

Deep Learning

Loo Tiang Kuan, Leonard

A0098795B

supervised by

Associate Professor Ji Hui

co-supervised by

Assistant Professor Alvina Goh Siew Wee

May 5, 2017

Today's Agenda

- 1 Introduction to Deep Learning
- 2 Neural Networks
- 3 Modifications
- 4 Regularization

Outline

1 Introduction to Deep Learning

- Recent Developments
- Overview

2 Neural Networks

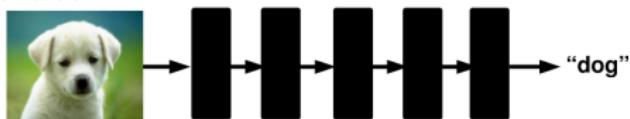
3 Modifications

4 Regularization

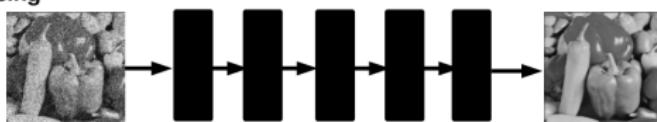
Recent developments

Supervised Deep Learning

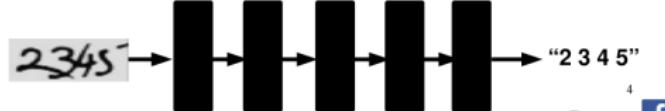
Classification



Denoising



OCR



4

Ranzato

CNN Photo

placed by $\|\mathbf{w} - \hat{\mathbf{w}}\|_1$, and resulting the following function:

placed by $\|\mathbf{w} - \hat{\mathbf{w}}\|_1$, and resulting the following function:

$$\min_{\mathbf{w}} \|\mathbf{w} - \hat{\mathbf{w}}\|_1 + C \sum_{i=1}^N f(\max(0, 1 - \mathbf{y}_i \mathbf{w}))$$

present a simple method to deal with the artifacts. we show that the proposed algorithm can also effectively process natural blurred images, and low illumination images which are not handled by the adaptive SVM framework [16] so that the disparity between \mathbf{w} and $\hat{\mathbf{w}}$ can be constrained while minimizing the classification error over \mathcal{D} .

Specifically, the regularizer $\|\mathbf{w}\|_1$ in standard ℓ_1 -regularized linear SVM [16] is replaced by $\|\mathbf{w} - \hat{\mathbf{w}}\|_1$, and resulting the following objective function:

the adaptive SVM framework [16] so that the disparity between \mathbf{w} and $\hat{\mathbf{w}}$ can be constrained while minimizing the classification error over \mathcal{D} . Specifically, the regularizer $\|\mathbf{w}\|_1$ in standard ℓ_1 -regularized linear SVM [16] is replaced by $\|\mathbf{w} - \hat{\mathbf{w}}\|_1$, and resulting the following objective function:

Outline

1 Introduction to Deep Learning

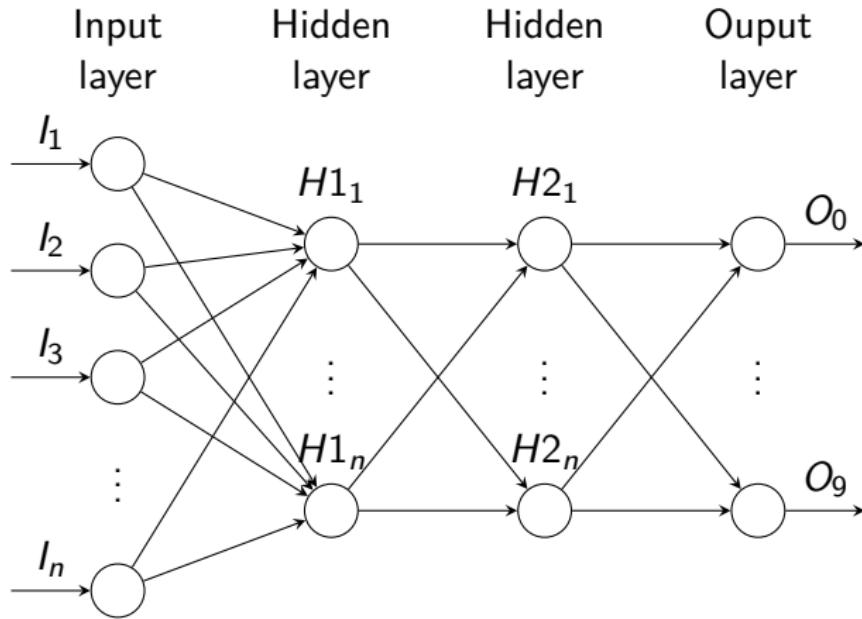
- Recent Developments
- Overview

2 Neural Networks

3 Modifications

4 Regularization

Overview



MNIST Data



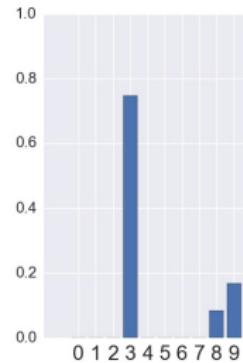
MNIST Handwritten Database

50000 Training Data

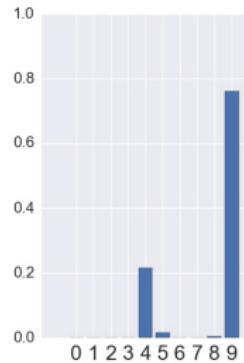
10000 Validation Data

10000 Test Data

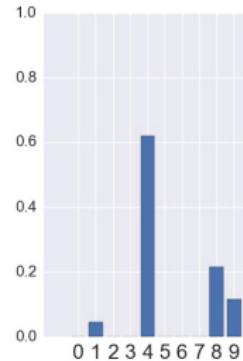
Overview



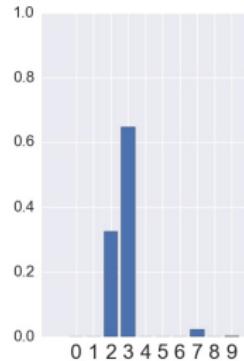
True Label: 3



True Label: 4



True Label: 9



True Label: 2



Overview

Architecture	Activation	Cost	Algorithm
Dense	Sigmoid	Quadratic	Gradient Descent (GD)
Convolutional	Softmax	Cross Entropy	Stochastic GD
	Tanh	Log-likelihood	Backpropagation
	ReLU		

Improvements

Regularization

Dropout

Initialisation

Overview

Session 1

Activations

Cost

Stochastic Gradient Descent

Backpropagation

Overfitting

Regularization

Session 2

Initialisation

Hyper-parameters

Variants of SGD

More Regularisers

Intro to ConvNets

Session 3

Popular ConvNet Architectures

AlexNet

ZFNet

VGG

GoogLeNet

ResNet

Denoising

Session 4/5

Autoencoders

Transposed Convolution

Generative Network

Generative Adversarial Network

Outline

1 Introduction to Deep Learning

2 Neural Networks

- Basics
 - Activation
 - Cost
- Learning

3 Modifications

4 Regularization

Outline

1 Introduction to Deep Learning

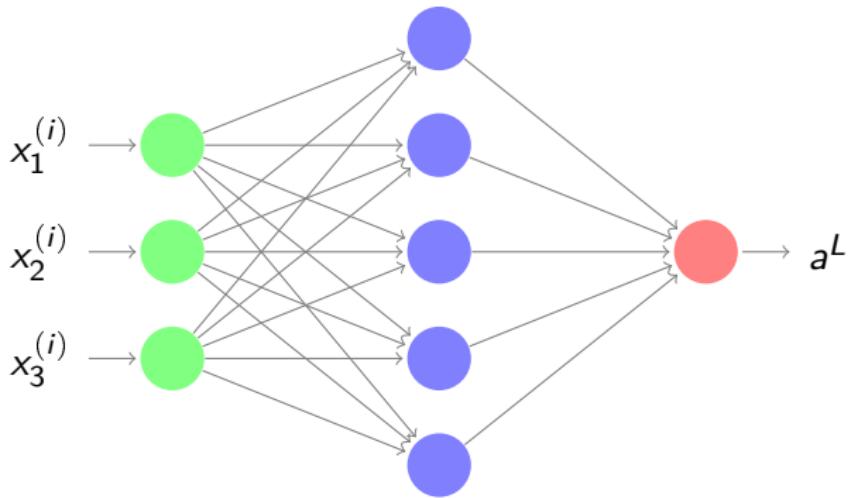
2 Neural Networks

- Basics
 - Activation
 - Cost
- Learning

3 Modifications

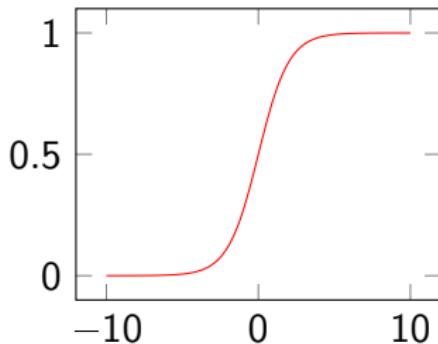
4 Regularization

Perceptron



$$a^I = \begin{cases} 1, & W^I a^{I-1} + b^I \geq 0 \\ 0, & W^I a^{I-1} + b^I < 0 \end{cases}$$

Sigmoid



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

$$a^I = \sigma(z^I) \approx \begin{cases} 1, & z \rightarrow \infty \\ 0, & z \rightarrow -\infty \end{cases}$$

Outline

1 Introduction to Deep Learning

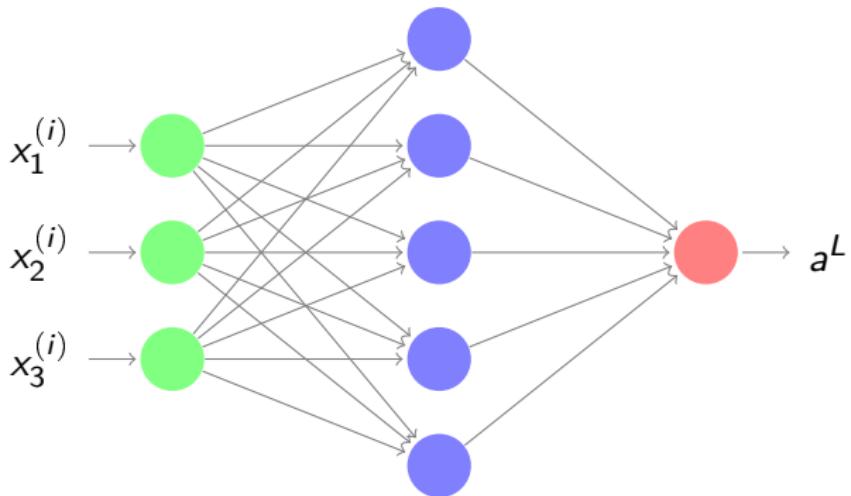
2 Neural Networks

- Basics
 - Activation
 - Cost
- Learning

3 Modifications

4 Regularization

Quadratic Cost



$$C(\mathbf{x}; \Theta) = C(\mathbf{x}, \mathbf{w}, \mathbf{b}) = \frac{1}{2n} \sum_{i=1}^n \|y(\mathbf{x}^{(i)}) - a^L(\mathbf{x}^{(i)})\|_2^2,$$

where $\mathbf{x} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}\}$

Outline

1 Introduction to Deep Learning

2 Neural Networks

- Basics

- Learning

- Gradient Descent
- Stochastic Gradient Descent
- Backpropagation

3 Modifications

4 Regularization

Outline

1 Introduction to Deep Learning

2 Neural Networks

- Basics
- Learning
 - Gradient Descent
 - Stochastic Gradient Descent
 - Backpropagation

3 Modifications

4 Regularization

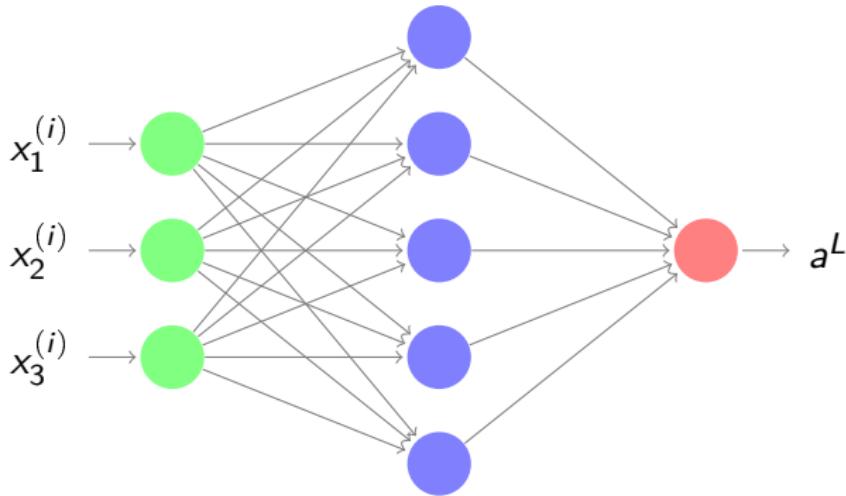
Gradient descent

For each $\mathbf{x}^{(i)}$,

$$\begin{aligned} C(\mathbf{x}^{(i)}, w, b) &= \frac{1}{2} \left[y(\mathbf{x}^{(i)}) - a^2(\mathbf{x}^{(i)}) \right]^2 \\ \frac{\partial C}{\partial w^{(2)}} &= \left(\frac{\partial C}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial z^{(2)}} \right) \frac{\partial z^{(2)}}{\partial w^{(2)}} \\ \frac{\partial C}{\partial b^{(2)}} &= \left(\frac{\partial C}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial z^{(2)}} \right) \frac{\partial z^{(2)}}{\partial b^{(2)}} \\ \frac{\partial C}{\partial w^{(1)}} &= \left(\frac{\partial C}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial z^{(2)}} \right) \left(\frac{\partial z^{(2)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial z^{(1)}} \right) \frac{\partial z^{(1)}}{\partial w^{(1)}} \\ \frac{\partial C}{\partial b^{(1)}} &= \left(\frac{\partial C}{\partial a^{(2)}} \frac{\partial a^{(2)}}{\partial z^{(2)}} \right) \left(\frac{\partial z^{(2)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial z^{(1)}} \right) \frac{\partial z^{(1)}}{\partial b^{(1)}} \end{aligned} \tag{1}$$

Backpropagation (3).

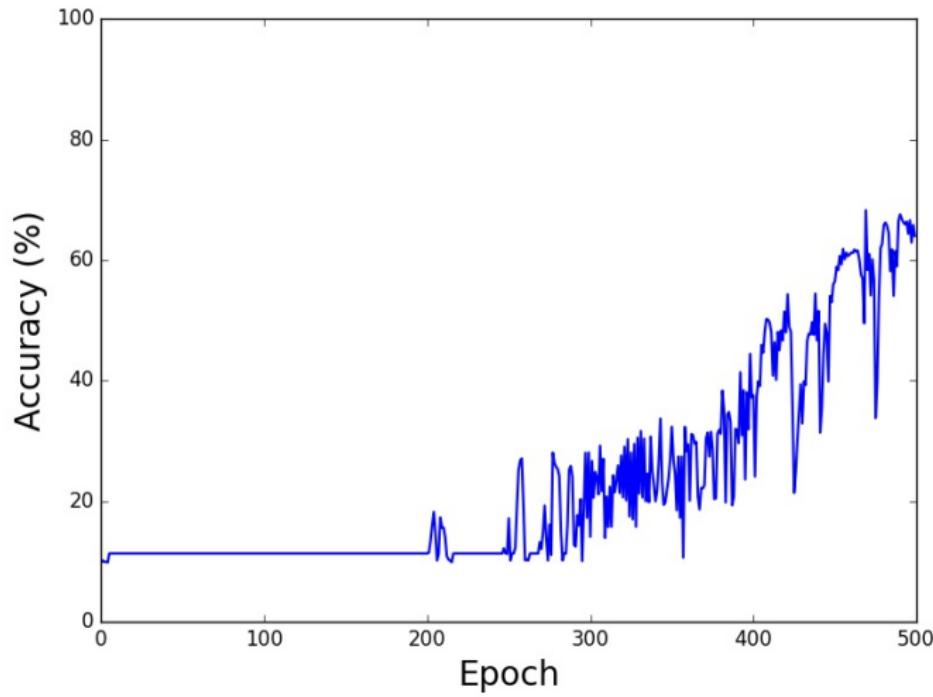
Gradient descent



$$C = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} |y - a^L|^2 \quad w^{(k+1)} = w^{(k)} - \eta \nabla_{w^{(k)}} C(\mathbf{x})$$

$$\nabla C(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \nabla C(\mathbf{x}^{(i)}) \quad b^{(k+1)} = b^{(k)} - \eta \nabla_{b^{(k)}} C(\mathbf{x})$$

Gradient Descent



Slow Learning!

Outline

1 Introduction to Deep Learning

2 Neural Networks

- Basics
- Learning
 - Gradient Descent
 - Stochastic Gradient Descent
 - Backpropagation

3 Modifications

4 Regularization

Stochastic Gradient Descent

Randomly sample into m mini-batches of size k :

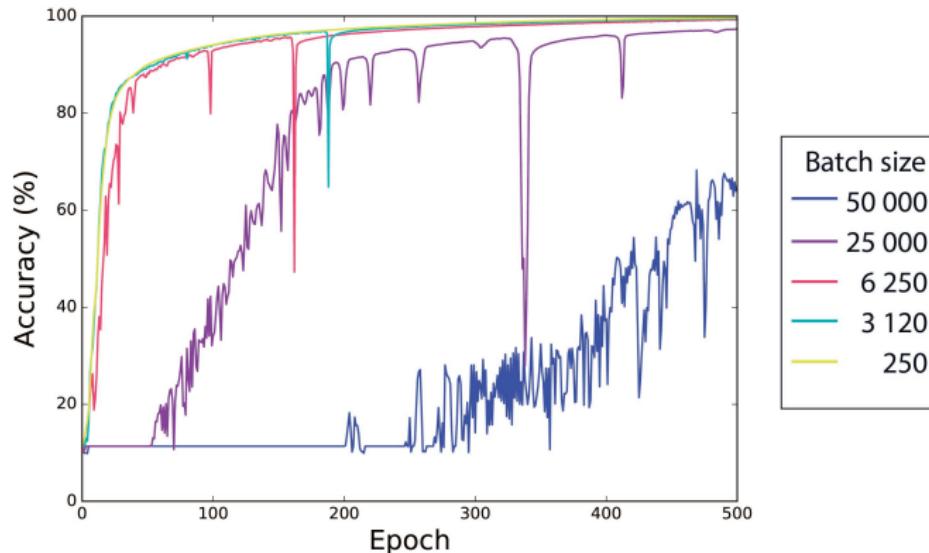
$$\begin{aligned}\mathbf{x} &= \left\{ \{x^{(1)}, x^{(2)}, \dots, x^{(k)}\}, \{x^{(k+1)}, x^{(k+2)}, \dots, x^{(2k)}\}, \dots, \{x^{(n-k+1)}, \dots, x^{(n)}\} \right\} \\ &= \{M_1, M_2, \dots, M_m\}\end{aligned}$$

$$\nabla C(\mathbf{x}) \approx \frac{1}{m} \sum_{i=1}^m \nabla C(\mathbf{x}_{M_j}), \text{ for each mini-batch } M_j$$

Parameters update for each mini-batch M_i :

$$\begin{aligned}w &\leftarrow w - \frac{\eta}{m} \frac{\partial C}{\partial w} \\ b &\leftarrow b - \frac{\eta}{m} \frac{\partial C}{\partial b}\end{aligned}\tag{2}$$

Stochastic Gradient Descent



1 Epoch - All $\{M_1, M_2, \dots\}$ are used.

Learning accelerates with more minibatches (smaller minibatch size).

Stochastic Gradient Descent

GPU ACCELERATION Training A Deep, Convolutional Neural Network			
Batch Size	Training Time CPU	Training Time GPU	GPU Speed Up
64 images	64 s	7.5 s	8.5X
128 images	124 s	14.5 s	8.5X
256 images	257 s	28.5 s	9.0X

<http://www.nvidia.com>

Deep Learning For Image Classification (Slide 28)

Remark:

- $w^{(k+1)} = w^{(k)} - \frac{\eta}{m} \frac{\partial C}{\partial w^{(k)}},$
need to adjust the learning rate η accordingly.
- Batch-size vs GPU computation

Outline

1 Introduction to Deep Learning

2 Neural Networks

- Basics
- Learning
 - Gradient Descent
 - Stochastic Gradient Descent
 - Backpropagation

3 Modifications

4 Regularization

Backpropagation

Recall the gradient descent in Eqn (1).

Definition (Error term δ_j^l)

$$\delta_j^l = \frac{\partial C}{\partial z_j^l}.$$

We have the four following Backpropagation equations

$$\begin{aligned}\delta_j^l &= \frac{\partial C}{\partial a_j^l} g'(z_j^l). \\ \delta^l &= \left([W^{l+1}]^T \delta^{l+1} \right) \odot g'(z^l). \\ \frac{\partial C}{\partial b_i^l} &= \delta_i^l. \\ \frac{\partial C}{\partial w_{ij}^l} &= a_j^{l-1} \delta_i^l.\end{aligned}\tag{3}$$

Outline

1 Introduction to Deep Learning

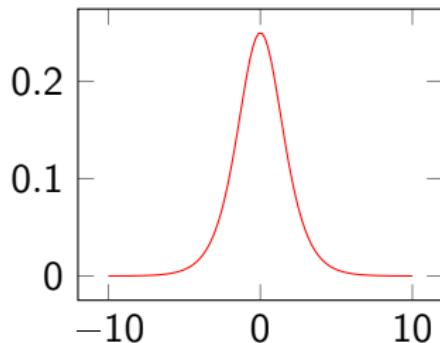
2 Neural Networks

3 Modifications

- Saturation : Sigmoid and Quadratic Cost
- Other Cost Functions
- Other Activations

4 Regularization

Saturation



$$C(\mathbf{x}, \mathbf{w}, \mathbf{b}) = \frac{1}{2n} \left\| \mathbf{y} - \mathbf{a}^L \right\|_2^2$$
$$\frac{\partial C}{\partial w^L} = a^{L-1} (a^L - y) \sigma'(z) \rightarrow 0$$
$$\frac{\partial C}{\partial b^L} = (a^L - y) \sigma'(z) \rightarrow 0$$

Outline

1 Introduction to Deep Learning

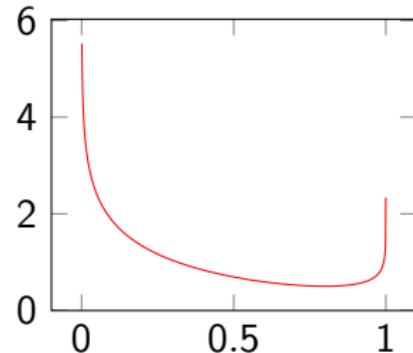
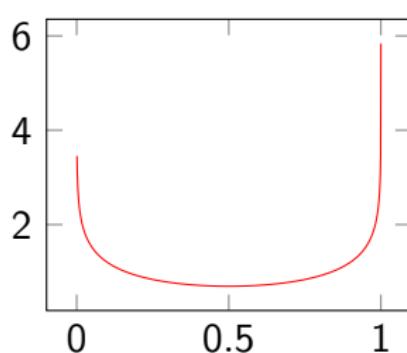
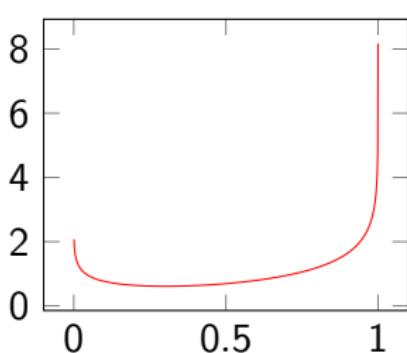
2 Neural Networks

3 Modifications

- Saturation : Sigmoid and Quadratic Cost
- Other Cost Functions
- Other Activations

4 Regularization

Cross Entropy Cost



$$C(\mathbf{x}, \mathbf{w}, \mathbf{b}) = - \left[\left(y \log a^L \right) + (1 - y) \log (1 - a^L) \right]$$

$$\nabla C(\mathbf{x}, \mathbf{w}, \mathbf{b}) = - \left[\frac{y}{a^L} - \frac{1 - y}{1 - a^L} \right]$$

$$\frac{\partial C}{\partial \mathbf{w}^L} = a^{L-1} (a^L - y)$$

$$\frac{\partial C}{\partial \mathbf{b}^L} = a^L - y$$

Other Activations and Cost Functions

Softmax function:

$$\text{softmax}(z^I) = \frac{e^{z_j^I}}{\sum_k e^{z_k^I}}$$

Log-likelihood Cost:

$$C = -\frac{1}{n} \sum_x \sum_k y_k \log a_k^L(x), \text{ for } k \text{ classes}$$

$$\begin{aligned}\frac{\partial C}{\partial w^L} &= a^{L-1}(a^L - y) \\ \frac{\partial C}{\partial b^L} &= a^L - y\end{aligned}$$

Outline

1 Introduction to Deep Learning

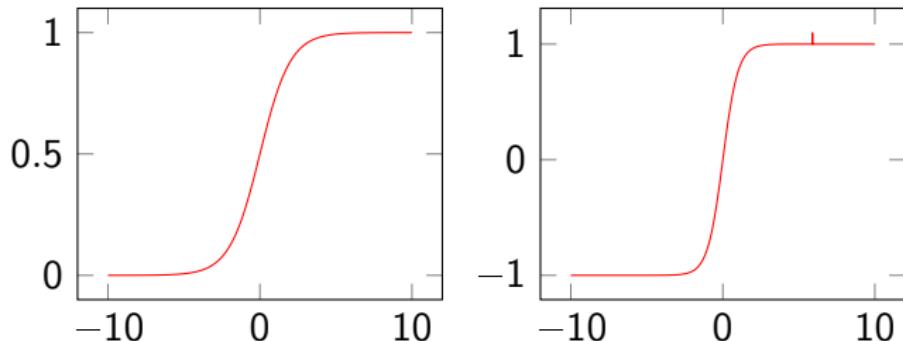
2 Neural Networks

3 Modifications

- Saturation : Sigmoid and Quadratic Cost
- Other Cost Functions
- Other Activations

4 Regularization

tanh



$$\sigma(z) = \frac{1+\tanh(\frac{z}{2})}{2}$$

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

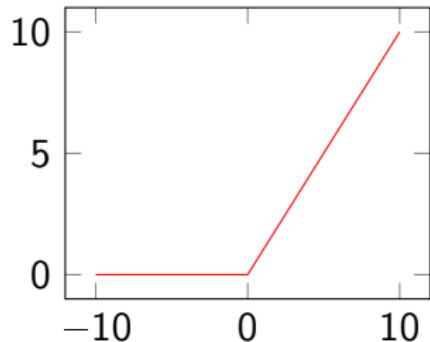
Pros

$\tanh(z) \in (-1, 1)$
weights can increase/decrease

Cons

Possibly saturate greater

ReLU



$$\text{ReLU}(z) = \max\{0, z\}$$

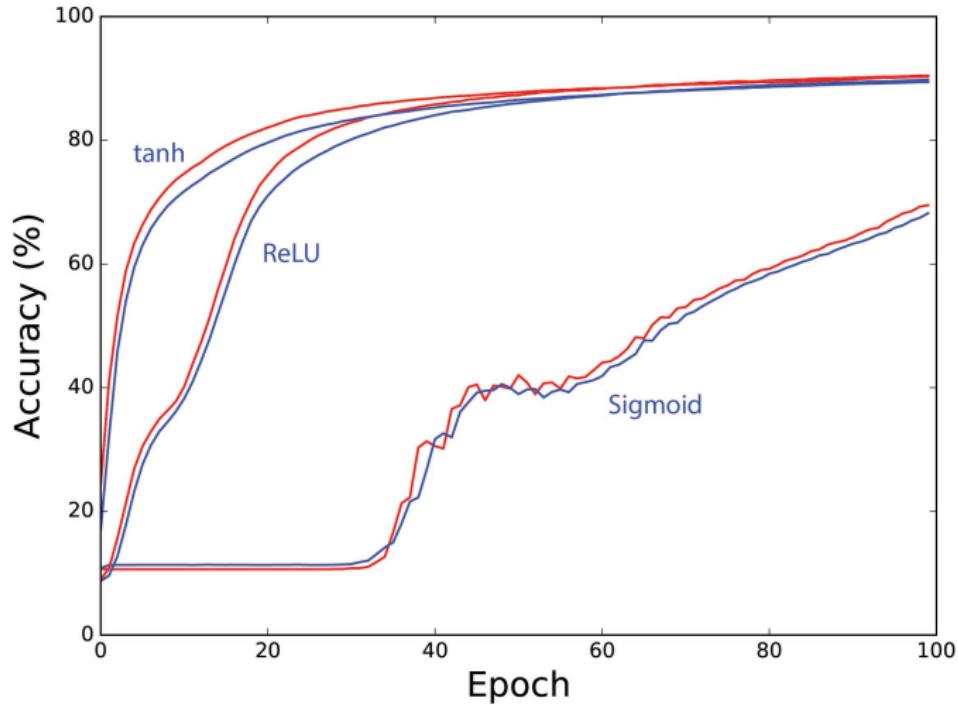
Pros

- Does not saturate at large z
- Perform well in image processing
- More computationally efficient

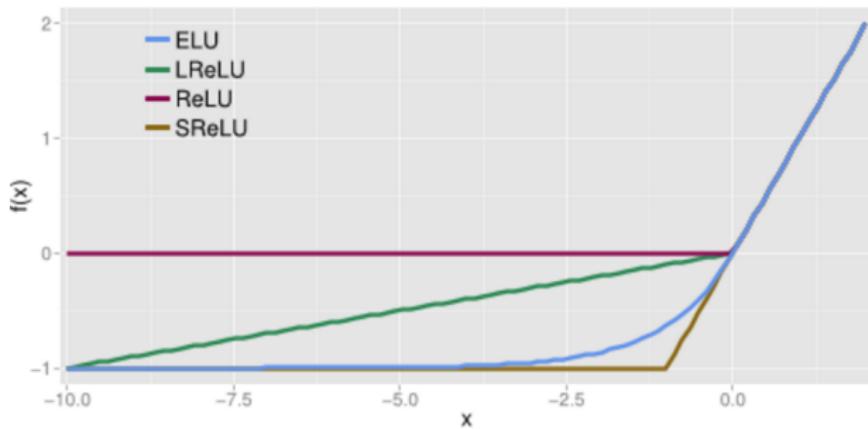
Cons

- gradient and value is 0 at negative z

Comparison



Even More Activations



$$\text{LReLU}(z) = \max(\alpha z, z)$$

$$\text{ELU}(z) = \max(\alpha(e^z - 1), z)$$

Outline

1 Introduction to Deep Learning

2 Neural Networks

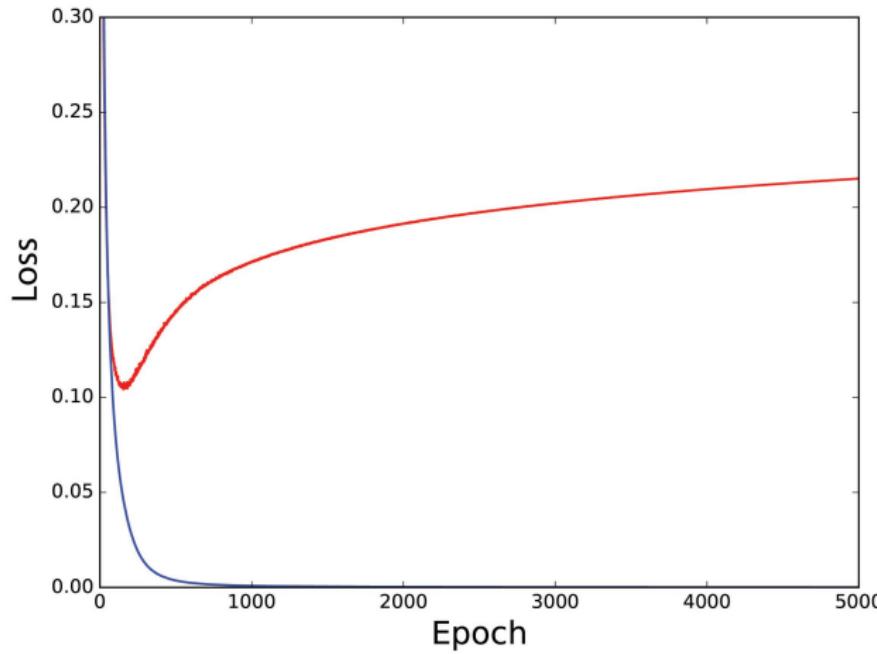
3 Modifications

4 Regularization

- Overfitting
- Regularization

Overfitting

- Too many parameters \implies Overfitting
- Regularize Weights



Outline

1 Introduction to Deep Learning

2 Neural Networks

3 Modifications

4 Regularization

- Overfitting
- Regularization
 - L2 Regularization
 - L1 Regularization
 - Early Stopping

Outline

1 Introduction to Deep Learning

2 Neural Networks

3 Modifications

4 Regularization

- Overfitting
- Regularization
 - L2 Regularization
 - L1 Regularization
 - Early Stopping

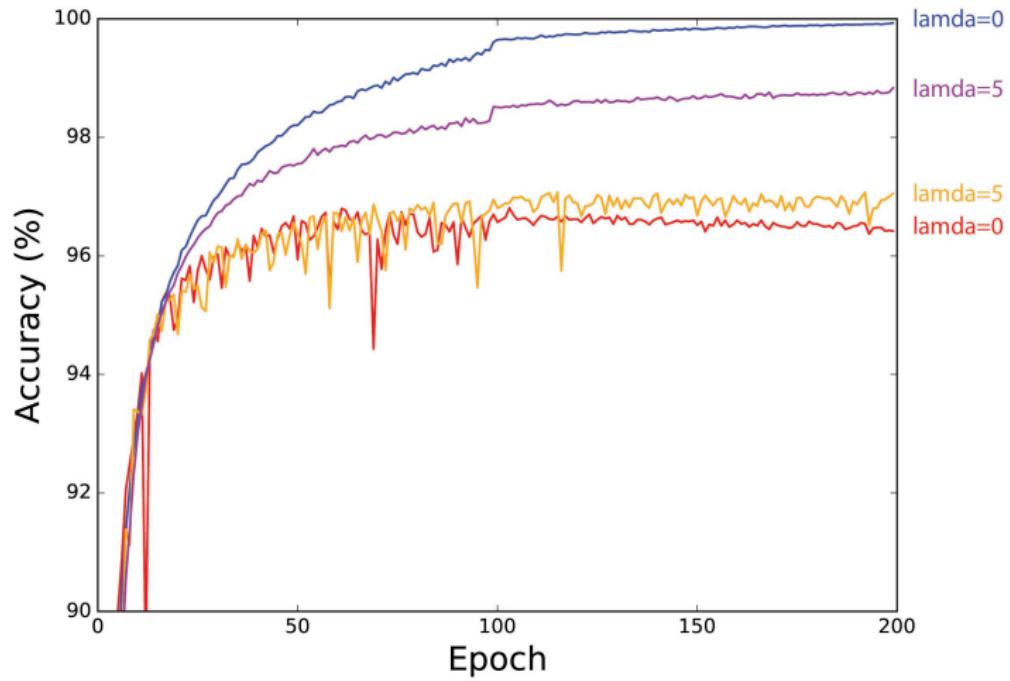
L2 Regularization

$$C = C_0 + \frac{1}{n} \frac{\lambda}{2} \|w\|_2^2 \quad (4)$$

$$\begin{aligned}\frac{\partial C}{\partial w} &= \frac{\partial C_0}{\partial w} + \frac{\lambda}{n} w \\ \frac{\partial C}{\partial b} &= \frac{\partial C_0}{\partial b}\end{aligned} \quad (5)$$

$$\begin{aligned}w_{k+1} &= \left(1 - \frac{\eta \lambda}{n}\right) w_k - \frac{\eta}{m} \frac{\partial C_0}{\partial w} \\ b_{k+1} &= b_k - \frac{\eta}{m} \frac{\partial C_0}{\partial b}\end{aligned} \quad (6)$$

L2 Regularization



Outline

1 Introduction to Deep Learning

2 Neural Networks

3 Modifications

4 Regularization

- Overfitting
- Regularization
 - L2 Regularization
 - L1 Regularization
 - Early Stopping

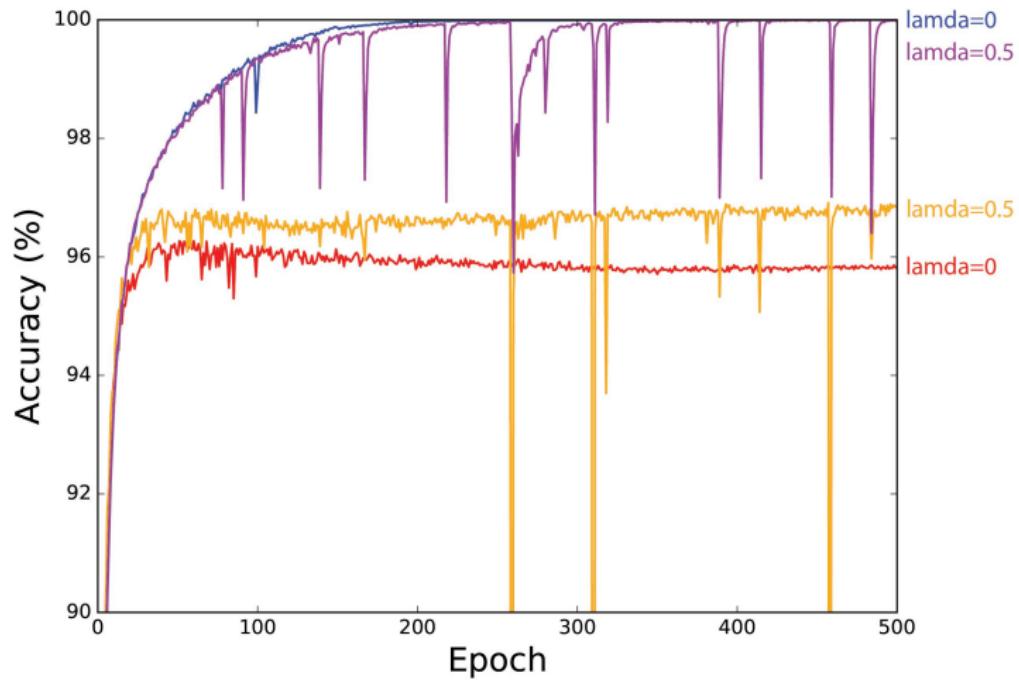
L1 Regularization

$$C = C_0 + \frac{\lambda}{n} ||w||_1 \quad (7)$$

$$\begin{aligned}\frac{\partial C}{\partial w} &= \frac{\partial C_0}{\partial w} + \frac{\lambda}{n} \text{sgn}(w) \\ \frac{\partial C}{\partial b} &= \frac{\partial C_0}{\partial b}\end{aligned} \quad (8)$$

$$\begin{aligned}w_{k+1} &= \left(w_k - \frac{\eta \lambda}{n} \text{sgn}(w_k) \right) - \frac{\eta}{m} \frac{\partial C_0}{\partial w} \\ b_{k+1} &= b_k - \frac{\eta}{m} \frac{\partial C_0}{\partial b}\end{aligned} \quad (9)$$

L1 Regularization



Outline

1 Introduction to Deep Learning

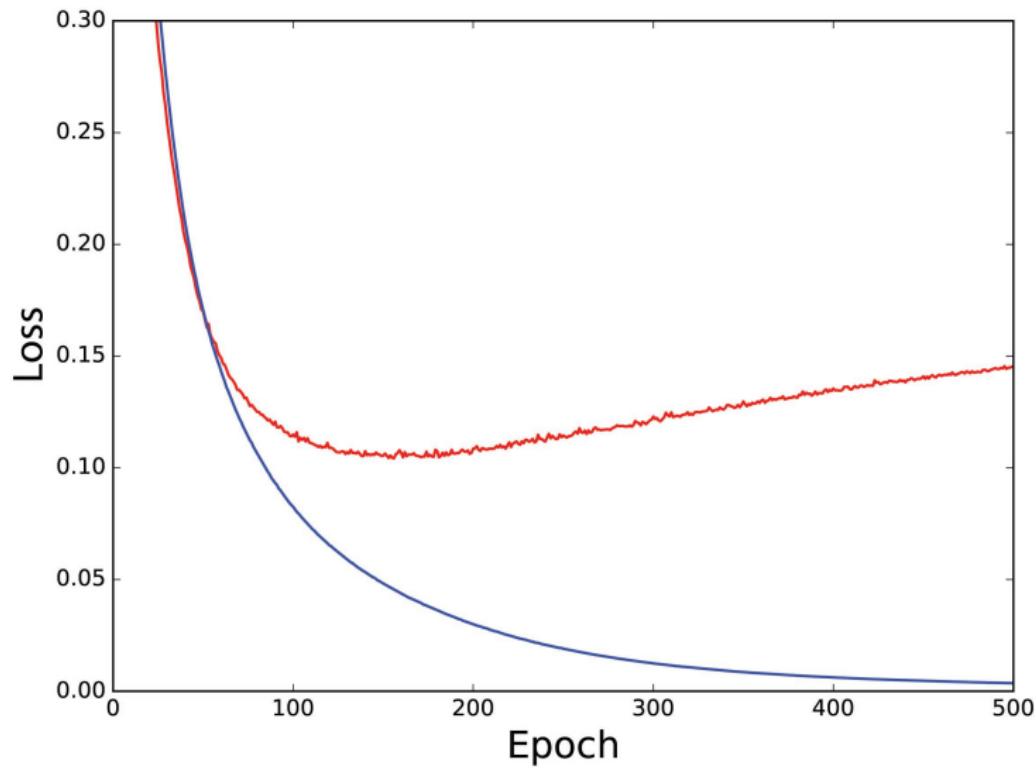
2 Neural Networks

3 Modifications

4 Regularization

- Overfitting
- Regularization
 - L2 Regularization
 - L1 Regularization
 - Early Stopping

Early Stopping



Thank You

References

-  www.cs.toronto.edu/~ranzato
-  neuralnetworksanddeeplearning.com/
-  Clevert, D. A., Unterthiner, T., & Hochreiter, S. (2015). *"Fast and accurate deep network learning by exponential linear units (ELUs)." arXiv preprint arXiv:1511.07289.*
-  Nair, V. and Hinton, G.E., 2010. *"Rectified linear units improve restricted boltzmann machines."* In Proceedings of the 27th International Conference on Machine Learning (ICML-10) (pp. 807-814).
-  Glorot, X., Bordes, A. and Bengio, Y., 2011, April. *"Deep Sparse Rectifier Neural Networks."* In Aistats (Vol. 15, No. 106, p. 275).

References

-  LeCun, Y.A., Bottou, L., Orr, G.B. and Müller, K.R., 2012. "*Efficient backprop.*" In Neural networks: Tricks of the trade (pp. 9-48). Springer Berlin Heidelberg.
-  Bengio, Y., 2009. "*Learning deep architectures for AI.*" Foundations and trends® in Machine Learning, 2(1), pp.1-127.