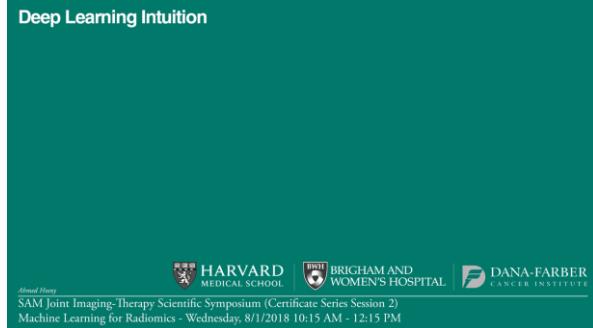
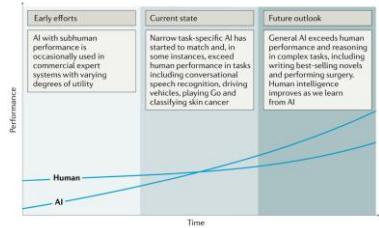


Deep Learning Intuition

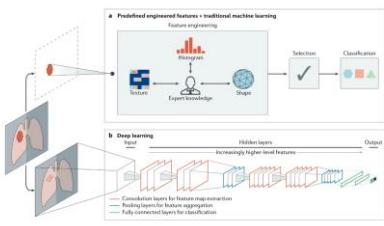


Deep Learning



Ahmed Hanif, Christian Partner, John Quackenbush, Lawrence H Schwartz and Hugo J W L Attwell
Artificial Intelligence in Radiology
Nature Reviews Cancer - 2018

Deep Learning



Ahmed Hanif, Christian Partner, John Quackenbush, Lawrence H Schwartz and Hugo J W L Attwell
Artificial Intelligence in Radiology
Nature Reviews Cancer - 2018

What is the intuition behind neural networks?

How do neural networks learn?

How to train neural networks?

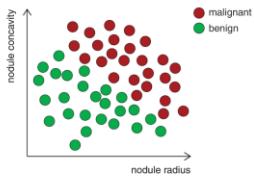
What is the intuition behind neural networks?

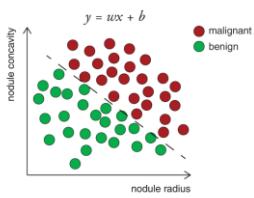
How do neural networks learn?

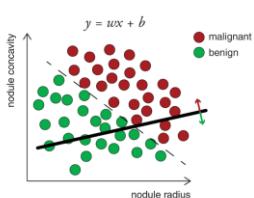
How to train neural networks?

Machine Learning: 4 Main Components

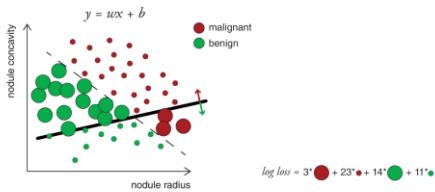
- Data
- Model/Representation
- Cost/Error/Loss of model
- Model Optimizer

Logistic Regression

Logistic Regression

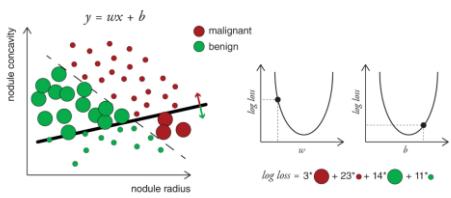
Logistic Regression

Logistic Regression



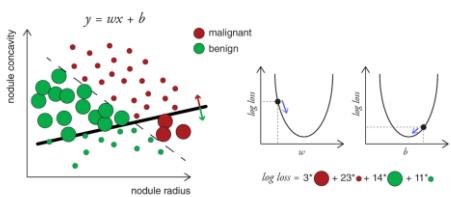
$$\log \text{loss} = 3^* \text{red} + 23^* \text{black} + 14^* \text{green} + 11^* \text{blue}$$

Logistic Regression



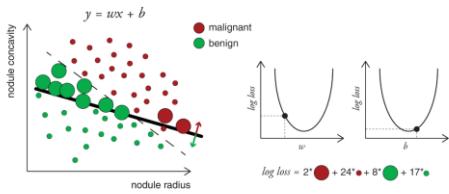
$$\log \text{loss} = 3^* \text{red} + 23^* \text{black} + 14^* \text{green} + 11^* \text{blue}$$

Logistic Regression

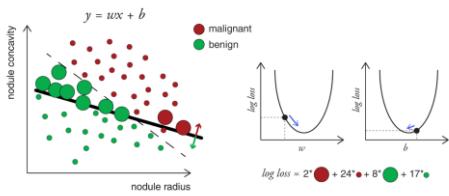


$$\log \text{loss} = 3^* \text{red} + 23^* \text{black} + 14^* \text{green} + 11^* \text{blue}$$

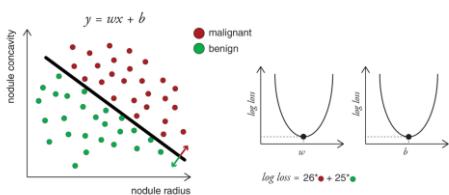
Logistic Regression

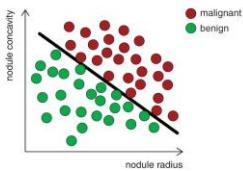


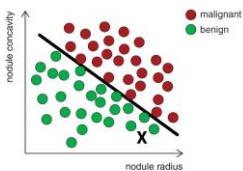
Logistic Regression

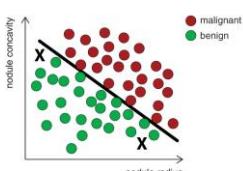


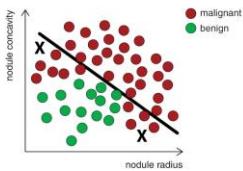
Logistic Regression

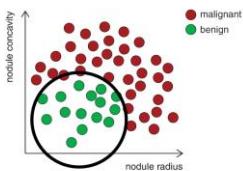


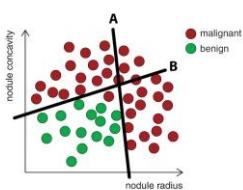
Dealing with Edge Conditions

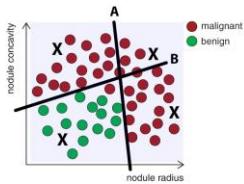
Dealing with Edge Conditions

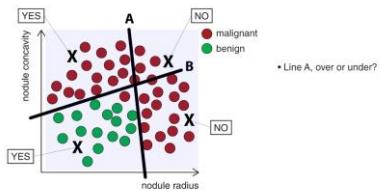
Dealing with Edge Conditions

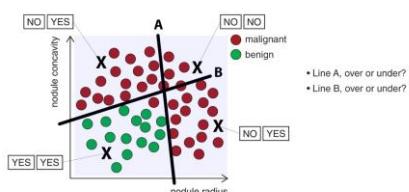
Dealing with Edge Conditions

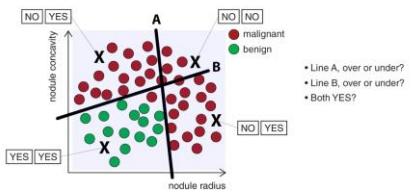
Dealing with Edge Conditions

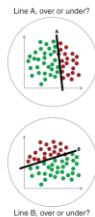
Dealing with Edge Conditions

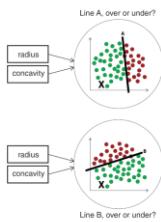
Quadrant Questioning

Quadrant Questioning

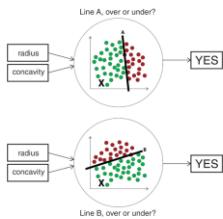
Quadrant Questioning

Quadrant Questioning

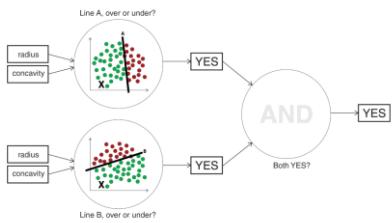
Graph Representation

Graph Representation

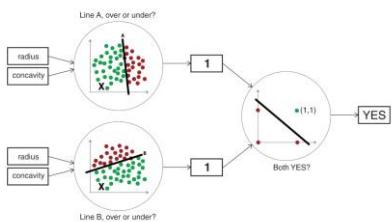
Graph Representation

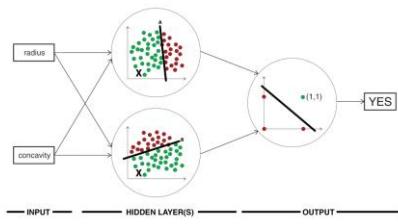


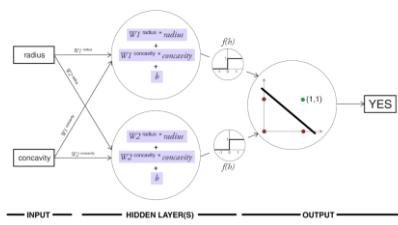
Graph Representation

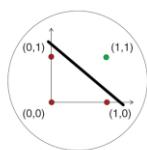


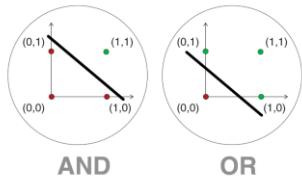
Graph Representation

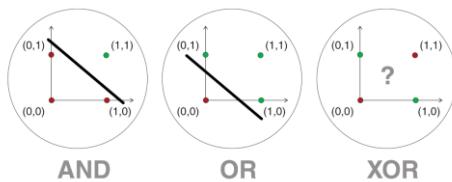


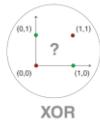
A Neural Network

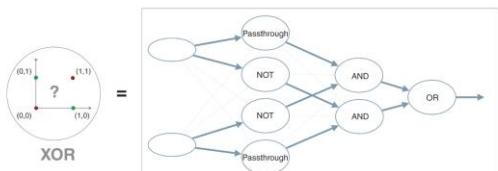
A Neural Network

XOR Perceptron**AND**

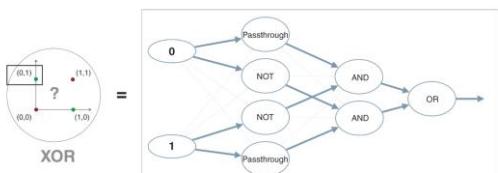
XOR Perceptron**AND****OR**

XOR Perceptron**AND****OR****XOR**

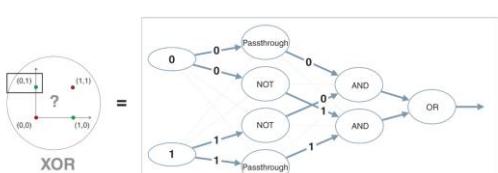
XOR Perceptron**XOR**

XOR Perceptron

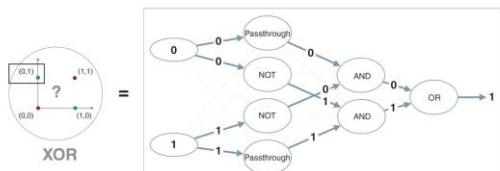
Lesson 2: Intro To Neural Networks
Udacity Deep Learning Nanodegree Program

XOR Perceptron

Lesson 2: Intro To Neural Networks
Udacity Deep Learning Nanodegree Program

XOR Perceptron

Lesson 2: Intro To Neural Networks
Udacity Deep Learning Nanodegree Program

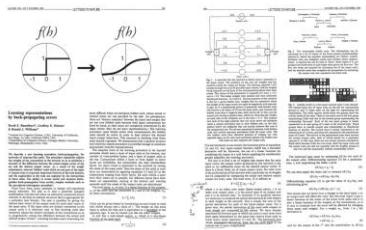
XOR Perceptron

Lesson 2: Intro To Neural Networks
Udacity Deep Learning Nanodegree Program

What is the intuition behind neural networks?

How do neural networks learn?

How to train neural networks?

Backpropagation & Gradient Descent in Neural Networks

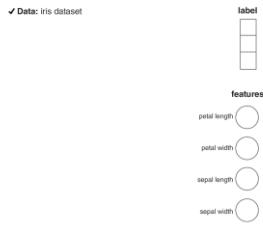
David E. Rumelhart, Geoffrey E. Hinton & Ronald J. Williams
Learning Representations by Back-propagating Errors
Nature - 1986

Iris Dataset

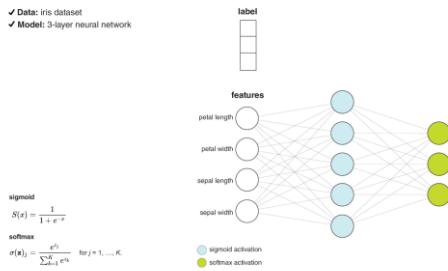


R.A. Fisher
The Use of Multiple Measurements in Taxonomic Problems
Annual Eugenics - 1936

How Neural Networks Learn

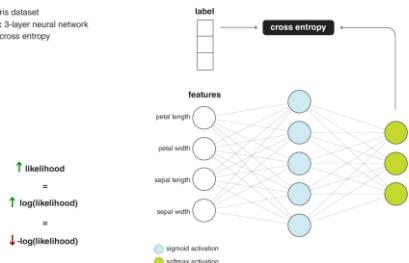


How Neural Networks Learn



How Neural Networks Learn

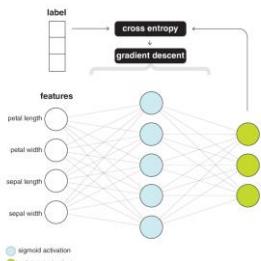
- ✓ Data: iris dataset
 - ✓ Model: 3-layer neural network
 - ✓ Loss: cross entropy



$$\begin{array}{l} \text{↑ likelihood} \\ = \\ \text{↑ log(likelihood)} \\ = \\ \text{↓ -log(likelihood)} \end{array}$$

How Neural Networks Learn

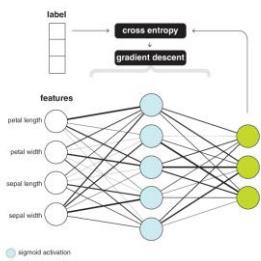
- ✓ Data: iris dataset
 - ✓ Model: 3-layer neural network
 - ✓ Loss: cross entropy
 - ✓ Optimizer: gradient descent



How Neural Networks Learn

- ✓ Data: iris dataset
 - ✓ Model: 3-layer neural network
 - ✓ Loss: cross entropy
 - ✓ Optimizer: gradient descent

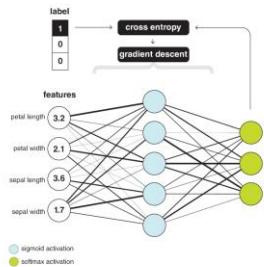
1. parameter initialization



How Neural Networks Learn

- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

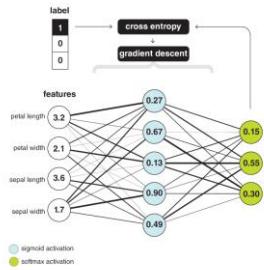
1. parameter initialization
2. data input



How Neural Networks Learn

- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

1. parameter initialization
2. data input
3. forward propagation



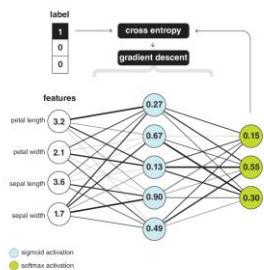
How Neural Networks Learn

- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation

current	
1	0.15
0	0.55
0	0.30

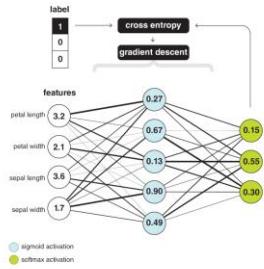
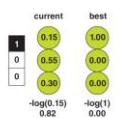
-log(0.15)
0.82



How Neural Networks Learn

- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

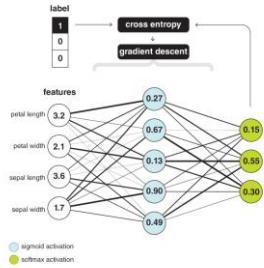
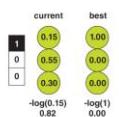
1. parameter initialization
2. data input
3. forward propagation
4. loss calculation



How Neural Networks Learn

- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation



How Neural Networks Learn

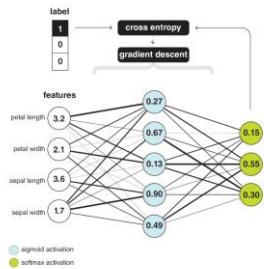
- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates

$\Delta w \leftarrow -\text{gradient}$

$$= -\frac{\partial L}{\partial w}$$

$$= -\eta \frac{\partial L}{\partial w}$$



How Neural Networks Learn

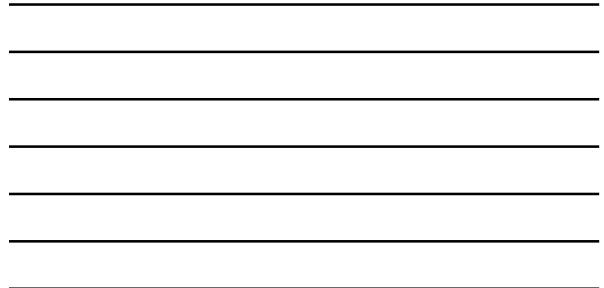
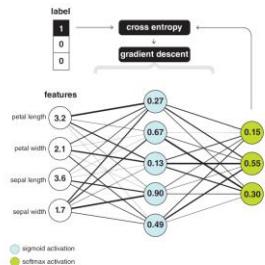
- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates

$$\Delta w \leftarrow -\text{gradient}$$

$$= -\frac{\partial L}{\partial w}$$

$$= -\eta \frac{\partial L}{\partial w}$$

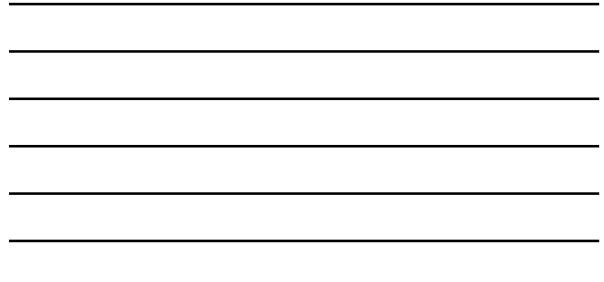
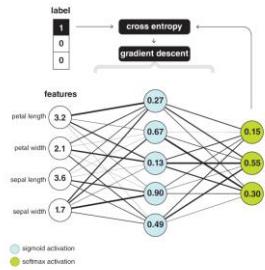


How Neural Networks Learn

- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates

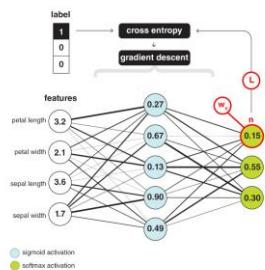
$f(f(\text{input}))$



How Neural Networks Learn

- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates



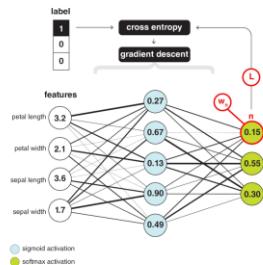
How Neural Networks Learn

- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates

$$\Delta w_n = -\eta \frac{\partial L}{\partial w_n}$$

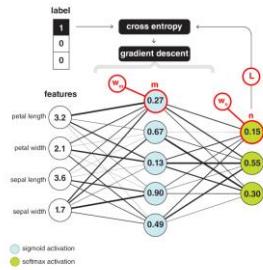
$$= -\eta \frac{\partial L}{\partial s_o} \frac{\partial s_o}{\partial n} \frac{\partial n}{\partial w_n}$$



How Neural Networks Learn

- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates



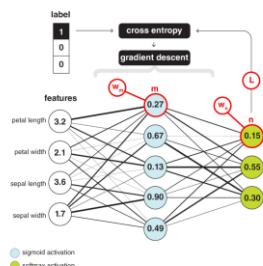
How Neural Networks Learn

- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates

$$\Delta w_m = -\eta \frac{\partial L}{\partial w_m}$$

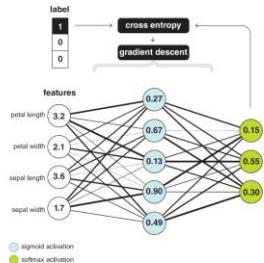
$$= -\eta \frac{\partial L}{\partial s_o} \frac{\partial s_o}{\partial n} \frac{\partial n}{\partial i} \frac{\partial i}{\partial m} \frac{\partial m}{\partial w_m}$$



How Neural Networks Learn

- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

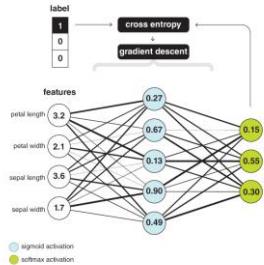
1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates



How Neural Networks Learn

- ✓ Data: iris dataset
- ✓ Model: 3-layer neural network
- ✓ Loss: cross entropy
- ✓ Optimizer: gradient descent

1. parameter initialization
2. data input
3. forward propagation
4. loss calculation
5. backpropagation + updates
6. repeat 2,3,4, & 5



Gradient Descent Flavors

vanilla gradient descent - entire dataset
stochastic gradient descent - random batch of samples (IID)
online gradient descent - (need not be IID)

Gradient Descent Flavors

vanilla gradient descent - entire dataset

stochastic gradient descent - random batch of samples (IID)

online gradient descent - (need not be IID)

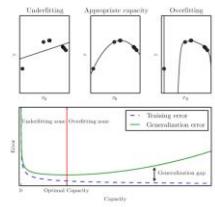
learning rate **batch size** **# of epochs**

What is the intuition behind neural networks?

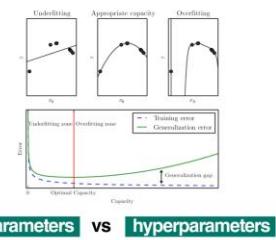
How do neural networks learn?

How to train neural networks?

The Perfect Fit



The Perfect Fit



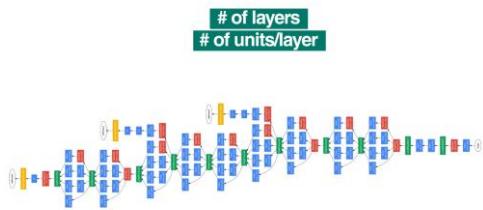
Jonathon Shlens, Yair Daon & Aaron Courville
Deep Learning
MIT Press - 2016

Hyperparameters

model-specific vs **optimizer-specific**

- | | |
|--|--|
| architecture
activations
initializations
loss functions
optimizers
regularizers | learning rate
batch size
of epochs |
|--|--|

Architecture



Christian Szegedy, Wei Liu, Yangqing Jia, et al.
Going Deeper with Convolutions (GoogleNet/Inception)
CVPR - 2015

Activations



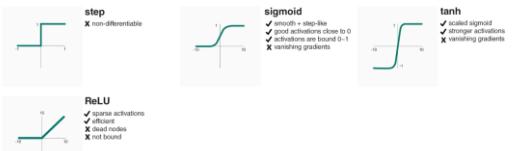
Activations



Activations



Activations



Hinrichs, Bar Shalom & Geoffrey E. Hinton
ImageNet Classification with Deep Convolutional Neural Networks
Advances in Neural Information Processing - NIPS 2012

Activations



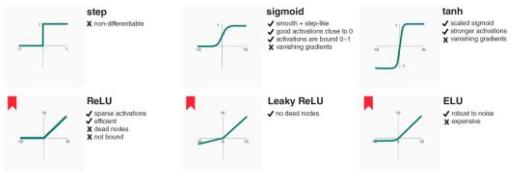
Kaiming He, Xiangyu Zhang, Shengyong Chen & Jian Sun
Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification
International Conference on Computer Vision - ICCV 2015

Activations



Dishant Chintala, Thomas Unterthiner & Sepp Hochreiter
Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs)
International Conference on Computer Vision - ICCV 2015

Activations

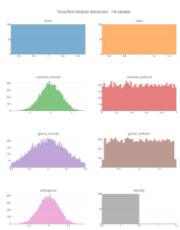


Initializations

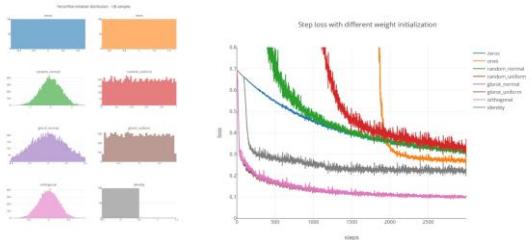
0 - stuck at a saddle point
constants - difficult to break the symmetry
large random values - small gradients, slow convergence



Initializations



Initializations



Initializations

Name	α	β	γ	Reference
Constant	$\alpha = 0$	$\beta = 0$	$\gamma \geq 0$	used by [ZF14]
Xavier/Glorot uniform	$\alpha = 0$	$\beta = \sqrt{\frac{6}{m_1 + m_2}}$	$\gamma = 0$	[GB10]
Xavier/Glorot normal	$\alpha = 0$	$\beta = \left(\frac{2}{m_1 + m_2}\right)^2$	$\gamma = 0$	[GB10]
He	$\alpha = 0$	$\beta = \frac{2}{m_1}$	$\gamma = 0$	[HZRS15a]
Orthogonal	—	—	$\gamma = 0$	[SMM13]
LSUV	—	—	$\gamma = 0$	[MM15]

Table B.2: Weight Initialization schemes of the form $w \sim \alpha \cdot \mathcal{U}[-1, 1] + \beta \cdot \mathcal{N}(0, 1) + \gamma$. m_1, m_2 are the number of units in the previous layer and the next layer. Typically, this is initialized with constant 0 and weight by one of the other schemes to prevent unit-exadaptation. However, dropout makes it possible to use constant initialization for all parameters. LSUV and Orthogonal initialization cannot be described with this simple pattern.

Mario [https://arxiv.org/pdf/1707.09725.pdf](#)

Loss Functions

regression - mean squared error

multiclass classification - categorical cross entropy

pixel classification - dice/Wasserstein dice coefficient

Optimizers**stochastic gradient descent + momentum**

Optimizers**stochastic gradient descent + momentum****adaptive gradient (AdaGrad)**

John Duchi, Elad Hazan & Yair Singer

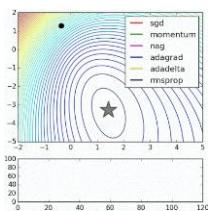
Adaptive Subgradient Methods for Online Learning and Stochastic Optimization

Journal of Machine Learning Research - 2011

Optimizers**stochastic gradient descent + momentum****adaptive gradient (AdaGrad)****root mean square propagation (RMSProp)**

*Geffrey Hinton*Coursera: Neural Networks for Machine Learning - Lecture 6
http://www.cs.toronto.edu/~tijmen/cs321/slides/lecture_slides_lec6.pdf

Optimizers



<https://imgur.com/gallery/qslp>

Regularizers

L1, L2 regularization

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2}|W|$$

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2} W^2$$

Regularizers

L1, L2 regularization

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2}|W|$$

$$\mathcal{L}_{new} = \mathcal{L} + \frac{\lambda}{2} W^2$$



<https://doi.org/10.1111/jcpp.12501> | Article first published online by Wiley Online Library: *Journal of Child Psychology and Psychiatry* 56(10), 1033–1042 © 2015 The Authors. *Journal of Child Psychology and Psychiatry* © 2015 Association for Child and Adolescent Mental Health.

Regularizers

dropout

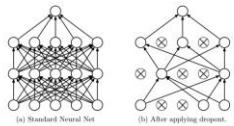


Figure 1: Dropout Neural Net Model. **Left:** A standard neural net with 2 hidden layers. **Right:** An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

Nicolas Sermatier, Geoffrey Hinton, Alex Krizhevsky, Ruslan Salakhutdinov

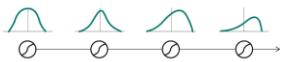
Dropout: A Simple Way to Prevent Neural Networks from Overfitting

Journal of Machine Learning Research - 2014



Regularizers

batch normalization



Jerome Boite & Christian Szegedy

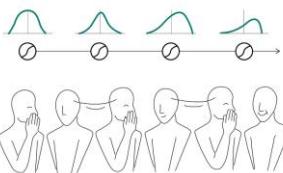
Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift

Journal of Machine Learning Research - 2014



Regularizers

batch normalization



Jerome Boite & Christian Szegedy

Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift

Journal of Machine Learning Research - 2014

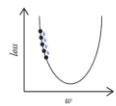


Optimizer-specific Hyperparameters

learning rate
0.1, 0.01, 0.001, 0.0001, ...

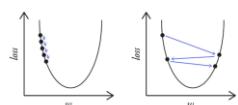
Optimizer-specific Hyperparameters

learning rate
0.1, 0.01, 0.001, 0.0001, ...



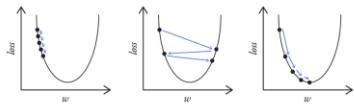
Optimizer-specific Hyperparameters

learning rate
0.1, 0.01, 0.001, 0.0001, ...



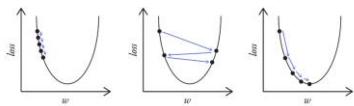
Optimizer-specific Hyperparameters

learning rate
0.1, 0.01, 0.001, 0.0001, ...



Optimizer-specific Hyperparameters

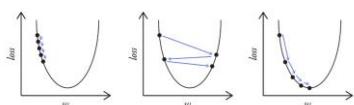
learning rate
0.1, 0.01, 0.001, 0.0001, ...



batch size

Optimizer-specific Hyperparameters

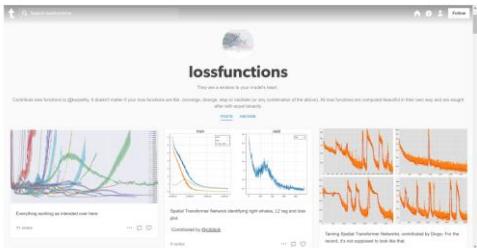
learning rate
0.1, 0.01, 0.001, 0.0001, ...



batch size
16, 32, 64, 128,...

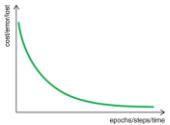
of epochs
early stopping

Babysitting your Network

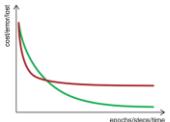


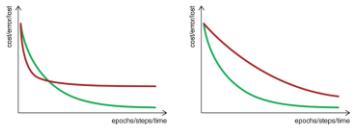
<https://lossfunctions.tumblr.com/>

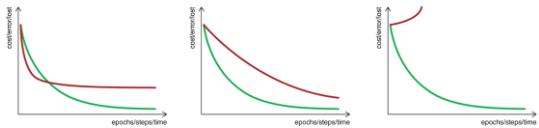
Debugging through Learning Curves

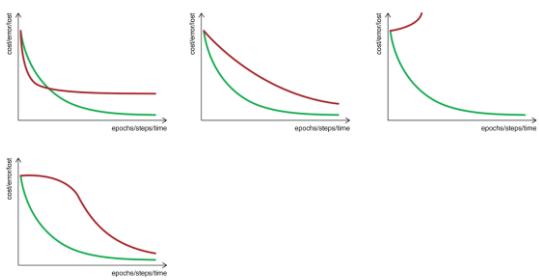


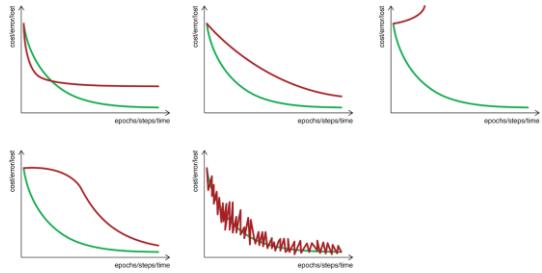
Debugging through Learning Curves



Debugging through Learning Curves

Debugging through Learning Curves

Debugging through Learning Curves

Debugging through Learning Curves**Debugging through Learning Curves**