: L
 - Calculate $z^{(l)} = W^{(l)} x x^{(l-1)} + b^{(l)}$
 - Calculate $x^{(l)} = f^{(l)} z^{(l)}$ (elementwise active
 - Backpropagate the error (by updating weighted)
 - Calculate $\delta^{(L)}$, Update $W^{(L)}$ and $b^{(L)}$ based
 - For layer I = L-1: 1
 - Calculate $\delta^{(l)} = (W^{(l+1)})^T \delta^{(l+1)} f'^{(l)}(z^{(l)})$
 - Update $W^{(L)}$ and $b^{(L)}$ based on $\frac{\partial \mathcal{L}}{\partial W^{(l)}} = \delta^{(l)}$
- Terminating condition (convergence, max itera

ation)
ghts and biases)
on
$$\frac{\partial \mathcal{L}}{\partial W^{(L)}} = \boldsymbol{\delta}^{(L)} (\boldsymbol{x}^{(L-1)})^T$$
 and $\frac{\partial \mathcal{L}}{\partial \boldsymbol{b}^{(L)}} = \boldsymbol{\delta}^{(L)}$

$$(x^{(l-1)})^T$$
 and $\frac{\partial \mathcal{L}}{\partial b^{(l)}} = \delta^{(l)}$
ation, etc.)





Neural Network as a Classifier

• Weakness

- Long training time
- Require a number of hyper-parameters typically best determined empirically, e.g., the network topology or "structure."
- Poor interpretability: Difficult to interpret the symbolic meaning behind the learned weights and of "hidden units" in the network

• Strength

- High tolerance to noisy data
- Successful on an array of real-world data, e.g., hand-written letters
- Algorithms are inherently parallel
- Techniques have recently been developed for the extraction of rules from trained neural networks
- Deep neural network is powerful





Example: Digits Recognition

• The architecture of the used neural network

• What each neurons are doing?







Summary

- Artificial Neural Networks (ANN)
 - Architecture
 - Activation function
 - Loss function
 - Optimization
 - Regularization
 - Training
 - Stochastic gradient descent + chain rule
 - Backpropagation

