

Introduction to Neural Networks and Deep Learning

Deep Forward Neural Networks

Andres Mendez-Vazquez

February 14, 2021

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

For this initial analysis

We will look at the paper by Bengio

- "Learning deep architectures for AI", Foundations and trends in Machine Learning 2, 1 (2009), pp. 1--127.

- After Shanon pointed out the fact they are useful to represent complex problems [1].

For this initial analysis

We will look at the paper by Bengio

- "Learning deep architectures for AI", Foundations and trends in Machine Learning 2, 1 (2009), pp. 1--127.

And for this, we will look at Boolean functions

- After Shanon pointed out the fact they are useful to represent complex problems [1].

Architecture

A two-layer circuit of logic gates can represent any boolean function
[2]

- Any boolean function can be written as a sum of products, disjunctive normal form:
 - ▶ AND gates on the first layer with optional negation of inputs,
 - ▶ And OR gate on the second layer

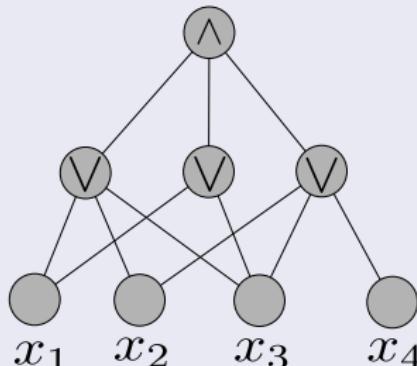
Example

Architecture

A two-layer circuit of logic gates can represent any boolean function [2]

- Any boolean function can be written as a sum of products, disjunctive normal form:
 - ▶ AND gates on the first layer with optional negation of inputs,
 - ▶ And OR gate on the second layer

Example



The Exponential Width

Here, we have a small problem

- There are functions computable with a polynomial-size logic gates circuit of depth k that require **exponential size** when restricted to depth $k - 1$ [3]
 - ▶ For Example

$$\text{parity} : (b_1, \dots, b_d) \in \{0, 1\}^d \mapsto \begin{cases} 1 & \text{if } \sum_{i=1}^d b_i \text{ is even} \\ -1 & \text{otherwise} \end{cases}$$

Block: The Input Block Learning

- Many of the results for boolean circuits can be generalized to architectures whose computational elements are linear threshold units

$$f(x) = 1_{wx+b>0}$$

- ▶ The fan-in of a circuit is the maximum number of inputs of a particular element.

The Exponential Width

Here, we have a small problem

- There are functions computable with a polynomial-size logic gates circuit of depth k that require **exponential size** when restricted to depth $k - 1$ [3]
 - ▶ For Example

$$\text{parity} : (b_1, \dots, b_d) \in \{0, 1\}^d \mapsto \begin{cases} 1 & \text{if } \sum_{i=1}^d b_i \text{ is even} \\ -1 & \text{otherwise} \end{cases}$$

How this impact Deep Learning?

- Many of the results for boolean circuits can be generalized to architectures whose computational elements are linear threshold units

$$f(x) = 1_{wx+b>0}$$

- ▶ The fan-in of a circuit is the maximum number of inputs of a particular element.

Therefore

How this impact shallow learning in Deep Learning?

- First, we define the concept of f_k function

- The function f_k is a function of N^{2k-2} variables. It is defined by a depth k circuit that is a tree. At the leaves of the tree there are unnegated variable, The i^{th} level from the bottom consists of \wedge -gates if i is even and otherwise it consists of \vee -gates.

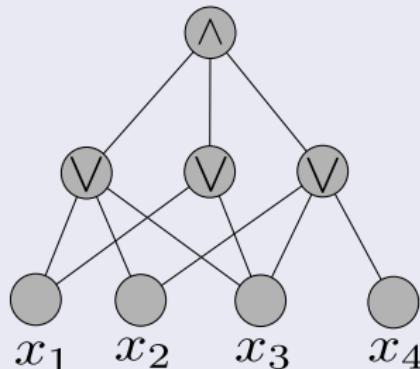
Therefore

How this impact shallow learning in Deep Learning?

- First, we define the concept of f_k function

Definition

- The function f_k is a function of N^{2k-2} variables. It is defined by a depth k circuit that is a tree. At the leaves of the tree there are unnegated variable, The i^{th} level from the bottom consists of \wedge -gates if i is even and otherwise it consists of \vee -gates.



An Important Theorem

Of particular interest is the following theorem

- Monotone weighted threshold circuits (i.e. multi-layer neural networks with linear threshold units and positive weights)

Theorem [1]

- A monotone weighted threshold circuit of depth $k - 1$ computing a function f_k has size at least 2^{cN} for some constant $c > 0$ and $N > N_0$.

An Important Theorem

Of particular interest is the following theorem

- Monotone weighted threshold circuits (i.e. multi-layer neural networks with linear threshold units and positive weights)

Theorem [4]

- A monotone weighted threshold circuit of depth $k - 1$ computing a function f_k has size at least 2^{cN} for some constant $c > 0$ and $N > N_0$.

Meaning

This theorem does not fail any type of architecture

- But the question arises, Are the depth 1, 2 and 3 architectures (many Machine Learning algorithms) too shallow to represent efficiently more complicated functions?

What happens in Deep Architectures?

- Bengio et al. argues that they can represent highly-varying functions

[5]

Meaning

This theorem does not fail any type of architecture

- But the question arises, Are the depth 1, 2 and 3 architectures (many Machine Learning algorithms) too shallow to represent efficiently more complicated functions?

What happens in Deep Architectures

- Bengio et al. argues that they can represent highly-varying functions [5]

Outline

1 Introduction

- Limitations of Shallow Architectures
- **Highly-varying functions**
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

Highly-varying functions

Meaning

- We say that a function is highly-varying when a piecewise approximation of that function would require a large number of pieces.

Example

- Deeper Architectures can handle such functions in a easier way than shallow ones.

Example

- The polynomial $\prod_{i=1}^n \sum_{j=1}^m a_{ij}x_j$ can be represented as a product of sums with only $O(nm)$ elements

Highly-varying functions

Meaning

- We say that a function is highly-varying when a piecewise approximation of that function would require a large number of pieces.

Clearly

- Deeper Architectures can handle such functions in a easier way than shallow ones.

Polynomial

- The polynomial $\prod_{i=1}^n \sum_{j=1}^m a_{ij}x_j$ can be represented as a product of sums with only $O(nm)$ elements

Highly-varying functions

Meaning

- We say that a function is highly-varying when a piecewise approximation of that function would require a large number of pieces.

Clearly

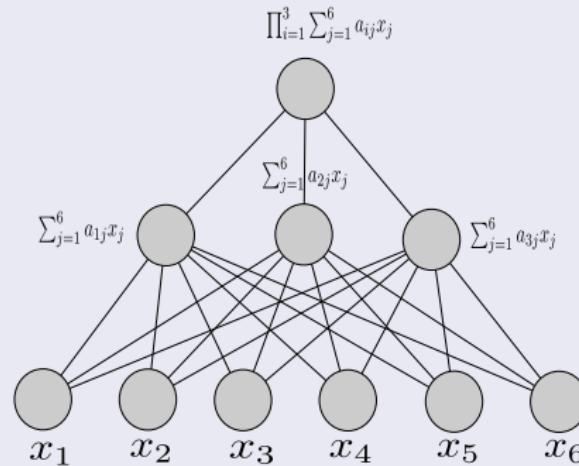
- Deeper Architectures can handle such functions in a easier way than shallow ones.

For Example

- The polynomial $\prod_{i=1}^n \sum_{j=1}^m a_{ij}x_j$ can be represented as a product of sums with only $O(nm)$ elements

Basically

We have a Perceptron Layer and a Product Second Layer

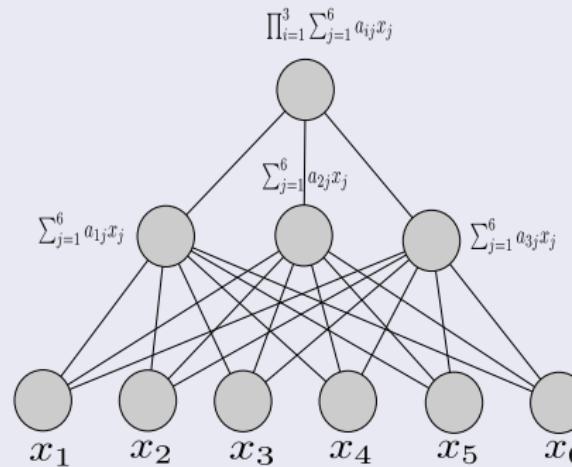


What will happen to the outputs?

- What will happen?

Basically

We have a Perceptron Layer and a Product Second Layer



What if I do a product of sums

- What will happen?

Ok, we have a problem

Because for our case

$$\prod_{i=1}^3 \sum_{j=1}^6 a_{ij} x_j = \sum_{j=1}^6 \prod_{i=1}^3 a_{ij} x_j$$

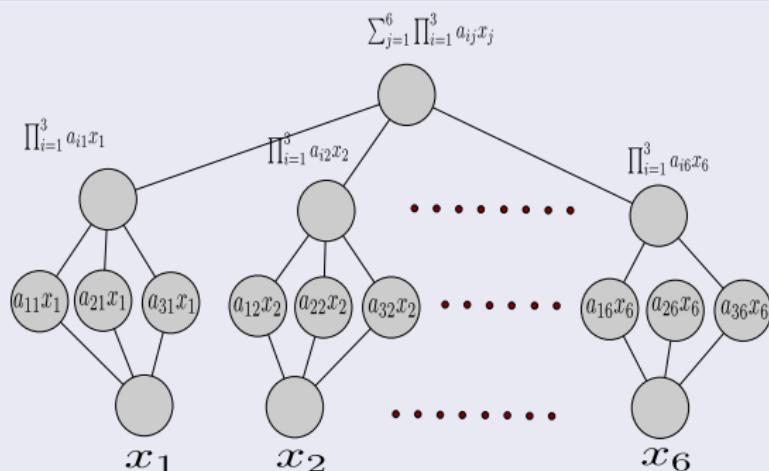
We have the following problem:

Ok, we have a problem

Because for our case

$$\prod_{i=1}^3 \sum_{j=1}^6 a_{ij} x_j = \sum_{j=1}^6 \prod_{i=1}^3 a_{ij} x_j$$

We have the following problem $O(n^m)$



Actually

You could claim

- Machine Learning shallow learning depends on complex computational units to handle complex functions

Deep learning

- Proposes simpler units but deeper structures to handle complex functions

What about both facts together?

- Complex adaptive units
- Deeper architectures to help such units
 - ▶ It seems to be the case of the human brain...!!!

Actually

You could claim

- Machine Learning shallow learning depends on complex computational units to handle complex functions

Deep Learning

- Proposes simpler units but deeper structures to handle complex functions

What about both kinds of filters?

- Complex adaptive units
- Deeper architectures to help such units
 - ▶ It seems to be the case of the human brain...!!!

Actually

You could claim

- Machine Learning shallow learning depends on complex computational units to handle complex functions

Deep Learning

- Proposes simpler units but deeper structures to handle complex functions

What about both ideas together

- Complex adaptive units
- Deeper architectures to help such units
 - ▶ It seems to be the case of the human brain...!!!

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

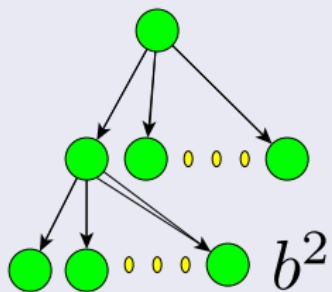
3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

Actually making some numbers

If we allow Fan Out problem

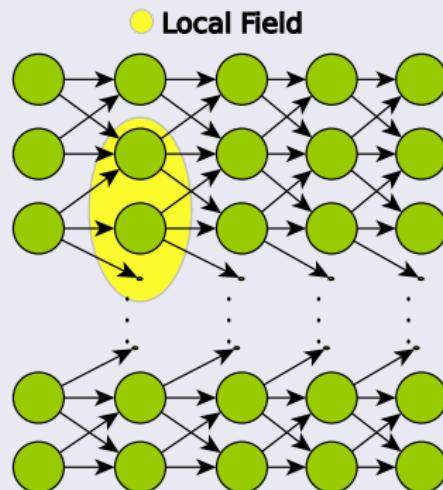
- Using a fully connected network with that allows a fanning out as



Yes, it allow more variability

This is what you try to do in shallow learning

- But other option avoid total fanout by using locality



Local vs Non-Local Generalization

We have smaller areas of information

- But we can go deeper

Round The Fair - What are Local Models?

- A local estimator partitions the input space in regions

Two Local-estimators-one based on matching local templates

- It can be thought of as having two levels...

The first level

- It is made of a set of templates which can be matched to the input.

Local vs Non-Local Generalization

We have smaller areas of information

- But we can go deeper

Avoid the Fan Out and Use Locality

- A local estimator partitions the input space in regions

↳ Local estimators are based on matching local templates

- It can be thought of as having two levels...

The first level

- It is made of a set of templates which can be matched to the input.

Local vs Non-Local Generalization

We have smaller areas of information

- But we can go deeper

Avoid the Fan Out and Use Locality

- A local estimator partitions the input space in regions

Thus, local estimators are based on matching local templates

- It can be thought of as having two levels...

The inner level

- It is made of a set of templates which can be matched to the input.

Local vs Non-Local Generalization

We have smaller areas of information

- But we can go deeper

Avoid the Fan Out and Use Locality

- A local estimator partitions the input space in regions

Thus, local estimators are based on matching local templates

- It can be thought of as having two levels...

The first level

- It is made of a set of templates which can be matched to the input.

Then

A template unit will output a value that indicates the degree of matching

$$K(x|\Theta)$$

The overall prediction is then:

- Typically a simple linear combination or product combination

$$L(x) = \sum_{i=1}^k K(x|\Theta_i)$$

For example, the kernel machine

$$f(x) = b + \sum_{i=1}^k \alpha_i K(x, x_i)$$

Then

A template unit will output a value that indicates the degree of matching

$$K(x|\Theta)$$

The second level combines these values

- Typically a simple linear combination or product combination

$$L(x) = \sum_{i=1}^k K(x|\Theta_i)$$

Classic example: the kernel machine

$$f(x) = b + \sum_{i=1}^k \alpha_i K(x, x_i)$$

Then

A template unit will output a value that indicates the degree of matching

$$K(x|\Theta)$$

The second level combines these values

- Typically a simple linear combination or product combination

$$L(x) = \sum_{i=1}^k K(x|\Theta_i)$$

Classic Example, the kernel machine

$$f(x) = b + \sum_{i=1}^k \alpha_i K(x, x_i)$$

As you can see

The Kernel has a local influence based on the support vectors

- For example the Gaussian Kernel

$$K(x, x_i) = \exp\left\{-\frac{\|x - x_i\|^2}{\sigma^2}\right\}$$

The Problem of Kernel

- The assumption that the target function is smooth or can be well approximated with a smooth function.

The Limitations of a fixed kernel Kernel like the Gaussian Kernel

- They have motivated a lot of research in designing kernels [6, 7]

As you can see

The Kernel has a local influence based on the support vectors

- For example the Gaussian Kernel

$$K(x, x_i) = \exp\left\{-\frac{\|x - x_i\|^2}{\sigma^2}\right\}$$

The Problem of Kernel

- The assumption that the target function is smooth or can be well approximated with a smooth function.

- They have motivated a lot of research in designing kernels [6, 7]

As you can see

The Kernel has a local influence based on the support vectors

- For example the Gaussian Kernel

$$K(x, x_i) = \exp\left\{-\frac{\|x - x_i\|^2}{\sigma^2}\right\}$$

The Problem of Kernel

- The assumption that the target function is smooth or can be well approximated with a smooth function.

The limitations of a fixed generic kernel such as the Gaussian kernel

- They have motivated a lot of research in designing kernels [6, 7]

For Example, in supervised learning

If we have the training example (x_i, y_i)

- We want to build predictor that output something near y_i when any other sample is near x_i

Especially when x_i is outlier

- Bengio and Le Cun claim this is not enough [8, 9]

Although it is possible to do

- That such highly varying space is due to a lack of the correct feature selection process.

For Example, in supervised learning

If we have the training example (x_i, y_i)

- We want to build predictor that output something near y_i when any other sample is near x_i

Basically the situation when regularizing

- Bengio and Le Cun claim this is not enough [8, 9]

Additional possible reason:

- That such highly varying space is due to a lack of the correct feature selection process.

For Example, in supervised learning

If we have the training example (x_i, y_i)

- We want to build predictor that output something near y_i when any other sample is near x_i

Basically the situation when regularizing

- Bengio and Le Cun claim this is not enough [8, 9]

Although, It is possible to argue

- That such highly varying space is due to a lack of the correct feature selection process.

However

If you look at the parity problem

$$\text{parity} : (b_1, \dots, b_d) \in \{0, 1\}^d \mapsto \begin{cases} 1 & \text{if } \sum_{i=1}^d b_i \text{ is even} \\ -1 & \text{otherwise} \end{cases}$$

Final note

- Let $f(\mathbf{x}) = b + \sum_{i=1}^{2^d} \alpha_i K(\mathbf{x}_i, \mathbf{x})$ be an affine combination of Gaussian with the same width σ centered on points $\mathbf{x}_i \in \{-1, 1\}^d$. If f solve the parity problem, then there are at least 2^{d-1} non-zero support vectors.

However

If you look at the parity problem

$$\text{parity} : (b_1, \dots, b_d) \in \{0, 1\}^d \mapsto \begin{cases} 1 & \text{if } \sum_{i=1}^d b_i \text{ is even} \\ -1 & \text{otherwise} \end{cases}$$

Theorem

- Let $f(\mathbf{x}) = b + \sum_{i=1}^{2^d} \alpha_i K(\mathbf{x}_i, \mathbf{x})$ be an affine combination of Gaussian with the same width σ centered on points $\mathbf{x}_i \in \{-1, 1\}^d$. If f solve the parity problem, then there are at least 2^{d-1} non-zero support vectors.

However

Although, this function is not a representative

- The kind of functions we are more interested in AI.

In suggest algorithm based systems

- They are not enough, but still not a conclusive result

What all

- More Memory could be added to those systems

However

Although, this function is not a representative

- The kind of functions we are more interested in AI.

It suggest that local based estimators

- They are not enough, but still not a conclusive result

- More Memory could be added to those systems

However

Although, this function is not a representative

- The kind of functions we are more interested in AI.

It suggest that local based estimators

- They are not enough, but still not a conclusive result

After all

- More Memory could be added to those systems

For example

Tensors have been used to add memory to SVM

$$\min_{U_i^{(m)}, \mathbf{K}^{(m)}, \boldsymbol{\beta}, b} \gamma \sum_{i=1}^N \left\| \mathcal{X}_i - \left[\mathbf{K}^{(1)} \mathbf{U}_i^{(1)}, \dots, \mathbf{K}^{(M)} \mathbf{U}_i^{(M)} \right] \right\|_F^2 + \dots$$
$$\lambda \boldsymbol{\beta}^T \widehat{\mathbf{K}} \boldsymbol{\beta} + \sum_{i=1}^N \left[1 - y_i \left(\widehat{\mathbf{k}}_i^T \boldsymbol{\beta} + b \right) \right]_+$$

- $\mathbf{K}^{(m)}$ are kernel matrices defined on each mode to capture the nonlinear part.
- $\mathbf{U}^{(m)} = [\mathbf{u}_1^{(m)}, \dots, \mathbf{u}_R^{(m)}]$ are factor matrices of size $I_m \times R_m$

However

A Problem

- You are limiting the Machine Learning operations to matrix additions and products and non-linear operations.
 - ▶ In a shallow way...

We need to add more complexity to our models

- After all deeper architectures construct complex functions layer by layer

However

A Problem

- You are limiting the Machine Learning operations to matrix additions and products and non-linear operations.
 - ▶ In a shallow way...

We need to add more complex functions

- After all deeper architectures construct complex functions layer by layer

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

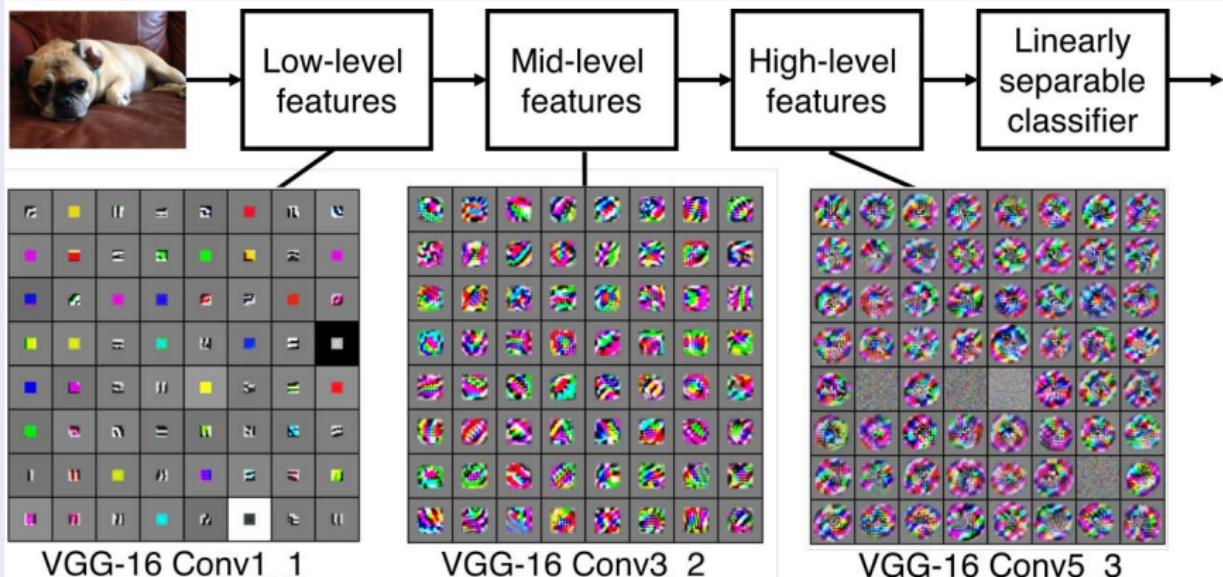
- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

By Using Weights in Certain Deep Learners

The Application of each Layer increase the complexity of the features



Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

● Introduction

- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

Some of the Models to be Reviewed

Convolutional Neural Networks

- The classic model that started the phenomena of Neural Networks.

RNN Encoder

- How to generate novel features by funneling.

Generative Models

- How to learn the underlying models

Some of the Models to be Reviewed

Convolutional Neural Networks

- The classic model that started the phenomena of Neural Networks.

Auto Encoder

- How to generate novel features by funneling.

Generative Models

- How to learn the underlying models

Some of the Models to be Reviewed

Convolutional Neural Networks

- The classic model that started the phenomena of Neural Networks.

Auto Encoder

- How to generate novel features by funneling.

Generative Models

- How to learn the underlying models

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
 - Auto Encoders
 - Generative Adversarial Networks
 - An Example, Boltzmann Machines
 - There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

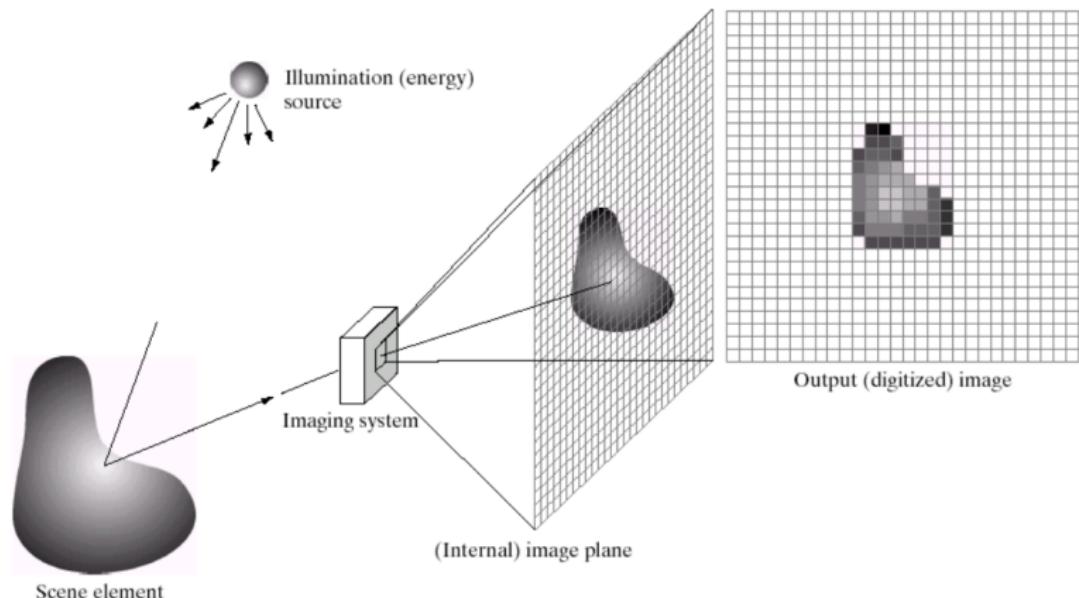
2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

Digital Images as pixels in a digitized matrix



Further

Pixel values typically represent

- Gray levels, colours, heights, opacities etc

Something Measurable

- Remember digitization implies that a digital image is an approximation of a real scene

Further

Pixel values typically represent

- Gray levels, colours, heights, opacities etc

Something Notable

- Remember digitization implies that a digital image is an approximation of a real scene

Therefore, we have the following process

Low Level Process

Input	Processes	Output
Image	Noise Removal	Improved Image
	Image Sharpening	

Example: Edge Detection

Therefore, we have the following process

Low Level Process

Input	Processes	Output
Image	Noise Removal	Improved Image
	Image Sharpening	

Example, Edge Detection



Then

Mid Level Process

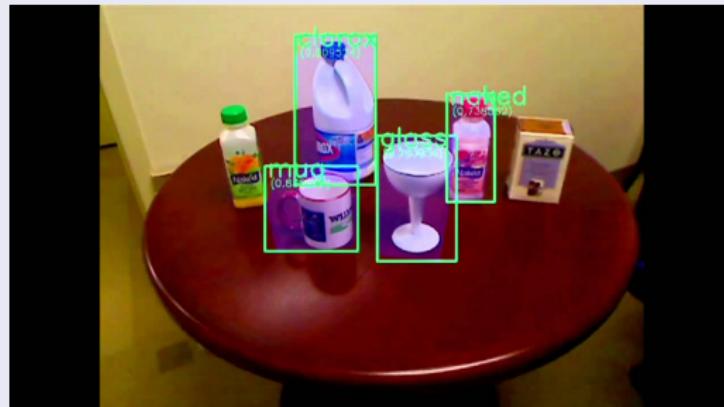
Input	Processes	Output
Image	Object Recognition	Attributes
	Segmentation	

Then

Mid Level Process

Input	Processes	Output
Image	Object Recognition	Attributes
	Segmentation	

Object Recognition



Therefore

It would be nice to automatize all these processes

- We would solve a lot of headaches when setting up such process

- By using a Neural Networks that replicates the process.

Therefore

It would be nice to automatize all these processes

- We would solve a lot of headaches when setting up such process

Why not to use the data sets

- By using a Neural Networks that replicates the process.

Convolutional Neural Networks History

Work by Hubel and Wiesel in the 1950s and 1960s

- They showed that cat and monkey visual cortices contain neurons that individually respond to small regions of the visual field.

- David H. Hubel and Torsten N. Wiesel (2005). Brain and visual perception: the story of a 25-year collaboration. Oxford University Press US, p. 106.

Convolutional Neural Networks History

Work by Hubel and Wiesel in the 1950s and 1960s

- They showed that cat and monkey visual cortices contain neurons that individually respond to small regions of the visual field.

After all more studies about the visual cortex happened

- David H. Hubel and Torsten N. Wiesel (2005). Brain and visual perception: the story of a 25-year collaboration. Oxford University Press US. p. 106.

Neurocognitron (Circa 1980)

Kunihiko Fukushima [10]

- Proposed a Hierarchical Network for image recognition with a convolution!!!

What is a Neurocognitron?

$$\varphi \left(\frac{1 + \sum_{k_{t-1}=1}^{K_{t-1}} \sum_{v \in S_t} a_t(k_{t-1}, v, k_t) u_{d-1}(k_{t-1}, n+v)}{1 + \frac{2}{1+r_t} b_t(k_t) v_{CI-1}(n)} - 1 \right)$$

What is a ReLU function?

$$\varphi(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$$

Neurocognitron (Circa 1980)

Kunihiko Fukushima [10]

- Proposed a Hierarchical Network for image recognition with a convolution!!!

But it used a function φ

$$\varphi \left(\frac{1 + \sum_{k_{t-1}=1}^{K_{t-1}} \sum_{v \in S_l} a_l(k_{t-1}, v, k_l) u_{cl-1}(k_{l=1}, n+v)}{1 + \frac{2r_l}{1+r_l} b_l(k_l) v_{Cl-1}(n)} - 1 \right)$$

What is ReLU function?

$$\varphi(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$$

Neurocognitron (Circa 1980)

Kunihiko Fukushima [10]

- Proposed a Hierarchical Network for image recognition with a convolution!!!

But it used a function φ

$$\varphi \left(\frac{1 + \sum_{k_{t-1}=1}^{K_{t-1}} \sum_{v \in S_l} a_l(k_{t-1}, v, k_l) u_{cl-1}(k_{l=1}, n+v)}{1 + \frac{2r_l}{1+r_l} b_l(k_l) v_{Cl-1}(n)} - 1 \right)$$

With a Relu function

$$\varphi(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases}$$

Furthermore (Circa 1993)

Weng et al. [11, 12]

- Proposed the use of Maxpooling to recognize 3D objects in 2D images

- The Beginning of the Dream!!!

Furthermore (Circa 1993)

Weng et al. [11, 12]

- Proposed the use of Maxpooling to recognize 3D objects in 2D images

Yan LeCunn finally proposed the use of backpropagation [13]

- The Beginning of the Dream!!!

Convolutional Neural Networks

Basically they are deep learners based in convolutions or its variants

$$(f * g)(i, j) = \sum_{k=-n}^{-n} \sum_{l=-n}^n f(k, l) \times g(i - k, j - l) \quad (1)$$

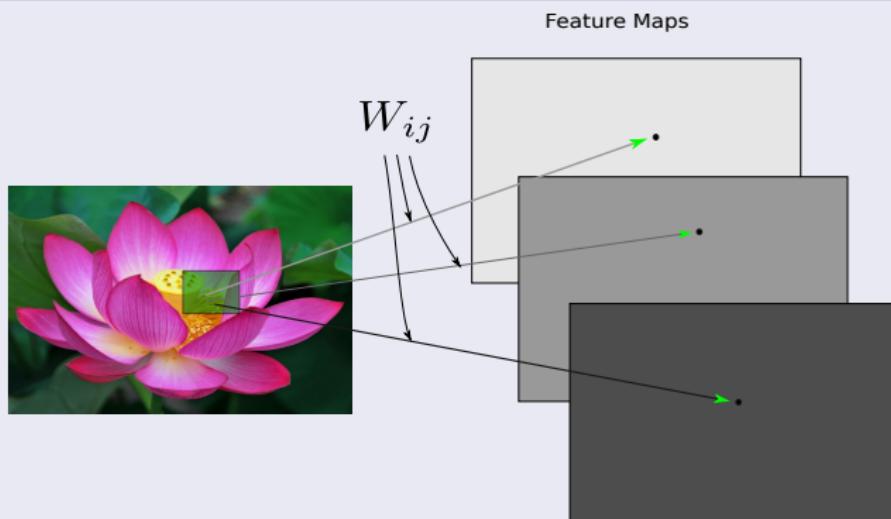
Kernel Filter

Convolutional Neural Networks

Basically they are deep learners based in convolutions or its variants

$$(f * g)(i, j) = \sum_{k=-n}^{+n} \sum_{l=-n}^{+n} f(k, l) \times g(i - k, j - l) \quad (1)$$

Basically Filters



Thus

Each Feature Map forms a 2D grid of features

That can be computed with a discrete convolution (*) of a kernel matrix k_{ij} which is the hidden weights matrix W_{ij} with rows and columns with its rows and columns flipped.

Thus

Each Feature Map forms a 2D grid of features

That can be computed with a discrete convolution (*) of a kernel matrix k_{ij} which is the hidden weights matrix W_{ij} with rows and columns with its rows and columns flipped.

In addition

- x_i is the i th channel of input.
- k_{ij} is the convolution kernel.
- y_j is the hidden layer output.

Thus

Each Feature Map forms a 2D grid of features

That can be computed with a discrete convolution (*) of a kernel matrix k_{ij} which is the hidden weights matrix W_{ij} with rows and columns with its rows and columns flipped.

In addition

- x_i is the i th channel of input.
- k_{ij} is the convolution kernel.

• y_j is the hidden layer output.

↳ Discrete convolution

$$y_j = \sum_i k_{ij} * x_i \quad (2)$$

Thus

Each Feature Map forms a 2D grid of features

That can be computed with a discrete convolution (*) of a kernel matrix k_{ij} which is the hidden weights matrix W_{ij} with rows and columns with its rows and columns flipped.

In addition

- x_i is the i th channel of input.
- k_{ij} is the convolution kernel.
- y_j is the hidden layer output.

↳ Discrete convolution

$$y_j = \sum_i k_{ij} * x_i \quad (2)$$

Thus

Each Feature Map forms a 2D grid of features

That can be computed with a discrete convolution (*) of a kernel matrix k_{ij} which is the hidden weights matrix W_{ij} with rows and columns with its rows and columns flipped.

In addition

- x_i is the i th channel of input.
- k_{ij} is the convolution kernel.
- y_j is the hidden layer output.

Thus the total output

$$y_j = \sum_i k_{ij} * x_i \quad (2)$$

Remark

We have that

- A Convolutional Neural Network (CNN) directly accepts raw images as input.

- Instead of assuming a certain comprehension of Computer Vision, one could think this is as a Silver Bullet.

- You still need to be aware of :

- ▶ The need of great quantities of data.
 - ▶ There initial works on how they work [14]

Remark

We have that

- A Convolutional Neural Network (CNN) directly accepts raw images as input.

Thus, their importance when training discrete filters

- Instead of assuming a certain comprehension of Computer Vision, one could think this is as a Silver Bullet.

- You still need to be aware of :

- ▶ The need of great quantities of data.
 - ▶ There initial works on how they work [14]

Remark

We have that

- A Convolutional Neural Network (CNN) directly accepts raw images as input.

Thus, their importance when training discrete filters

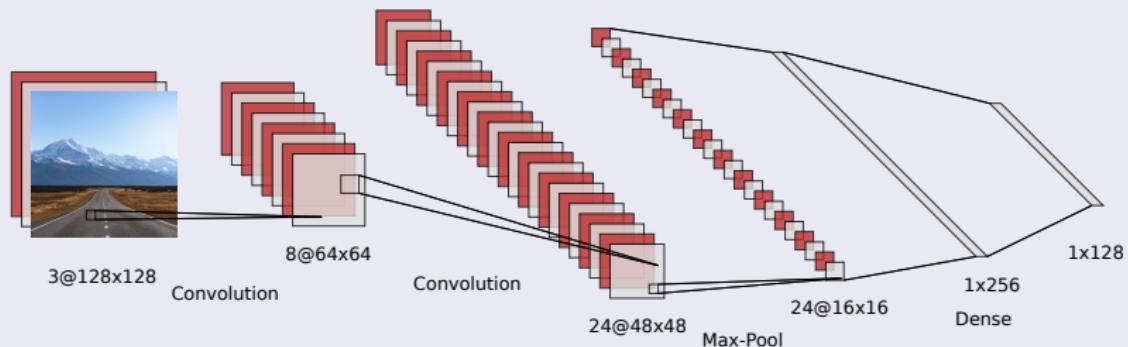
- Instead of assuming a certain comprehension of Computer Vision, one could think this is as a Silver Bullet.

However, you still

- You still need to be aware of :
 - ▶ The need of great quantities of data.
 - ▶ There initial works on how they work [14]

Example of CNN

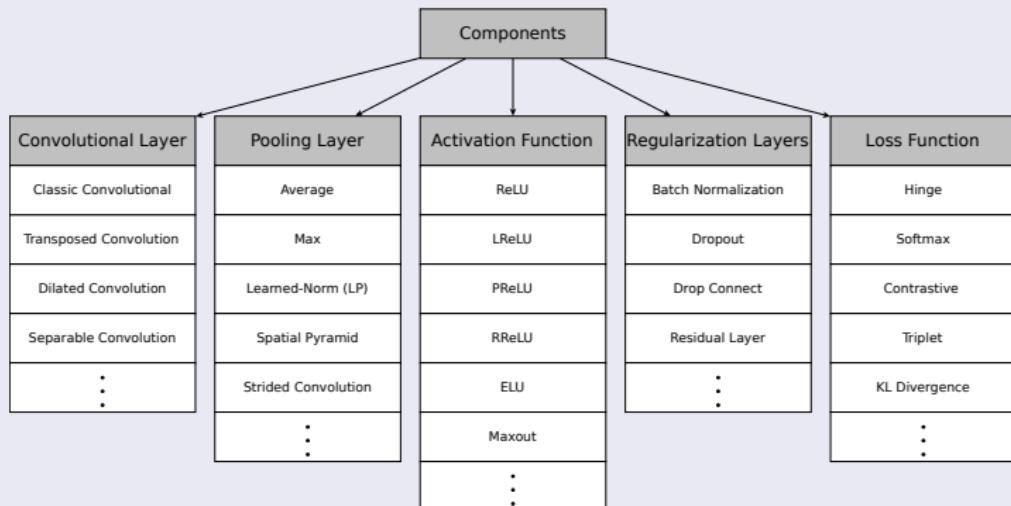
A Basic Convolutional Network



However

We will see that there are many possible architectures

- And different layers [15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25] :



Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

We know that

Many of the existing machine learning algorithms

- They depend on the quality of the input characteristics to generate a good model.

Dimensionality

- The amount of these variables is also important, given that performance tends to decline as the input dimensionality increases.

We know that

Many of the existing machine learning algorithms

- They depend on the quality of the input characteristics to generate a good model.

Not only that

- The amount of these variables is also important, given that performance tends to decline as the input dimensionality increases.

We have several techniques for that

Principal Component Analysis

$$L(\mathbf{u}_1) = \mathbf{u}_1^T S \mathbf{u}_1 + \lambda_1 \left(1 - \mathbf{u}_1^T \mathbf{u}_1\right)$$

Linear regression

$$\Phi(Y) = \sum_i \left| Y_i - \sum_j W_{ij} Y_j \right|^2$$

Manifolds

- Uniform Manifold Approximation and Projection for Dimension Reduction [26]

We have several techniques for that

Principal Component Analysis

$$L(\mathbf{u}_1) = \mathbf{u}_1^T S \mathbf{u}_1 + \lambda_1 \left(1 - \mathbf{u}_1^T \mathbf{u}_1\right)$$

Linear Locally Embeddings

$$\Phi(Y) = \sum_i \left| Y_i - \sum_j W_{ij} Y_j \right|^2$$

And finally

- Uniform Manifold Approximation and Projection for Dimension Reduction [26]

We have several techniques for that

Principal Component Analysis

$$L(\mathbf{u}_1) = \mathbf{u}_1^T S \mathbf{u}_1 + \lambda_1 \left(1 - \mathbf{u}_1^T \mathbf{u}_1\right)$$

Linear Locally Embeddings

$$\Phi(Y) = \sum_i \left| Y_i - \sum_j W_{ij} Y_j \right|^2$$

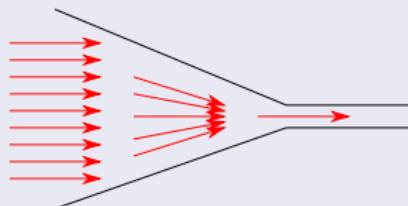
And recently

- Uniform Manifold Approximation and Projection for Dimension Reduction [26]

Therefore

We have the need to codify the original feature into better ones

- This can be done by a series of mappings that act as funnels, How?



Essentially, we have a series of mapping

$$x \in \mathbb{R}^{n_1} \rightarrow f_1(x) \in \mathbb{R}^{n_2} \rightarrow f_2(x_1) \in \mathbb{R}^{n_3} \dots \rightarrow f_m(x_{m-1}) \in \mathbb{R}^{n_{m+1}}$$

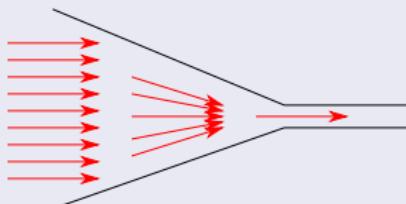
Where

$$n_1 < n_2 < \dots < n_m < n_{m+1}$$

Therefore

We have the need to codify the original feature into better ones

- This can be done by a series of mappings that act as funnels, How?



Basically, we have a series of mappings

$$x \in \mathbb{R}^{n_1} \rightarrow f_1(x) \in \mathbb{R}^{n_2} \rightarrow f_2(x_1) \in \mathbb{R}^{n_3} \dots \rightarrow f_m(x_{m-1}) \in \mathbb{R}^{n_{m+1}}$$

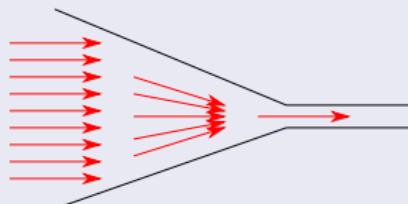
Where

$$n_1 < n_2 < \dots < n_m < n_{m+1}$$

Therefore

We have the need to codify the original feature into better ones

- This can be done by a series of mappings that act as funnels, How?



Basically, we have a series of mappings

$$x \in \mathbb{R}^{n_1} \rightarrow f_1(x) \in \mathbb{R}^{n_2} \rightarrow f_2(x_1) \in \mathbb{R}^{n_3} \dots \rightarrow f_m(x_{m-1}) \in \mathbb{R}^{n_{m+1}}$$

Where

$$n_1 < n_2 < \dots < n_m < n_{m+1}$$

Then, we can use linear mappings for this

With the following matrix functions

$$\sigma [f_{A_{i+1}}(x_i)] = \sigma(A_{i+1}x)$$

- Therefore, we have the following architecture.

Then, we can use linear mappings for this

With the following matrix functions

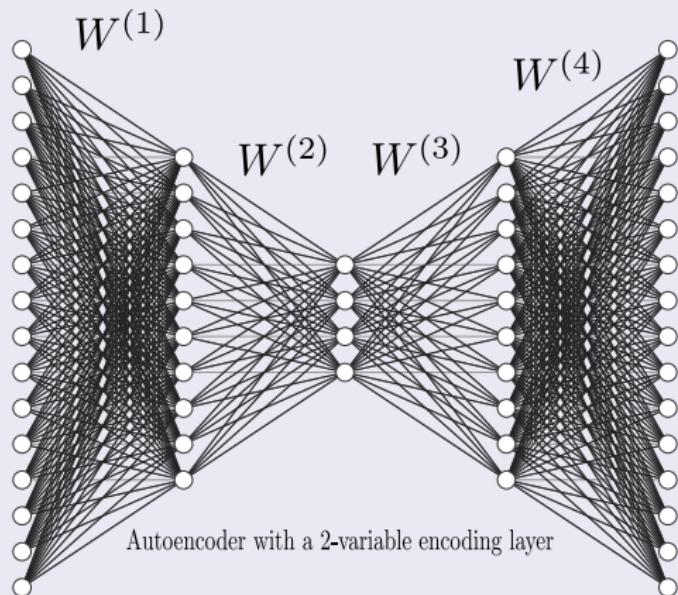
$$\sigma [f_{A_{i+1}}(x_i)] = \sigma(A_{i+1}x)$$

Therefore

- Therefore, we have the following architecture.

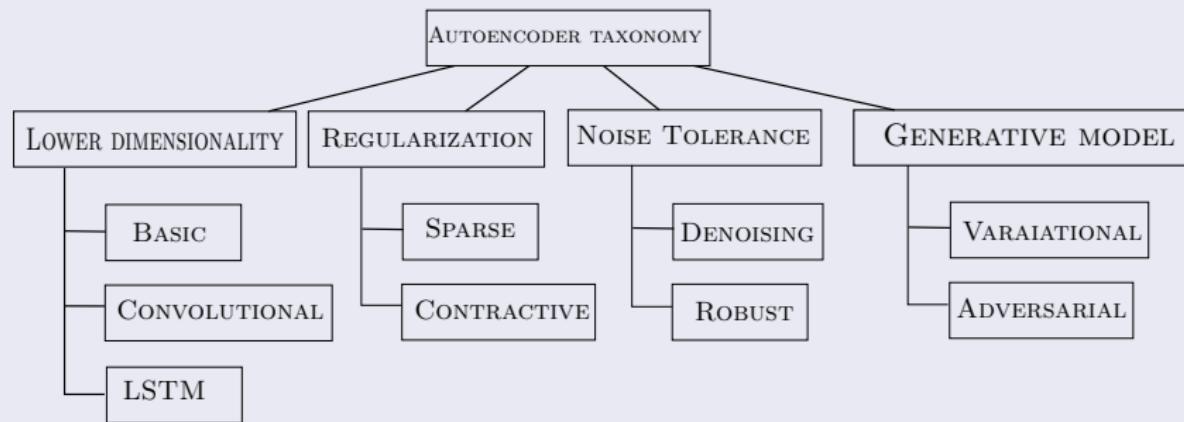
The Basic Auto Encoder Architecture

We have



Taxonomy

Most popular Auto Encoders



Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks**
 - An Example, Boltzmann Machines
 - There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

Generative Adversarial Networks

They can be seen as an Accept-Reject MCMC Model

- However, they do not require Markov Chains with the classic problem:
 - ▶ The independence between the samples to generate ergodic probabilities (The real one)

Both networks compete

- The generator network tries to produce realistic-looking samples
- The discriminator network tries to figure out whether an image came from the training set or the generator network

Generative Adversarial Networks

They can be seen as an Accept-Reject MCMC Model

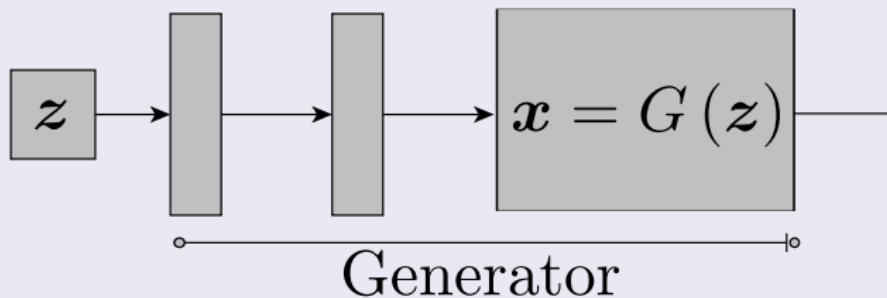
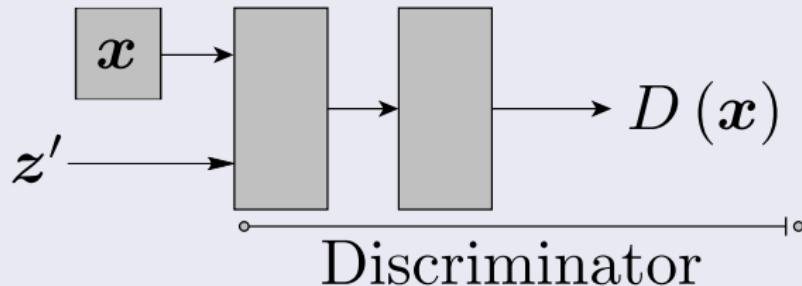
- However, they do not require Markov Chains with the classic problem:
 - ▶ The independence between the samples to generate ergodic probabilities (The real one)

As in the Accept-Reject

- The **generator network** tries to produce realistic-looking samples
- The **discriminator network** tries to figure out whether an image came from the training set or the generator network

Graphically

We have the following Basic Model



Here

There is a need to join both functions

- So, we can use the idea of Backpropagation to obtain the desired minimization.

- We can define a sensible learning criterion when the dataset is not linearly separable

For this, we can use the logistic cross-entropy loss (We will explain more about this later)

$$\mathcal{L}_{LCE}(z, t) = LCE(\sigma(z), t) = t \log(1 + e^{-z}) + (1 - t) \log(1 + e^z)$$

Here

There is a need to join both functions

- So, we can use the idea of Backpropagation to obtain the desired minimization.

How can we do this?

- We can define a sensible learning criterion when the dataset is not linearly separable

For this we can use the logistic cross-entropy loss function and more about this later

$$\mathcal{L}_{LCE}(z, t) = LCE(\sigma(z), t) = t \log(1 + e^{-z}) + (1 - t) \log(1 + e^z)$$

Here

There is a need to join both functions

- So, we can use the idea of Backpropagation to obtain the desired minimization.

How can we do this?

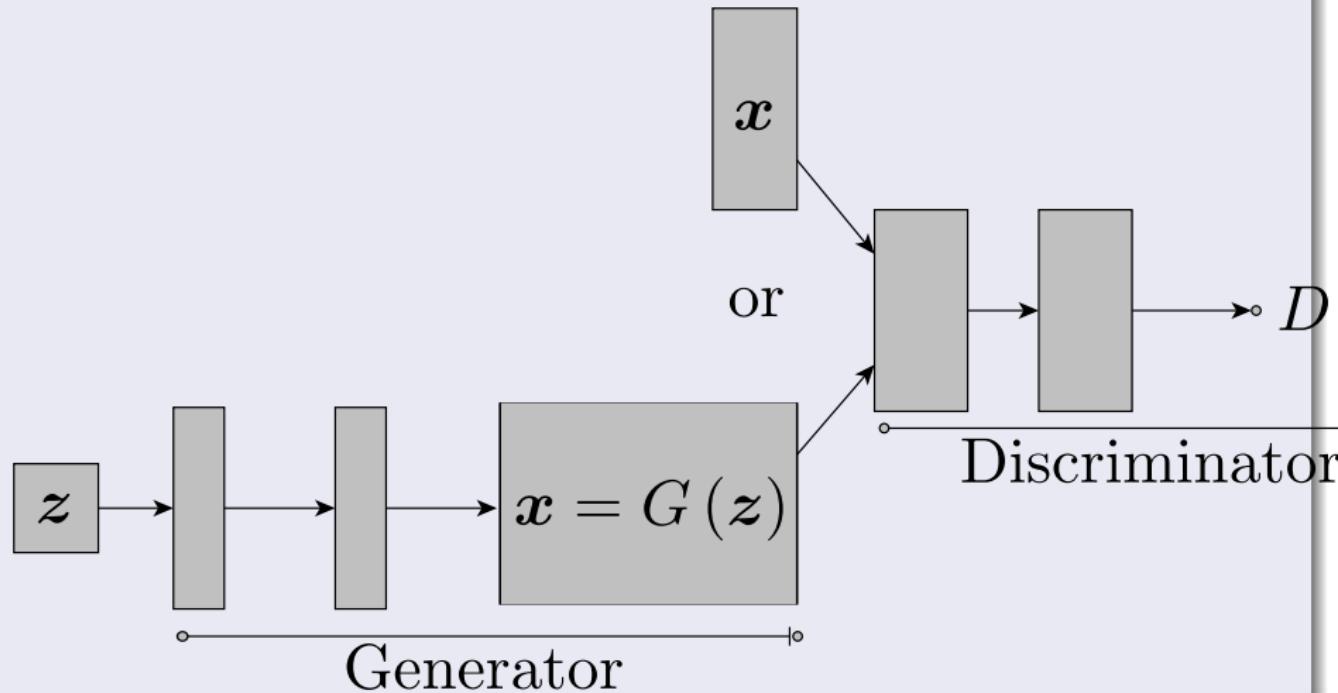
- We can define a sensible learning criterion when the dataset is not linearly separable

For this, we can use the **logistic cross-entropy loss** (We will explain more about this later)

$$\mathcal{L}_{LCE}(z, t) = L_{CE}(\sigma(z), t) = t \log(1 + e^{-z}) + (1 - t) \log(1 + e^z)$$

Therefore, we have

The following architecture use this idea



In this basic Generator

D denote the discriminator's predicted probability of being data

$$\mathcal{J}_D = E_{\mathbf{x} \sim \mathcal{D}} [-\log D(\mathbf{x})] + E_{\mathbf{z}} [-\log (1 - D(G(\mathbf{z})))]$$

Can't we do the same for the generator?

$$\mathcal{J}_G = -\mathcal{J}_D = const + E_{\mathbf{z}} [\log (1 - D(G(\mathbf{z})))]$$

In this basic Generator

D denote the discriminator's predicted probability of being data

$$\mathcal{J}_D = E_{\mathbf{x} \sim \mathcal{D}} [-\log D(\mathbf{x})] + E_{\mathbf{z}} [-\log (1 - D(G(\mathbf{z})))]$$

One possible cost function for the generator

$$\mathcal{J}_G = -\mathcal{J}_D = const + E_{\mathbf{z}} [\log (1 - D(G(\mathbf{z})))]$$

Then using both functions

The minimax formulation

- Since the generator and discriminator are playing a zero-sum game against each other.

Discriminator

$$\max_G \min_D \mathcal{J}_D$$

Then we either consider using the loss $\mathbb{E}[\cdot]$

$$\mathcal{J}_G = \frac{1}{N} \sum_{i=1}^N [G(z) - x_i]^2$$

Then using both functions

The minimax formulation

- Since the generator and discriminator are playing a zero-sum game against each other.

Basically

$$\max_G \min_D \mathcal{J}_D$$

Let's take a look at the generator loss

$$\mathcal{J}_G = \frac{1}{N} \sum_{i=1}^N [G(z) - x_i]^2$$

Then using both functions

The minimax formulation

- Since the generator and discriminator are playing a zero-sum game against each other.

Basically

$$\max_G \min_D \mathcal{J}_D$$

There are other examples using the LSE [27]

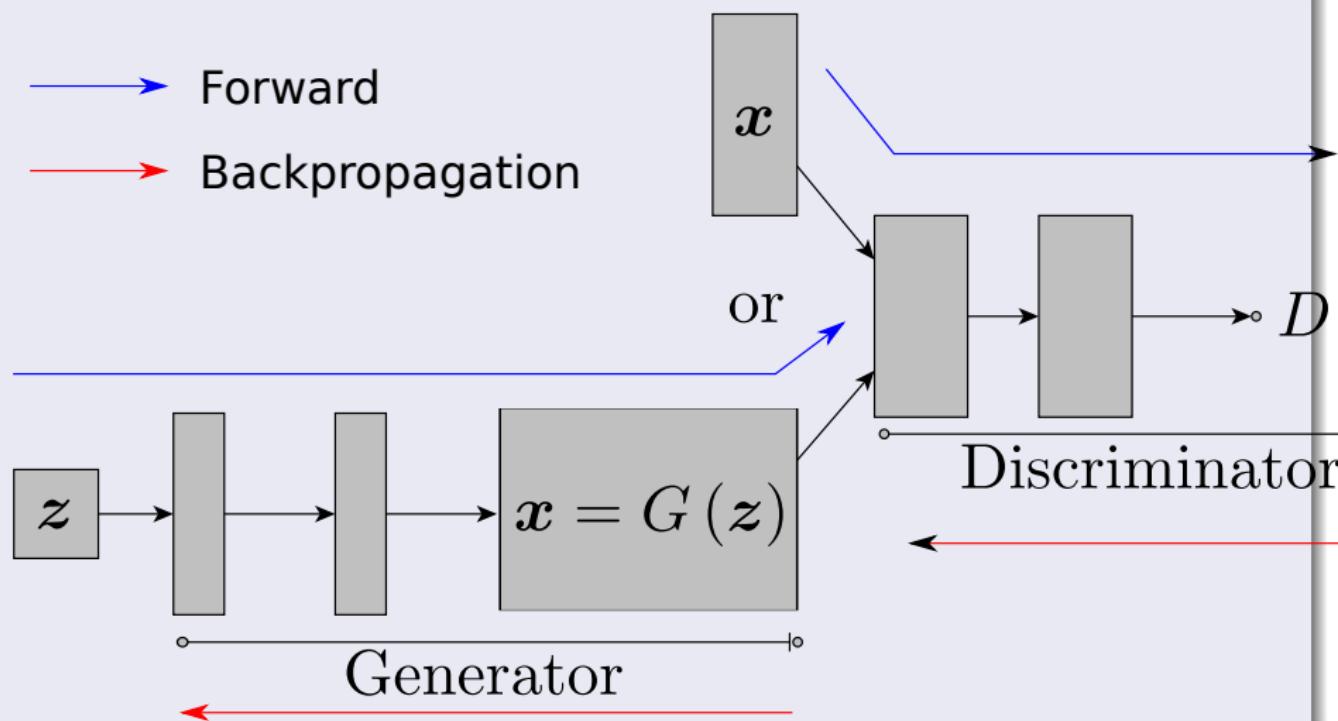
$$\mathcal{J}_G = \frac{1}{N} \sum_{i=1}^N [G(\mathbf{z}) - \mathbf{x}]^2$$

Therefore, we have two updates

First update the Discriminator

→ Forward

→ Backpropagation



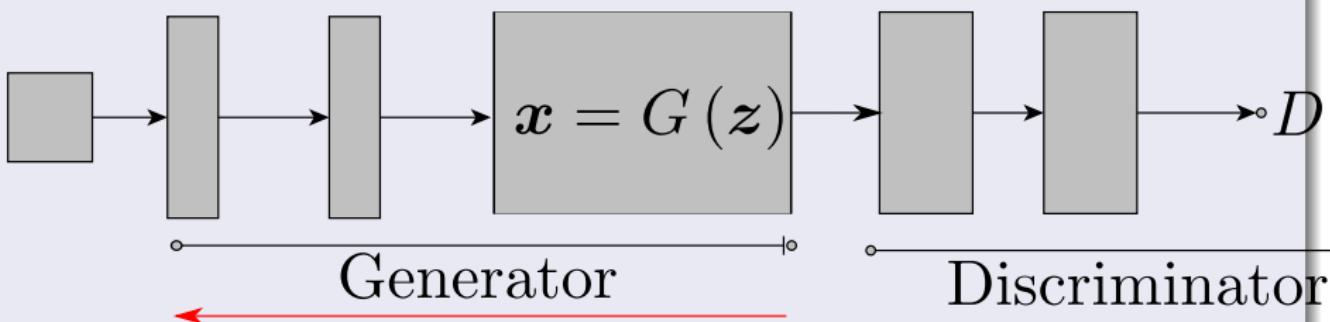
Now

Update the Generator

Backprop Derivatives Through the Discriminator, but do not change variables on it... only in the generator

→ Forward

→ Backpropagation



Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

The Basic Energy Models

We have that the Boltzmann Machines

- A Boltzmann machine is a network of units that are connected to each other

Binary values of the connected units

- Each unit takes a binary value in $\{0, 1\}$

Represented by a random variable X_i , $i = 1, \dots, N$.

Additional parameters

- Bias b_i

- Weight w_{ij} between unit i and unit j , $(i, j) \in [1, N - 1] \times [i + 1, N]$

The Basic Energy Models

We have that the Boltzmann Machines

- A Boltzmann machine is a network of units that are connected to each other

Here, we have N be the number of units

- Each unit takes a binary value in $\{0, 1\}$
 - ▶ Represented by a random variable X_i , $i = 1, \dots, N$.

parameters

- Bias b_i

• Weight w_{ij} between unit i and unit j , $(i, j) \in [1, N - 1] \times [i + 1, N]$

The Basic Energy Models

We have that the Boltzmann Machines

- A Boltzmann machine is a network of units that are connected to each other

Here, we have N be the number of units

- Each unit takes a binary value in $\{0, 1\}$
 - ▶ Represented by a random variable X_i , $i = 1, \dots, N$.

Additionally, it has parameters

- Bias b_i
- Weight w_{ij} between unit i and unit j , $(i, j) \in [1, N - 1] \times [i + 1, N]$

The Energy Based Structure

The energy of the Boltzmann machine is defined by

$$E_{W,b}[\mathbf{x}] = -\sum_{i=1}^N b_i x_i - \sum_{i=1}^{N-1} \sum_{j=i+1}^N w_{ij} x_i x_j = -\mathbf{b}^T \mathbf{x} - \mathbf{x}^T W \mathbf{x}$$

This allows us to define a probability distribution:

$$p_{W,b}(\mathbf{x}) = \frac{\exp(-E_{W,b}[\mathbf{x}])}{\sum_{\tilde{\mathbf{x}}} \exp(-E_{W,b}[\tilde{\mathbf{x}}])}$$

The Energy Based Structure

The energy of the Boltzmann machine is defined by

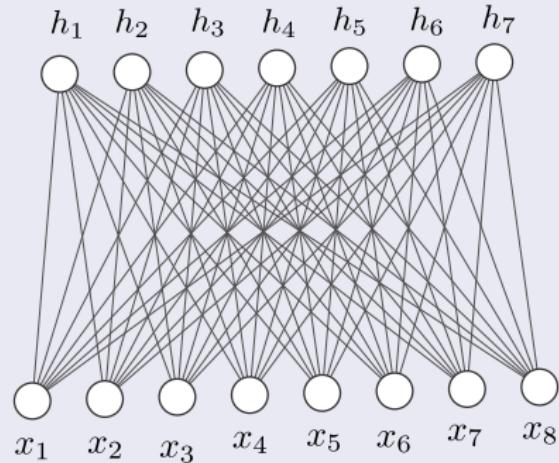
$$E_{W,b} [\mathbf{x}] = - \sum_{i=1}^N b_i x_i - \sum_{i=1}^{N-1} \sum_{j=i+1}^N w_{ij} x_i x_j = -\mathbf{b}^T \mathbf{x} - \mathbf{x}^T W \mathbf{x}$$

This allows to define a probability distribution

$$\mathbb{P}_{W,b} (\mathbf{x}) = \frac{\exp(-E_{W,b} [\mathbf{x}])}{\sum_{\tilde{\mathbf{x}}} \exp(-E_{W,b} [\tilde{\mathbf{x}}])}$$

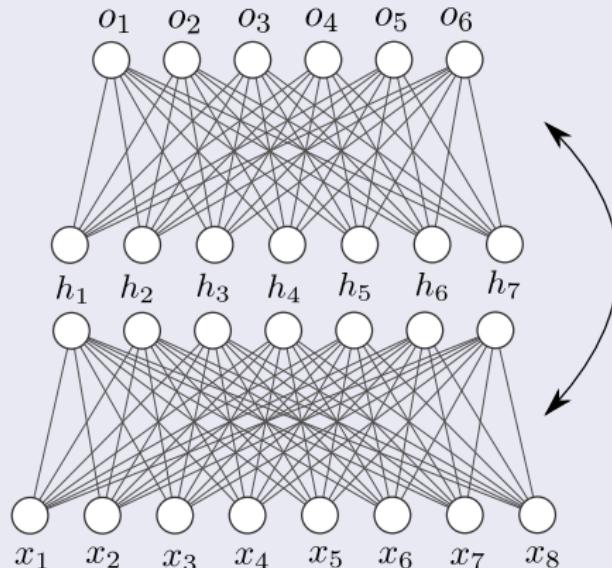
Example

Restricted Boltzmann Machines where the connectivity is layer by layer



Thus, using it as a basic model

We can stack them into a multiple layer model



Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

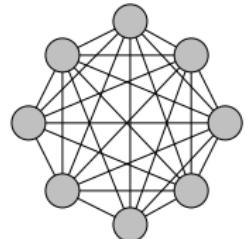
- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

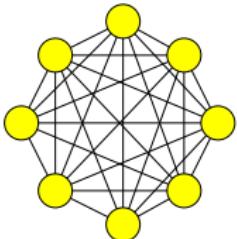
- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

There Are Many More!!! Here a few more...

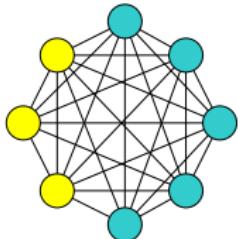
Markov Chain



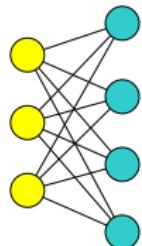
Hopfield Network



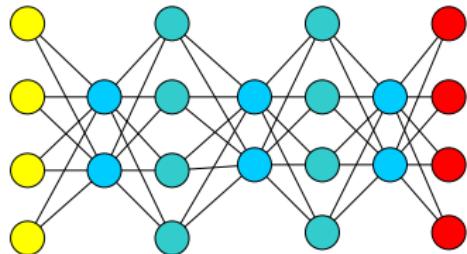
Boltzmann Machine



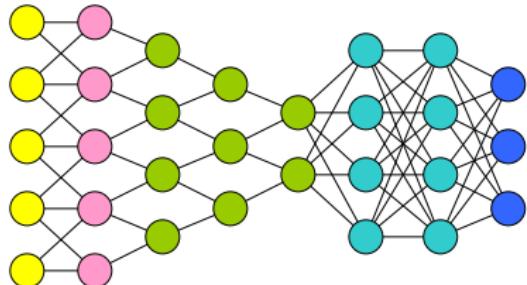
Restricted BM



Deep Belief Network

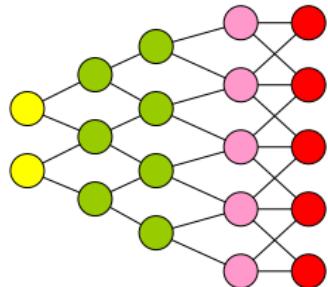


Convolutional Network

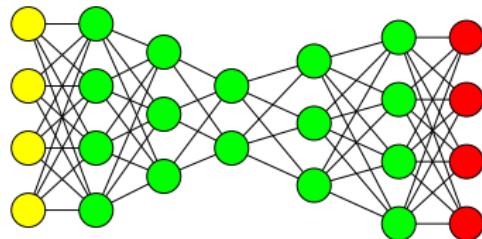


Furthermore

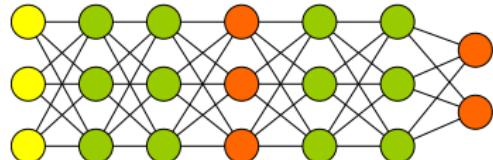
Deconvolutional Network



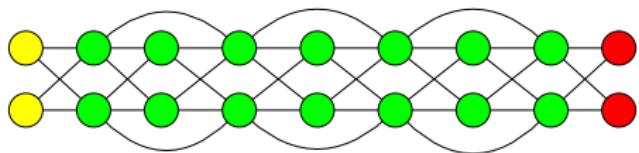
Autoencoder



Generative Adversarial Network



Deep Residual Network



Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

● The Degradation Problem

- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

Definition

Degradation Problem

- With the network depth increasing, accuracy gets saturated (which might be unsurprising) and then degrades rapidly.

Something Mysterious

- Unexpectedly, such degradation is not caused by overfitting,

Something Else

- to a suitably deep model leads to higher training error,

Definition

Degradation Problem

- With the network depth increasing, accuracy gets saturated (which might be unsurprising) and then degrades rapidly.

Something Notable

- Unexpectedly, such degradation is not caused by overfitting,

- to a suitably deep model leads to higher training error,

Definition

Degradation Problem

- With the network depth increasing, accuracy gets saturated (which might be unsurprising) and then degrades rapidly.

Something Notable

- Unexpectedly, such degradation is not caused by overfitting,

and adding more layers

- to a suitably deep model leads to higher training error,

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- **The Vanishing and Exploding Gradients**
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

As We know

In Recurrent Neural Networks, we have the problem

- Vanishing and Exploding Gradients

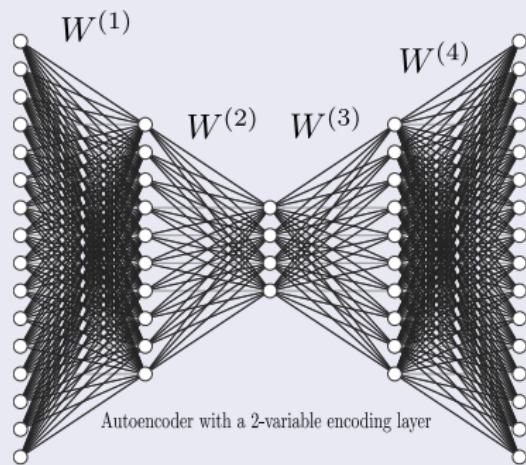
In this section, we will discuss some countermeasures to have much better training.

As We know

In Recurrent Neural Networks, we have the problem

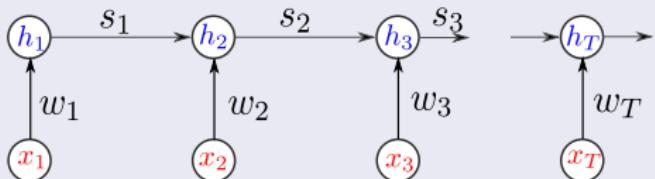
- Vanishing and Exploding Gradients

In the Deeper Architectures as encoder-decoder we have such phenomena



Consider a simple encoder encoder network

We have this simplified version

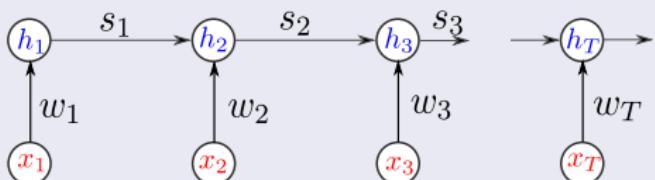


$$h_t = w_t x_t + z_{t-1}$$

$$z_t = s_t h_t$$

Consider a simple encoder encoder network

We have this simplified version



We have the following structure

$$h_t = w_t x_t + z_{t-1}$$

$$z_t = s_t h_t$$

Backpropagation Rules

Then, we get the following backpropagation rules

$$\frac{\partial h_t}{\partial w_i} = \frac{\partial h_t}{\partial h_{t-1}} \times \frac{\partial h_{t-1}}{\partial h_{t-2}} \times \dots \times \frac{\partial h_i}{\partial w_i}$$
$$\frac{\partial h_t}{\partial s_i} = \frac{\partial h_t}{\partial h_{t-1}} \times \frac{\partial h_{t-1}}{\partial h_{t-2}} \times \dots \times \frac{\partial h_{i+1}}{\partial s_i}$$

Then, we have

By Using Our simplifying assumption that

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial (w_t x_t + s_{t-1} h_{t-1})}{\partial h_{t-1}} = s_{t-1}$$

and for i

$$\frac{\partial h_t}{\partial w_i} = x_t$$

Finally, we have that

$$\frac{\partial h_t}{\partial w_i} = x_t \begin{bmatrix} & & & \\ & \text{H} & & \\ & & & \\ & & & s_k \\ i-1 & & & \\ \hline & & & \\ k-t-1 & & & \end{bmatrix}$$

Then, we have

By Using Our simplifying assumption that

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial (w_t x_t + s_{t-1} h_{t-1})}{\partial h_{t-1}} = s_{t-1}$$

And for $\frac{\partial h_i}{\partial w_i}$

$$\frac{\partial h_i}{\partial w_i} = x_t$$

Finally, we have that

$$\frac{\partial h_i}{\partial w_i} = x_t \begin{bmatrix} & & \\ & \vdots & \\ & s_k & \\ & \vdots & \\ & & \end{bmatrix}_{k=t-1}^{t-1}$$

Then, we have

By Using Our simplifying assumption that

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial (w_t x_t + s_{t-1} h_{t-1})}{\partial h_{t-1}} = s_{t-1}$$

And for $\frac{\partial h_i}{\partial w_i}$

$$\frac{\partial h_i}{\partial w_i} = x_t$$

Finally, we have that

$$\frac{\partial h_t}{\partial w_i} = x_t \left[\prod_{k=t-1}^{i-1} s_k \right]$$

It is clear that

Unless the s_k 's are near to 1

- You have the vanishing gradient if $s_k \in [0, 1)$ for all k .
- You have the exploding gradient if $s_k \in (1, +\infty]$ for all k .

Even though it is not mentioned

- These terms tend to appear in the Deep Learners when Backpropagation is done

In the backpropagation

- We have many activation function that squash the signal...

It is clear that

Unless the s_k 's are near to 1

- You have the vanishing gradient if $s_k \in [0, 1)$ for all k .
- You have the exploding gradient if $s_k \in (1, +\infty]$ for all k .

Even with activation functions

- These terms tend to appear in the Deep Learners when Backpropagation is done

In the backpropagation

- We have many activation function that squash the signal...

It is clear that

Unless the s_k 's are near to 1

- You have the vanishing gradient if $s_k \in [0, 1]$ for all k .
- You have the exploding gradient if $s_k \in (1, +\infty]$ for all k .

Even with activation functions

- These terms tend to appear in the Deep Learners when Backpropagation is done

In the case of Forward

- We have many activation function that squash the signal...

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively**
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

Instead of doing this

Let us do the following

$$f(x) = 3.5x(1 - x)$$

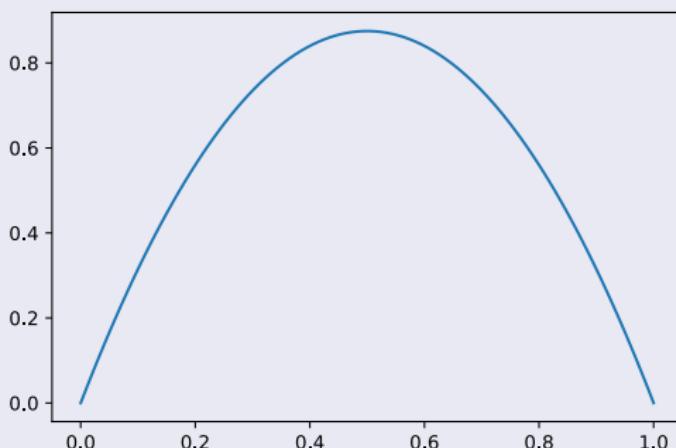
Find the composition rule

Instead of doing this

Let us do the following

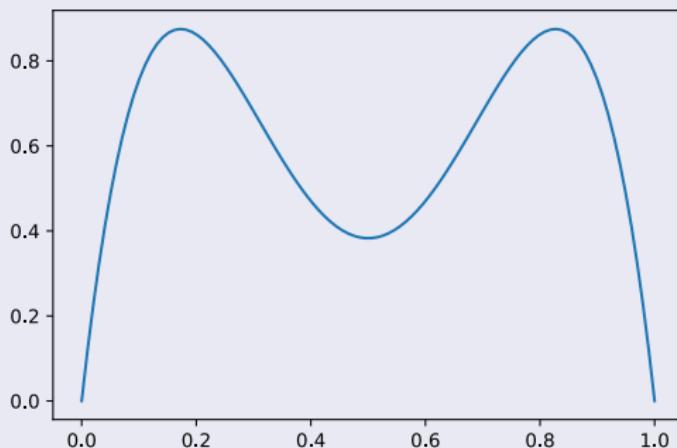
$$f(x) = 3.5x(1 - x)$$

In the first composition, we get



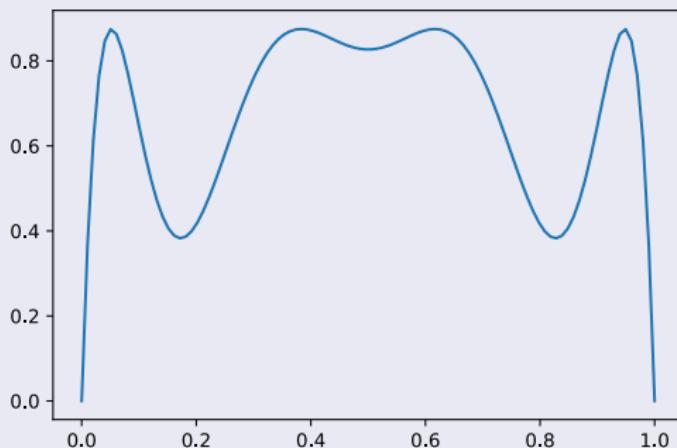
Now, as we compound the function

Second one, $y = f \circ f (x)$



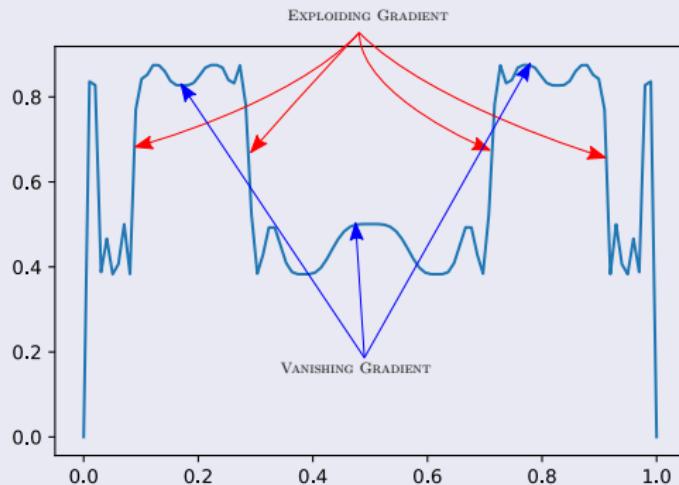
Now, as we increment iterations

Third one, $y = f \circ f \circ f (x)$



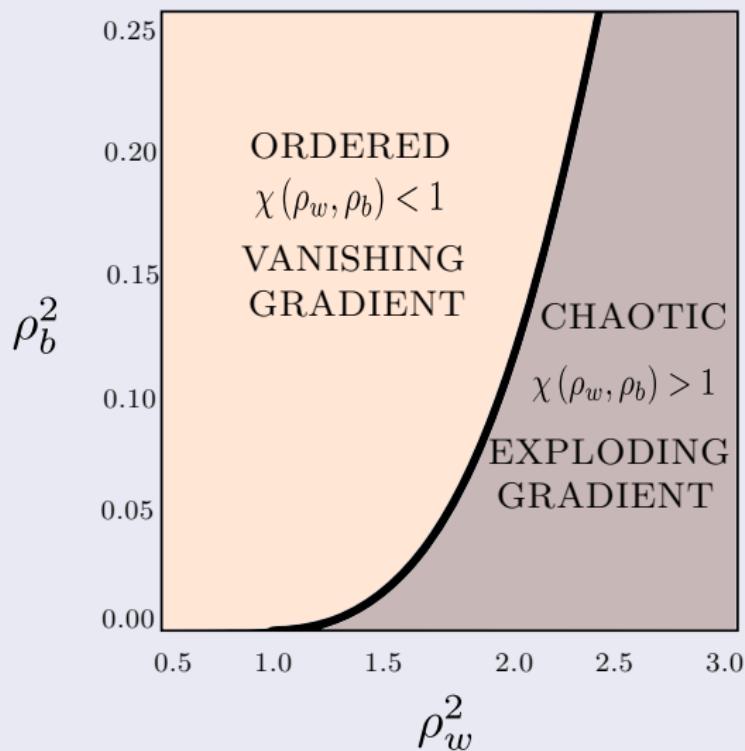
Finally

We see the increment in the gradient part negative or positive



Actually, we have

A Frontier defining the Vanishing and Exploding Gradient [28]



Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points**
 - Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
 - Conclusions

Actually

Eventually, the iterates go to infinity or zero OR

- They wind up at a fixed point...

• Fixed Point

$$x = f(x)$$

Actually

Eventually, the iterates go to infinity or zero OR

- They wind up at a fixed point...

A Fixed Point?

$$x = f(x)$$

Basically

The fixed points can be thought

- Some fixed points repel the iterates; **these are called sources.**
- Other fixed points attract the iterates; **these are called sinks.**

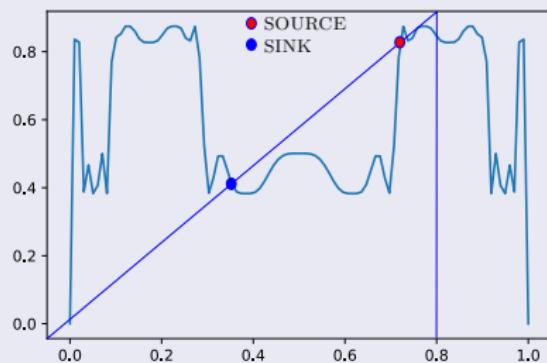
[Source: https://en.wikipedia.org/wiki/Attractor_(mathematics)]

Basically

The fixed points can be thought

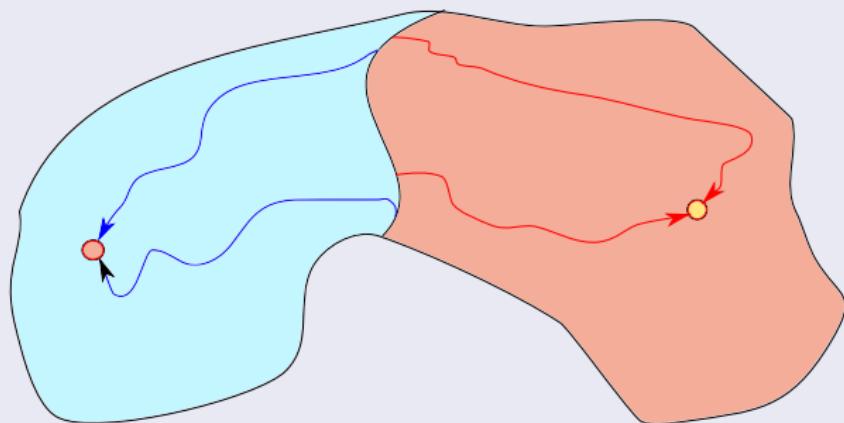
- Some fixed points repel the iterates; **these are called sources**.
- Other fixed points attract the iterates; **these are called sinks**.

Basically $f'(x) < 1$ are sinks and $f'(x) > 1$ are sources



Areas of attraction

Basically, we have that there are areas the pull in the iterations of the function



These fixed points

In Deep Structures as RNN without sigmoid functions

$$\begin{aligned}\mathbf{h}_t &= W_{sd} \mathbf{x}_t + U_{ss} \mathbf{h}_{t-1} \\ \mathbf{y}_t &= V_{os} \mathbf{h}_t\end{aligned}$$

We have

$$\mathbf{x}_t = V_{os} [W_{sd} \mathbf{x}_t + U_{ss} \mathbf{h}_{t-1}]$$

The eigenvalues of V_{os}

- Then, we have that

$$\mathbf{x}_t = V_{os} W_{sd} \mathbf{x}_t + V_{os} U_{ss} \mathbf{h}_{t-1} = I \mathbf{x}_t + 0$$

These fixed points

In Deep Structures as RNN without sigmoid functions

$$\begin{aligned}\mathbf{h}_t &= W_{sd} \mathbf{x}_t + U_{ss} \mathbf{h}_{t-1} \\ \mathbf{y}_t &= V_{os} \mathbf{h}_t\end{aligned}$$

We have

$$\mathbf{x}_t = V_{os} [W_{sd} \mathbf{x}_t + U_{ss} \mathbf{h}_{t-1}]$$

The original problem is:

- Then, we have that

$$\mathbf{x}_t = V_{os} W_{sd} \mathbf{x}_t + V_{os} U_{ss} \mathbf{h}_{t-1} = I \mathbf{x}_t + 0$$

These fixed points

In Deep Structures as RNN without sigmoid functions

$$\begin{aligned}\mathbf{h}_t &= W_{sd} \mathbf{x}_t + U_{ss} \mathbf{h}_{t-1} \\ \mathbf{y}_t &= V_{os} \mathbf{h}_t\end{aligned}$$

We have

$$\mathbf{x}_t = V_{os} [W_{sd} \mathbf{x}_t + U_{ss} \mathbf{h}_{t-1}]$$

Therefore if $\mathbf{b} = V_{os} U_{ss} \mathbf{h}_{t-1}$

- Then, we have that

$$\mathbf{x}_t = V_{os} W_{sd} \mathbf{x}_t + V_{os} U_{ss} \mathbf{h}_{t-1} = I \mathbf{x}_t + \mathbf{0}$$

Therefore

We have that

$$V_{os}W_{sd} \approx I \text{ and } \mathbf{h}_{t-1} \approx 0$$

They define an area

Where V_{os} and W_{sd}

- They are the inverse of each other

↳ [View the full document on my personal website](#)

- Basically they fixed point converts a RNN without activation functions in a linear model

They define an area

Where V_{os} and W_{sd}

- They are the inverse of each other

And the hidden state is almost zero

- Basically they fixed point converts a RNN without activation functions in a linear model

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network**
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- **Stabilizing the Network**
 - **Gradient Clipping**
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

Gradient Clipping

We prevent gradient from blowing up by rescaling to a certain value

$$\|\nabla_{\theta} L\| > \eta \implies \nabla_{\theta} L = \frac{\eta \nabla_{\theta} L}{\|\nabla_{\theta} L\|}$$

We have a series of slides on loss [2]

$$\min_{x \in \mathbb{R}^d} f(x)$$

Furthermore, we define a sphere

$$S = \{x | \exists y \text{ such that } f(y) \leq f(x_0), \text{ and } \|x - y\| \leq 1\}$$

Gradient Clipping

We prevent gradient from blowing up by rescaling to a certain value

$$\|\nabla_{\theta}L\| > \eta \implies \nabla_{\theta}L = \frac{\eta \nabla_{\theta}L}{\|\nabla_{\theta}L\|}$$

We have a series of nice analysis [29]

$$\min_{x \in \mathbb{R}^d} f(x)$$

Furthermore, we define a set

$$S = \{x | \exists y \text{ such that } f(y) \leq f(x_0), \text{ and } \|x - y\| \leq 1\}$$

Gradient Clipping

We prevent gradient from blowing up by rescaling to a certain value

$$\|\nabla_{\theta}L\| > \eta \implies \nabla_{\theta}L = \frac{\eta \nabla_{\theta}L}{\|\nabla_{\theta}L\|}$$

We have a series of nice analysis [29]

$$\min_{x \in \mathbb{R}^d} f(x)$$

Furthermore, we define a space

$$S = \{x | \exists y \text{ such that } f(y) \leq f(x_o), \text{ and } \|x - y\| \leq 1\}$$

We have then for S

In \mathbb{R}^2 the following example

Assumptions

Assumption 1

- Function f is lower bounded by f^*

- Function f is twice differentiable

Assumptions

Assumption 1

- Function f is lower bounded by f^*

Assumption 2

- Function f is twice differentiable

Then, there are the following proposals

The ordinary gradient descent

$$x_{k+1} = x_k - \eta \nabla f(x_k)$$

The Cauchy gradient descent

$$x_{k+1} = x_k - h_c \nabla f(x_k), \text{ where } h_c = \min \left\{ \eta_c, \frac{\eta_c}{\|\nabla f(x)\|} \right\}$$

Normalized gradient descent proposal

$$x_{k+1} = x_k - h_n \nabla f(x_k), \text{ where } h_n = \frac{\eta_c}{\|\nabla f(x)\| + \beta}$$

Then, there are the following proposals

The ordinary gradient descent

$$x_{k+1} = x_k - \eta \nabla f(x_k)$$

The Clipped Gradient Descent (CGD)

$$x_{k+1} = x_k - h_c \nabla f(x_k), \text{ where } h_c = \min \left\{ \eta_c, \frac{\gamma \eta_c}{\|\nabla f(x)\|} \right\}$$

Normalized gradient descent method

$$x_{k+1} = x_k - h_n \nabla f(x_k), \text{ where } h_n = \frac{\eta_c}{\|\nabla f(x)\| + \beta}$$

Then, there are the following proposals

The ordinary gradient descent

$$x_{k+1} = x_k - \eta \nabla f(x_k)$$

The Clipped Gradient Descent (CGD)

$$x_{k+1} = x_k - h_c \nabla f(x_k), \text{ where } h_c = \min \left\{ \eta_c, \frac{\gamma \eta_c}{\|\nabla f(x)\|} \right\}$$

Normalized Gradient Descent (NGD)

$$x_{k+1} = x_k - h_n \nabla f(x_k), \text{ where } h_n = \frac{\eta_c}{\|\nabla f(x)\| + \beta}$$

Remark

Clipped GD and NGD are almost equivalent

- If we set $\gamma\eta_c = \eta_n$ and $\eta_c = \frac{\eta_n}{\beta}$ then

$$\frac{1}{2}h_c \leq h_n \leq 2h_c$$

A Natural Question

Definition

- The objective f is called L -smooth if
$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\| \text{ for all } x, y \in \mathbb{R}^d$$

This is equivalent under some differentiable

$$\|\nabla^2 f(x)\| \leq L$$

This is known as strong convexity

$$f(y) \approx f(x) + \nabla^T f(x)(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(x)(y - x)$$

A Natural Question

Definition

- The objective f is called L -smooth if $\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$ for all $x, y \in \mathbb{R}^d$

This is equivalent under a twice differentiable f

$$\|\nabla^2 f(x)\| \leq L$$

Then you get the following approximation

$$f(y) \approx f(x) + \nabla^T f(x)(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(x)(y - x)$$

A Natural Question

Definition

- The objective f is called L -smooth if $\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$ for all $x, y \in \mathbb{R}^d$

This is equivalent under a twice differentiable f

$$\|\nabla^2 f(x)\| \leq L$$

Then, you get the following upper-bound

$$f(y) \approx f(x) + \nabla^T f(x)(y - x) + \frac{1}{2}(y - x)^T \nabla^2 f(x)(y - x)$$

Then, it is possible to use the 3 Assumption

We have that

$$f(y) \leq f(x) + \nabla^T f(x)(y - x) + \frac{1}{2}L\|y - x\|^2$$

Then taking all the other variables and assuming $y = x + h\mathbf{1}$

$$h^* = \arg \min_h \left[f(x) + h\|\nabla f(x)\|^2 + \frac{1}{2}Lh^2\|\nabla f(x)\|^2 \right] = \frac{1}{L}$$

Basically

- This choice of h leads to GD with a fixed step,

Then, it is possible to use the 3 Assumption

We have that

$$f(y) \leq f(x) + \nabla^T f(x)(y - x) + \frac{1}{2}L\|y - x\|^2$$

Then fixing all the other variables and assuming $y = x - h\nabla f(x)$

$$h^* = \arg \min_h \left[f(x) - h \|\nabla f(x)\|^2 + \frac{1}{2}Lh^2 \|\nabla f(x)\|^2 \right] = \frac{1}{L}$$

Exercise

- This choice of h leads to GD with a fixed step,

Then, it is possible to use the 3 Assumption

We have that

$$f(y) \leq f(x) + \nabla^T f(x)(y - x) + \frac{1}{2}L\|y - x\|^2$$

Then fixing all the other variables and assuming $y = x - h\nabla f(x)$

$$h^* = \arg \min_h \left[f(x) - h \|\nabla f(x)\|^2 + \frac{1}{2}Lh^2 \|\nabla f(x)\|^2 \right] = \frac{1}{L}$$

Basically

- This choice of h leads to GD with a fixed step,

Now

Question

- “Is clipped gradient descent optimized for a different smoothness condition?”

$$f(y) \leq f(x) + \nabla^T f(x)(y - x) + \frac{1}{2}L\|y - x\|^2$$

Answer

$$\eta^* = \frac{\eta}{\|\nabla f(x)\| + \beta}$$

Now

Question

- “Is clipped gradient descent optimized for a different smoothness condition?”

Inspired in the equation

$$f(y) \leq f(x) + \nabla^T f(x)(y - x) + \frac{1}{2}L\|y - x\|^2$$

ANSWER

$$\eta^* = \frac{\eta}{\|\nabla f(x)\| + \beta}$$

Now

Question

- “Is clipped gradient descent optimized for a different smoothness condition?”

Inspired in the equation

$$f(y) \leq f(x) + \nabla^T f(x)(y - x) + \frac{1}{2}L \|y - x\|^2$$

Assume

$$h^* = \frac{\eta}{\|\nabla f(x)\| + \beta}$$

Then, we have

Assume that such value optimize the equation

$$f(x) - h \|\nabla f(x)\|^2 + \frac{1}{2} L h^2 \|\nabla f(x)\|^2$$

Then we have

$$L(x) = \frac{\|\nabla f(x)\| + \beta}{\eta}$$

Assumption of the function f is smooth

- (L_0, L_1) -smoothness. f is (L_0, L_1) -smooth, if there exist positive L_0 and L_1 such that $\|\nabla^2 f(x)\| \leq L_0 + L_1 \|\nabla f(x)\|$
 - » $\nabla^2 f(x)$ is the Hessian

Then, we have

Assume that such value optimize the equation

$$f(x) - h \|\nabla f(x)\|^2 + \frac{1}{2} L h^2 \|\nabla f(x)\|^2$$

Then, we have

$$L(x) = \frac{\|\nabla f(x)\| + \beta}{\eta}$$

Assumption 3: the function $L(x)$ is smooth

- (L_0, L_1) -smoothness. f is (L_0, L_1) -smooth, if there exist positive L_0 and L_1 such that $\|\nabla^2 f(x)\| \leq L_0 + L_1 \|\nabla f(x)\|$
 - » $\nabla^2 f(x)$ is the Hessian

Then, we have

Assume that such value optimize the equation

$$f(x) - h \|\nabla f(x)\|^2 + \frac{1}{2} L h^2 \|\nabla f(x)\|^2$$

Then, we have

$$L(x) = \frac{\|\nabla f(x)\| + \beta}{\eta}$$

Assumption 3 by using $\|\nabla^2 f(x)\| \leq L$

- (L_0, L_1) -smoothness. f is (L_0, L_1) -smooth, if there exist positive L_0 and L_1 such that $\|\nabla^2 f(x)\| \leq L_0 + L_1 \|\nabla f(x)\|$
 - ▶ $\nabla^2 f(x)$ is the Hessian

The final Theorem

Theorem (CGD) [29]

- Assume that Assumptions 1, 2, and 3 hold in set S . With parameters

$$\eta_c = \frac{1}{10L_o} \text{ and } \gamma = \min \left\{ \frac{1}{\eta_c}, \frac{1}{10L_o\eta_c} \right\},$$

- Then Clipped GD terminates in

$$\frac{20L_0(f(x_0) - f^*)}{\epsilon^2} + \frac{20 \max \{1, L_1^2\} (f(x_0) - f^*)}{L_0} \text{ iterations}$$

Remarks

The paper

- It points out to a high correlation between the Jacobian and the Hessian

Please read the paper...

Remarks

The paper

- It points out to a high correlation between the Jacobian and the Hessian

There are more work to be done

- Please read the paper...

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network**
 - Gradient Clipping
 - Normalizing your Data**
 - Normalization Layer AKA Batch Normalization
- Conclusions

Another way to stabilize the network

Data Normalization

- Standardization is the most popular form of preprocessing
 - ▶ Normally mean subtraction and subsequent scaling by the standard deviation.

Mean subtraction

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i \text{ then } x'_i = x_i - \mu$$

Final

- Standardization refers to altering the data dimensions such that they are of approximately the same scale.

Another way to stabilize the network

Data Normalization

- Standardization is the most popular form of preprocessing
 - ▶ Normally mean subtraction and subsequent scaling by the standard deviation.

Mean subtraction

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i \text{ then } x_i^c = x_i - \mu$$

Final

- Standardization refers to altering the data dimensions such that they are of approximately the same scale.

Another way to stabilize the network

Data Normalization

- Standardization is the most popular form of preprocessing
 - ▶ Normally mean subtraction and subsequent scaling by the standard deviation.

Mean subtraction

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i \text{ then } x_i^c = x_i - \mu$$

Finally

- Standardization refers to altering the data dimensions such that they are of approximately the same scale.

Therefore, we have that

Standardization

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mu)^2$$
$$x_i^s = \frac{x_i - \mu}{\sigma}$$

However, there is one more step to do.

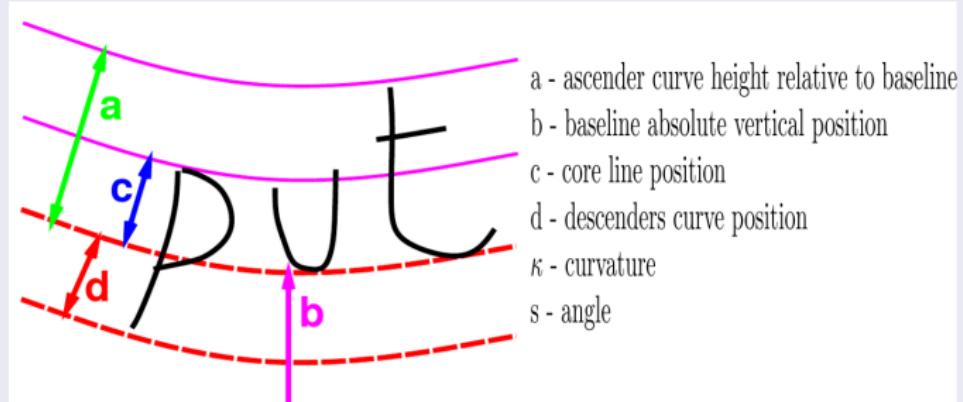
Therefore, we have that

Standardization

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mu)^2$$

$$x_i^s = \frac{x_i - \mu}{\sigma}$$

However, there other tricks, Bengio et al [30]



Softmax Scaling

Thus

- All new features have zero mean and unit variance.

However

- Other linear techniques limit the feature values in the range of $[0, 1]$ or $[-1, 1]$ by proper scaling.

However

- We can non-linear mapping. For example the softmax scaling

Softmax Scaling

Thus

- All new features have zero mean and unit variance.

Further

- Other linear techniques limit the feature values in the range of $[0, 1]$ or $[-1, 1]$ by proper scaling.

• We can non-linear mapping. For example the softmax scaling

Softmax Scaling

Thus

- All new features have zero mean and unit variance.

Further

- Other linear techniques limit the feature values in the range of $[0, 1]$ or $[-1, 1]$ by proper scaling.

However

- We can non-linear mapping. For example the softmax scaling.

Steps of Softmax Scaling

Softmax Scaling

- It consists of two steps

First step:

$$y_{ik} = \frac{x_{ik} - \bar{x}_k}{\sigma} \quad (3)$$

Second step:

$$\hat{x}_{ik} = \frac{1}{1 + \exp \{-y_{ik}\}} \quad (4)$$

Steps of Softmax Scaling

Softmax Scaling

- It consists of two steps

First one

$$y_{ik} = \frac{x_{ik} - \bar{x}_k}{\sigma} \quad (3)$$

$$\hat{x}_{ik} = \frac{1}{1 + \exp \{-y_{ik}\}} \quad (4)$$

Steps of Softmax Scaling

Softmax Scaling

- It consists of two steps

First one

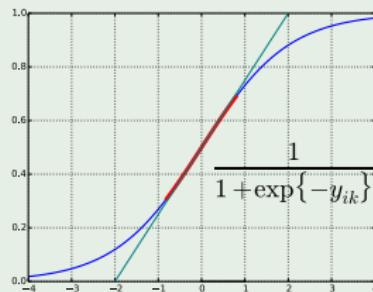
$$y_{ik} = \frac{x_{ik} - \bar{x}_k}{\sigma} \quad (3)$$

Second one

$$\hat{x}_{ik} = \frac{1}{1 + \exp \{-y_{ik}\}} \quad (4)$$

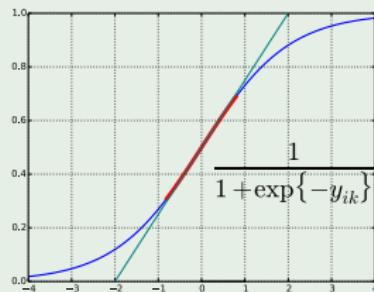
Explanation

Notice the red area is almost flat!!!



Explanation

Notice the red area is almost flat!!!



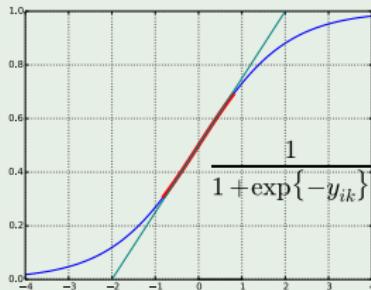
Thus, we have that

- The red region represents values of y inside of the region defined by the mean and variance (small values of y).

Then, if we have those values y behaves as a linear function.

Explanation

Notice the red area is almost flat!!!



Thus, we have that

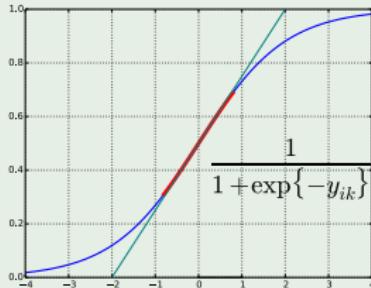
- The red region represents values of y inside of the region defined by the mean and variance (small values of y).
- Then, if we have those values x behaves as a linear function.

• They are squashed by the exponential part of the function.

• They are squashed by the exponential part of the function.

Explanation

Notice the red area is almost flat!!!



Thus, we have that

- The red region represents values of y inside of the region defined by the mean and variance (small values of y).
- Then, if we have those values x behaves as a linear function.

And values too away from the mean

- They are squashed by the exponential part of the function.

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network**
 - Gradient Clipping
 - Normalizing your Data
- Normalization Layer AKA Batch Normalization**
- Conclusions

Here, the people at Google [22] around 2015

They commented in the “Internal Covariate Shift Phenomena”

- Due to the change in the distribution of each layer’s input

Deep learning

- The min-batch forces to have those changes which impact on the learning capabilities of the network.

Internal Covariate Shift

- Internal Covariate Shift as the change in the distribution of network activations due to the change in network parameters during training.

Here, the people at Google [22] around 2015

They commented in the “Internal Covariate Shift Phenomena”

- Due to the change in the distribution of each layer’s input

They claim

- The min-batch forces to have those changes which impact on the learning capabilities of the network.

- Internal Covariate Shift as the change in the distribution of network activations due to the change in network parameters during training.

Here, the people at Google [22] around 2015

They commented in the “Internal Covariate Shift Phenomena”

- Due to the change in the distribution of each layer’s input

They claim

- The min-batch forces to have those changes which impact on the learning capabilities of the network.

In Neural Networks, they define this

- Internal Covariate Shift as the change in the distribution of network activations due to the change in network parameters during training.

They gave the following reasons

Consider a layer with the input u that adds the learned bias b

- Then, it normalizes the result by subtracting the mean of the activation over the training data:

$$\hat{x} = x - E[x]$$

- $\mathcal{X} = \{x, \dots, x_N\}$ the data samples and $E[x] = \frac{1}{N} \sum_{i=1}^N x_i$

Now, take the partial derivative. The dependence is on b , Δb , and x .

- Then $b = b + \Delta b$ where $\Delta b \propto -\frac{\partial l}{\partial x}$

$$u + (b + \Delta b) - E[u + (b + \Delta b)] = u + b - E[u + b]$$

They gave the following reasons

Consider a layer with the input u that adds the learned bias b

- Then, it normalizes the result by subtracting the mean of the activation over the training data:

$$\hat{x} = x - E[x]$$

- $\mathcal{X} = \{x, \dots, x_N\}$ the data samples and $E[x] = \frac{1}{N} \sum_{i=1}^N x_i$

Now, if the gradient ignores the dependence of $E[x]$ on b

- Then $b = b + \Delta b$ where $\Delta b \propto -\frac{\partial l}{\partial \hat{x}}$

$$u + (b + \Delta b) - E[u + (b + \Delta b)] = u + b - E[u + b]$$

They gave the following reasons

Consider a layer with the input u that adds the learned bias b

- Then, it normalizes the result by subtracting the mean of the activation over the training data:

$$\hat{\mathbf{x}} = \mathbf{x} - E[\mathbf{x}]$$

- $\mathcal{X} = \{\mathbf{x}, \dots, \mathbf{x}_N\}$ the data samples and $E[\mathbf{x}] = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$

Now, if the gradient ignores the dependence of $E[\mathbf{x}]$ on b

- Then $b = b + \Delta b$ where $\Delta b \propto -\frac{\partial l}{\partial \hat{\mathbf{x}}}$

Finally

$$u + (b + \Delta b) - E[u + (b + \Delta b)] = u + b - E[u + b]$$

Then

The following will happen

- The update to b leads to **no change** in the output of the layer.

- We need to integrate the normalization into the process of training.

Then

The following will happen

- The update to b leads to **no change** in the output of the layer.

Therefore

- We need to integrate the normalization into the process of training.

Normalization via Mini-Batch Statistic

It is possible to describe the normalization as a transformation layer

$$\hat{x} = \text{Norm}(x, \mathcal{X})$$

- Which depends on all the training samples \mathcal{X} which also depends on the layer parameters

For back-propagation we will need to compute the following terms

$$\frac{\partial \text{Norm}(x, \mathcal{X})}{\partial x} \text{ and } \frac{\partial \text{Norm}(x, \mathcal{X})}{\partial \mathcal{X}}$$

Normalization via Mini-Batch Statistic

It is possible to describe the normalization as a transformation layer

$$\hat{\mathbf{x}} = \text{Norm}(\mathbf{x}, \mathcal{X})$$

- Which depends on all the training samples \mathcal{X} which also depends on the layer parameters

For back-propagation, we will need to generate the following terms

$$\frac{\partial \text{Norm}(\mathbf{x}, \mathcal{X})}{\partial \mathbf{x}} \text{ and } \frac{\partial \text{Norm}(\mathbf{x}, \mathcal{X})}{\partial \mathcal{X}}$$

Normalization via Mini-Batch Statistic

Problem!!!

- whitening the layer inputs is expensive, as it requires computing the covariance matrix

$$Cov[\mathbf{x}] = E_{\mathbf{x} \in \mathcal{X}} [\mathbf{x}\mathbf{x}^T] \text{ and } E[\mathbf{x}] E[\mathbf{x}]^T$$

- ▶ To produce the whitened activations

Therefore

A Better Options, we can normalize each dimension

$$\hat{\mathbf{x}}^{(k)} = \frac{\mathbf{x}^{(k)} - \mu}{\sigma}$$

- with $\mu = E [\mathbf{x}^{(k)}]$ and $\sigma^2 = Var [\mathbf{x}^{(k)}]$

• Normalization has two effects:

- Simply normalizing each input of a layer may change what the layer can represent.

• We need to keep track of the mean and the variance

- Which can represent the identity transform

Therefore

A Better Options, we can normalize each dimension

$$\hat{\mathbf{x}}^{(k)} = \frac{\mathbf{x}^{(k)} - \mu}{\sigma}$$

- with $\mu = E [\mathbf{x}^{(k)}]$ and $\sigma^2 = Var [\mathbf{x}^{(k)}]$

This allows to speed up convergence

- Simply normalizing each input of a layer may change what the layer can represent.

- Which can represent the identity transform

Therefore

A Better Options, we can normalize each dimension

$$\hat{\mathbf{x}}^{(k)} = \frac{\mathbf{x}^{(k)} - \mu}{\sigma}$$

- with $\mu = E [\mathbf{x}^{(k)}]$ and $\sigma^2 = Var [\mathbf{x}^{(k)}]$

This allows to speed up convergence

- Simply normalizing each input of a layer may change what the layer can represent.

So, we need to insert a transformation in the network

- Which can represent the identity transform

The Transformation

The Linear transformation

$$\mathbf{y}^{(k)} = \gamma^{(k)} \hat{\mathbf{x}}^{(k)} + \beta^{(k)}$$

The Inverse transformation

- This allow to recover the identity by setting $\gamma^{(k)} = \sqrt{\text{Var}[\mathbf{x}^{(k)}]}$ and $\beta^{(k)} = E[\mathbf{x}^{(k)}]$ if necessary.

The Transformation

The Linear transformation

$$\mathbf{y}^{(k)} = \gamma^{(k)} \hat{\mathbf{x}}^{(k)} + \beta^{(k)}$$

The parameters $\gamma^{(k)}, \beta^{(k)}$

- This allow to recover the identity by setting $\gamma^{(k)} = \sqrt{Var [\mathbf{x}^{(k)}]}$ and $\beta^{(k)} = E [\mathbf{x}^{(k)}]$ if necessary.

Finally

Batch Normalizing Transform

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(x_i)\}$

- ➊ $\mu_B = \frac{1}{m} \sum_{i=1}^m x_i$
- ➋ $\sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu)^2$
- ➌ $\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$
- ➍ $y_i = \gamma^{(i)} \hat{x}_i + \beta = BN_{\gamma, \beta}(x_i)$

Finally

Batch Normalizing Transform

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(x_i)\}$

① $\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m x_i$

② $\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \mu)^2$

③ $\hat{x}_i = \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$

④ $y_i = \gamma^{(i)} \hat{x}_i + \beta = BN_{\gamma, \beta}(x_i)$

Finally

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_1 \dots \mathbf{x}_m\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

① $\mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$

② $\sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu)^2$

③ $\hat{\mathbf{x}}_i = \frac{\mathbf{x}_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$

④ $y_i = \gamma^{(i)} \hat{\mathbf{x}}_i + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$

Finally

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_1 \dots \mathbf{x}_m\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

$$① \mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$$

$$② \sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu)^2$$

$$③ \hat{\mathbf{x}} = \frac{\mathbf{x}_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$④ y_i = \gamma \hat{\mathbf{x}} + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$$

Finally

Batch Normalizing Transform

Input: Values of \mathbf{x} over a mini-batch: $\mathcal{B} = \{\mathbf{x}_1 \dots \mathbf{x}_m\}$, Parameters to be learned: γ, β

Output: $\{y_i = BN_{\gamma, \beta}(\mathbf{x}_i)\}$

$$① \mu_{\mathcal{B}} = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$$

$$② \sigma_{\mathcal{B}}^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu)^2$$

$$③ \hat{\mathbf{x}} = \frac{\mathbf{x}_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$④ \mathbf{y}_i = \gamma^{(k)} \hat{\mathbf{x}}_i + \beta = BN_{\gamma, \beta}(\mathbf{x}_i)$$

Backpropagation

We have the following equations by using the loss function l

$$① \frac{\partial l}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \times \gamma$$

$$② \frac{\partial l}{\partial \mu_B} = \sum_{i=1}^m \frac{\partial l}{\partial x_i} \times (x_i - \mu_B) \times \left(-\frac{1}{2}\right) \times (\sigma_B^2 + \epsilon)^{-\frac{3}{2}}$$

$$③ \frac{\partial l}{\partial \sigma_B^2} = \left(\sum_{i=1}^m \frac{\partial l}{\partial x_i} \times \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_B^2} \times \frac{\sum_{i=1}^m -2 \times (x_i - \mu_B)}{m}$$

$$④ \frac{\partial l}{\partial x_i} = \frac{\partial l}{\partial y_i} \times \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_B^2} \times \frac{2 \times (x_i - \mu_B)}{m} + \frac{\partial l}{\partial \mu_B} \times \frac{1}{m}$$

$$⑤ \frac{\partial l}{\partial \gamma} = \sum_{i=1}^m \frac{\partial l}{\partial y_i} \times \hat{x}_i$$

$$⑥ \frac{\partial l}{\partial \epsilon} = \sum_{i=1}^m \frac{\partial l}{\partial y_i}$$

Backpropagation

We have the following equations by using the loss function l

$$① \frac{\partial l}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \times \gamma$$

$$② \frac{\partial l}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \times (x_i - \mu_B) \times \left(-\frac{1}{2}\right) \times (\sigma_B^2 + \epsilon)^{-\frac{3}{2}}$$

$$③ \frac{\partial l}{\partial \mu_B} = \left(\sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \times \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_B^2} \times \frac{\sum_{i=1}^m -2 \times (x_i - \mu_B)}{m}$$

$$④ \frac{\partial l}{\partial x_i} = \frac{\partial l}{\partial y_i} \times \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_B^2} \times \frac{2 \times (x_i - \mu_B)}{m} + \frac{\partial l}{\partial \mu_B} \times \frac{1}{m}$$

$$⑤ \frac{\partial l}{\partial \gamma} = \sum_{i=1}^m \frac{\partial l}{\partial y_i} \times \hat{x}_i$$

$$⑥ \frac{\partial l}{\partial \beta} = \sum_{i=1}^m \frac{\partial l}{\partial y_i}$$

Backpropagation

We have the following equations by using the loss function l

$$① \frac{\partial l}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \times \gamma$$

$$② \frac{\partial l}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \times (\mathbf{x}_i - \mu_B) \times \left(-\frac{1}{2}\right) \times (\sigma_B^2 + \epsilon)^{-\frac{3}{2}}$$

$$③ \frac{\partial l}{\partial \mu_B} = \left(\sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \times \frac{-1}{\sqrt{\sigma_B^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_B^2} \times \frac{\sum_{i=1}^m -2 \times (\mathbf{x}_i - \mu_B)}{m}$$

$$\textcircled{1} \quad \frac{\partial l}{\partial x_i} = \frac{\partial l}{\partial y_i} \times \frac{1}{\sqrt{\sigma_B^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_B^2} \times \frac{2 \times (\mathbf{x}_i - \mu_B)}{m} + \frac{\partial l}{\partial \mu_B} \times \frac{1}{m}$$

$$\textcircled{2} \quad \frac{\partial l}{\partial \gamma} = \sum_{i=1}^m \frac{\partial l}{\partial y_i} \times \hat{x}_i$$

$$\textcircled{3} \quad \frac{\partial l}{\partial \beta} = \sum_{i=1}^m \frac{\partial l}{\partial y_i}$$

Backpropagation

We have the following equations by using the loss function l

$$① \quad \frac{\partial l}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \times \gamma$$

$$② \quad \frac{\partial l}{\partial \sigma_{\mathcal{B}}^2} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \times (\mathbf{x}_i - \mu_{\mathcal{B}}) \times \left(-\frac{1}{2}\right) \times (\sigma_{\mathcal{B}}^2 + \epsilon)^{-\frac{3}{2}}$$

$$③ \quad \frac{\partial l}{\partial \mu_{\mathcal{B}}} = \left(\sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \times \frac{-1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_{\mathcal{B}}^2} \times \frac{\sum_{i=1}^m -2 \times (\mathbf{x}_i - \mu_{\mathcal{B}})}{m}$$

$$④ \quad \frac{\partial l}{\partial \mathbf{x}_i} = \frac{\partial l}{\partial \hat{x}_i} \times \frac{1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_{\mathcal{B}}^2} \times \frac{2 \times (\mathbf{x}_i - \mu_{\mathcal{B}})}{m} + \frac{\partial l}{\partial \mu_{\mathcal{B}}} \times \frac{1}{m}$$

$$\text{⑤ } \frac{\partial l}{\partial \gamma} = \sum_{i=1}^m \frac{\partial l}{\partial y_i} \times \hat{y}_i$$

$$\text{⑥ } \frac{\partial l}{\partial \beta} = \sum_{i=1}^m \frac{\partial l}{\partial b}$$

Backpropagation

We have the following equations by using the loss function l

$$① \frac{\partial l}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \times \gamma$$

$$② \frac{\partial l}{\partial \sigma_{\mathcal{B}}^2} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \times (\mathbf{x}_i - \mu_{\mathcal{B}}) \times \left(-\frac{1}{2}\right) \times (\sigma_{\mathcal{B}}^2 + \epsilon)^{-\frac{3}{2}}$$

$$③ \frac{\partial l}{\partial \mu_{\mathcal{B}}} = \left(\sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \times \frac{-1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_{\mathcal{B}}^2} \times \frac{\sum_{i=1}^m -2 \times (\mathbf{x}_i - \mu_{\mathcal{B}})}{m}$$

$$④ \frac{\partial l}{\partial \mathbf{x}_i} = \frac{\partial l}{\partial \hat{x}_i} \times \frac{1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_{\mathcal{B}}^2} \times \frac{2 \times (\mathbf{x}_i - \mu_{\mathcal{B}})}{m} + \frac{\partial l}{\partial \mu_{\mathcal{B}}} \times \frac{1}{m}$$

$$⑤ \frac{\partial l}{\partial \gamma} = \sum_{i=1}^m \frac{\partial l}{\partial y_i} \times \hat{\mathbf{x}}_i$$

Backpropagation

We have the following equations by using the loss function l

$$① \frac{\partial l}{\partial \hat{x}_i} = \frac{\partial l}{\partial y_i} \times \gamma$$

$$② \frac{\partial l}{\partial \sigma_{\mathcal{B}}^2} = \sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \times (\mathbf{x}_i - \mu_{\mathcal{B}}) \times \left(-\frac{1}{2}\right) \times (\sigma_{\mathcal{B}}^2 + \epsilon)^{-\frac{3}{2}}$$

$$③ \frac{\partial l}{\partial \mu_{\mathcal{B}}} = \left(\sum_{i=1}^m \frac{\partial l}{\partial \hat{x}_i} \times \frac{-1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \right) + \frac{\partial l}{\partial \sigma_{\mathcal{B}}^2} \times \frac{\sum_{i=1}^m -2 \times (\mathbf{x}_i - \mu_{\mathcal{B}})}{m}$$

$$④ \frac{\partial l}{\partial \mathbf{x}_i} = \frac{\partial l}{\partial \hat{x}_i} \times \frac{1}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \frac{\partial l}{\partial \sigma_{\mathcal{B}}^2} \times \frac{2 \times (\mathbf{x}_i - \mu_{\mathcal{B}})}{m} + \frac{\partial l}{\partial \mu_{\mathcal{B}}} \times \frac{1}{m}$$

$$⑤ \frac{\partial l}{\partial \gamma} = \sum_{i=1}^m \frac{\partial l}{\partial y_i} \times \hat{\mathbf{x}}_i$$

$$⑥ \frac{\partial l}{\partial \beta} = \sum_{i=1}^m \frac{\partial l}{\partial y_i}$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{x^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- 1. $N_{BN}^{tr} = N // \text{Training BN network}$
- 2. for $k = 1..K$ do
 - 3. Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$ to N_{BN}^{tr}
 - 4. Modify each layer in N_{BN}^{tr} with input $x^{(k)}$ to take $y^{(k)}$ instead
 - 5. Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
 - 6. $N_{BN}^{inf} = N_{BN}^{tr} // \text{Inference BN network with frozen parameters}$
 - 7. for $k = 1..K$ do
 - 8. Process multiple training mini-batches B , each of size m , and average over them
 - 9. $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$
 - 10. In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with
 - 11. $y = \frac{\gamma}{\sqrt{Var[x]+\epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x]+\epsilon}} \right]$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{x^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- ① $N_{BN}^{tr} = N$ // Training BN network
- for $k = 1..K$ do
 - Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$ to N_{BN}^{tr}
 - Modify each layer in N_{BN}^{tr} with input $x^{(k)}$ to take $y^{(k)}$ instead
 - Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- for $k = 1..K$ do
 - Process multiple training mini-batches B , each of size m , and average over them
 - $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1}E_B[\sigma_B^2]$
 - In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with
$$y = \frac{\gamma}{\sqrt{Var[x]+\epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x]+\epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{x^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- ① $N_{BN}^{tr} = N$ // Training BN network
- ② for $k = 1 \dots K$ do
 - ➊ Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$ to N_{BN}^{tr}
 - ➋ Modify each layer in N_{BN}^{tr} with input $x^{(k)}$ to take $y^{(k)}$ instead
 - ➌ Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
 - ➍ $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
 - ➎ for $k = 1 \dots K$ do
 - ➏ Process multiple training mini-batches B , each of size m , and average over them
 - ➐ $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$
 - ➑ In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with
 - ➒ $y = \frac{\gamma}{\sqrt{Var[x]+\epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x]+\epsilon}} \right]$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- ➊ $N_{BN}^{tr} = N$ // Training BN network
- ➋ for $k = 1 \dots K$ do
 - ➌ Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
 - ➍ Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $y^{(k)}$ instead
 - ➎ Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
 - ➏ $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
 - ➐ for $k = 1 \dots K$ do
 - ➑ Process multiple training mini-batches B , each of size m , and average over them
 - ➒ $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$
 - ➓ In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with
 - ➔ $y = \frac{\gamma}{\sqrt{Var[x]+\epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x]+\epsilon}} \right]$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- ① $N_{BN}^{tr} = N$ // Training BN network
- ② for $k = 1 \dots K$ do
- ③ Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
- ④ Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $y^{(k)}$ instead

- ⑤ Train N_{BN}^{tr} to optimize the parameters $\{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- ⑥ $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- ⑦ for $k = 1 \dots K$ do
- ⑧ Process multiple training mini-batches B , each of size m , and average over them
- ⑨ $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$
- ⑩ In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with
$$y = \frac{\gamma}{\sqrt{Var[x]+\epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x]+\epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- ➊ $N_{BN}^{tr} = N$ // Training BN network
- ➋ for $k = 1 \dots K$ do
 - ➌ Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
 - ➍ Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $y^{(k)}$ instead
 - ➎ Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- ➏ $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- ➐ for $k = 1 \dots K$ do
 - ➑ Process multiple training mini-batches B , each of size m , and average over them
 - ➒ $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$
 - ➓ In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with
 - ➔ $y = \frac{\gamma}{\sqrt{Var[x]+\epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x]+\epsilon}} \right]$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- ➊ $N_{BN}^{tr} = N$ // Training BN network
- ➋ for $k = 1 \dots K$ do
 - ➌ Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
 - ➍ Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $y^{(k)}$ instead
 - ➎ Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
 - ➏ $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
 - ➐ for $k = 1 \dots K$ do
 - ➑ Process multiple training mini-batches B , each of size m , and average over them
 - ➒ $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$
 - ➓ In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with
 - ➔ $y = \frac{\gamma}{\sqrt{Var[x]+\epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x]+\epsilon}} \right]$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- ① $N_{BN}^{tr} = N$ // Training BN network
- ② for $k = 1 \dots K$ do
- ③ Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
- ④ Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $y^{(k)}$ instead
- ⑤ Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- ⑥ $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- ⑦ for $k = 1 \dots K$ do
 - ⑧ Process multiple training mini-batches B , each of size m , and average over them
 - ⑨ $E[x] = E_B[\mu_B]$ and $Var[x] = \frac{m}{m-1} E_B[\sigma_B^2]$
 - ⑩ In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(x)$ with
$$y = \frac{\gamma}{\sqrt{Var[x]+\epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x]+\epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- ① $N_{BN}^{tr} = N$ // Training BN network
- ② for $k = 1 \dots K$ do
- ③ Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
- ④ Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $y^{(k)}$ instead
- ⑤ Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- ⑥ $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- ⑦ for $k = 1 \dots K$ do
- ⑧ Process multiple training mini-batches \mathcal{B} , each of size m , and average over them

• $E[x] = \bar{x}_B$ [avg] and $Var[x] = \frac{m}{m-1} \bar{s}_B^2$

• In N_{BN}^{tr} , replace the transform $y = BN_{\gamma, \beta}(x)$ with

• $y = \frac{\gamma}{\sqrt{Var[x]+\epsilon}} \times x + \left[\beta - \frac{\gamma E[x]}{\sqrt{Var[x]+\epsilon}} \right]$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- ① $N_{BN}^{tr} = N$ // Training BN network
- ② for $k = 1 \dots K$ do
- ③ Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
- ④ Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $y^{(k)}$ instead
- ⑤ Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- ⑥ $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- ⑦ for $k = 1 \dots K$ do
- ⑧ Process multiple training mini-batches \mathcal{B} , each of size m , and average over them
- ⑨ $E[\mathbf{x}] = E_{\mathcal{B}}[\mu_{\mathcal{B}}]$ and $Var[\mathbf{x}] = \frac{m}{m-1} \mathcal{B} [\sigma_{\mathcal{B}}^2]$

⑩ In N_{BN}^{tr} , replace the transform $y = BN_{\gamma, \beta}(\mathbf{x})$ with

$$y = \frac{1}{\sqrt{Var[\mathbf{x}]+\epsilon}} \times \mathbf{x} + \left[\beta - \frac{\gamma E[\mathbf{x}]}{\sqrt{Var[\mathbf{x}]+\epsilon}} \right]$$

Training Batch Normalization Networks

Input: Network N with trainable parameters Θ ; subset of activations $\{\mathbf{x}^{(k)}\}_{k=1}^K$

Output: Batch-normalized network for inference N_{BN}^{inf}

- ① $N_{BN}^{tr} = N$ // Training BN network
- ② for $k = 1 \dots K$ do
- ③ Add transformation $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(\mathbf{x}^{(k)})$ to N_{BN}^{tr}
- ④ Modify each layer in N_{BN}^{tr} with input $\mathbf{x}^{(k)}$ to take $y^{(k)}$ instead
- ⑤ Train N_{BN}^{tr} to optimize the parameters $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- ⑥ $N_{BN}^{inf} = N_{BN}^{tr}$ // Inference BN network with frozen parameters
- ⑦ for $k = 1 \dots K$ do
- ⑧ Process multiple training mini-batches \mathcal{B} , each of size m , and average over them
- ⑨ $E[\mathbf{x}] = E_{\mathcal{B}}[\mu_{\mathcal{B}}]$ and $Var[\mathbf{x}] = \frac{m}{m-1} \mathcal{B} [\sigma_{\mathcal{B}}^2]$
- ⑩ In N_{BN}^{inf} , replace the transform $y = BN_{\gamma, \beta}(\mathbf{x})$ with
- ⑪
$$\mathbf{y} = \frac{\gamma}{\sqrt{Var[\mathbf{x}]+\epsilon}} \times \mathbf{x} + \left[\beta - \frac{\gamma E[\mathbf{x}]}{\sqrt{Var[\mathbf{x}]+\epsilon}} \right]$$

However

Santurkar et al. [23]

- They found that is not the covariance shift the one affected by it!!!

Batch Normalization

- Batch normalization has been arguably one of the most successful architectural innovations in deep learning.

ReLU vs. Standardized ReLU on CIFAR-10

- on CIFAR-10 with and without BatchNorm

However

Santurkar et al. [23]

- They found that is not the covariance shift the one affected by it!!!

Santurkar et al. recognize that

- Batch normalization has been arguably one of the most successful architectural innovations in deep learning.

They also find that

- on CIFAR-10 with and without BatchNorm

However

Santurkar et al. [23]

- They found that is not the covariance shift the one affected by it!!!

Santurkar et al. recognize that

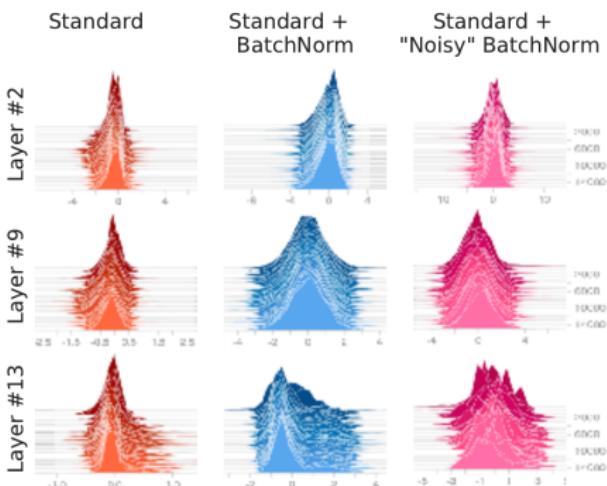
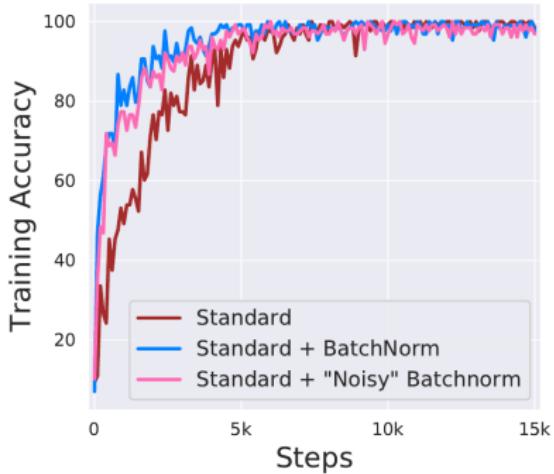
- Batch normalization has been arguably one of the most successful architectural innovations in deep learning.

They used a standard Very deep convolutional network

- on CIFAR-10 with and without BatchNorm

They found something quite interesting

The following facts



Actually Batch Normalization

It does not do anything to the Internal Covariate Shift

- Actually smooth the optimization manifold
 - ▶ It is not the only way to achieve it!!!

- "This suggests that the positive impact of BatchNorm on training might be somewhat serendipitous."

Actually Batch Normalization

It does not do anything to the Internal Covariate Shift

- Actually smooth the optimization manifold
 - ▶ It is not the only way to achieve it!!!

They suggest that

- “This suggests that the positive impact of BatchNorm on training might be somewhat serendipitous.”

They actually have a connected result

To the analysis of gradient clipping!!!

- They are the same group

Lemma: The effect of BN differs on the Lipschitzness of the loss

- For a BatchNorm network with loss $\hat{\mathcal{L}}$ and an identical non-BN network with (identical) loss \mathcal{L} ,

$$\left\| \nabla_{y_j} \hat{\mathcal{L}} \right\|^2 \leq \frac{\gamma^2}{\sigma_j^2} \left[\left\| \nabla_{y_j} \mathcal{L} \right\|^2 + \frac{1}{m} \langle \mathbf{1}, \nabla_{y_j} \mathcal{L} \rangle^2 - \frac{1}{\sqrt{m}} \langle \nabla_{y_j} \mathcal{L}, \hat{y}_j \rangle^2 \right]$$

They actually have a connected result

To the analysis of gradient clipping!!!

- They are the same group

Theorem (The effect of BatchNorm on the Lipschitzness of the loss)

- For a BatchNorm network with loss $\hat{\mathcal{L}}$ and an identical non-BN network with (identical) loss \mathcal{L} ,

$$\left\| \nabla_{y_j} \hat{\mathcal{L}} \right\|^2 \leq \frac{\gamma^2}{\sigma_j^2} \left[\left\| \nabla_{y_j} \mathcal{L} \right\|^2 - \frac{1}{m} \langle \mathbf{1}, \nabla_{y_j} \mathcal{L} \rangle^2 - \frac{1}{\sqrt{m}} \langle \nabla_{y_j} \mathcal{L}, \hat{\mathbf{y}}_j \rangle^2 \right]$$

Outline

1 Introduction

- Limitations of Shallow Architectures
- Highly-varying functions
- Local vs Non-Local Generalization
- From Simpler Features to More Complex Features

2 Deep Forward Architectures

- Introduction
- Convolutional Neural Networks
 - Image Processing
- Auto Encoders
- Generative Adversarial Networks
 - An Example, Boltzmann Machines
- There Are Many More

3 Problems and Possible Solutions with Deeper Architectures

- The Degradation Problem
- The Vanishing and Exploding Gradients
- Reasoning Iteratively
- Fixed Points
- Stabilizing the Network
 - Gradient Clipping
 - Normalizing your Data
 - Normalization Layer AKA Batch Normalization
- Conclusions

We have seen many concepts

Deep Forward Networks

- Although a simple idea

The main idea is very simple:

- Basically... From Lower Complexity Features toward more complex more informative!!!

In conclusion

- Deep Forward Networks look to have more expressibility than shallow learners.

We have seen many concepts

Deep Forward Networks

- Although a simple idea

They represent a rich field of study

- Basically... From Lower Complexity Features toward more complex more informative!!!

- Deep Forward Networks look to have more expressibility than shallow learners.

We have seen many concepts

Deep Forward Networks

- Although a simple idea

They represent a rich field of study

- Basically... From Lower Complexity Features toward more complex more informative!!!

In conclusion

- Deep Forward Networks look to have more expressibility than shallow learners.

-  C. E. Shannon, "A symbolic analysis of relay and switching circuits," *Electrical Engineering*, vol. 57, no. 12, pp. 713–723, 1938.
-  E. Mendelson, *Introduction to mathematical logic*. Chapman and Hall/CRC, 2009.
-  J. Hastad, "Almost optimal lower bounds for small depth circuits," in *Proceedings of the eighteenth annual ACM symposium on Theory of computing*, pp. 6–20, Citeseer, 1986.
-  J. Håstad and M. Goldmann, "On the power of small-depth threshold circuits," *Computational Complexity*, vol. 1, no. 2, pp. 113–129, 1991.
-  Y. Bengio *et al.*, "Learning deep architectures for ai," *Foundations and trends® in Machine Learning*, vol. 2, no. 1, pp. 1–127, 2009.
-  M. Gönen and E. Alpaydın, "Multiple kernel learning algorithms," *Journal of machine learning research*, vol. 12, no. Jul, pp. 2211–2268, 2011.

- ❑ G. R. G. Lanckriet, N. Cristianini, P. Bartlett, L. E. Ghaoui, and M. I. Jordan, “Learning the kernel matrix with semidefinite programming,” *J. Mach. Learn. Res.*, vol. 5, pp. 27–72, Dec. 2004.
- ❑ Y. Bengio, O. Delalleau, and N. L. Roux, “The curse of highly variable functions for local kernel machines,” in *Advances in neural information processing systems*, pp. 107–114, 2006.
- ❑ Y. Bengio, Y. LeCun, *et al.*, “Scaling learning algorithms towards ai,” *Large-scale kernel machines*, vol. 34, no. 5, pp. 1–41, 2007.
- ❑ K. Fukushima, “Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position,” *Biological cybernetics*, vol. 36, no. 4, pp. 193–202, 1980.
- ❑ J. J. Weng, N. Ahuja, and T. S. Huang, “Learning recognition and segmentation of 3-d objects from 2-d images,” in *1993 (4th) International Conference on Computer Vision*, pp. 121–128, IEEE, 1993.

-  J. J. Weng, N. Ahuja, and T. S. Huang, "Learning recognition and segmentation using the cresceptron," *International Journal of Computer Vision*, vol. 25, no. 2, pp. 109–143, 1997.
-  Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural computation*, vol. 1, no. 4, pp. 541–551, 1989.
-  T. Wiatowski and H. Bölcskei, "A mathematical theory of deep convolutional neural networks for feature extraction," *IEEE Transactions on Information Theory*, vol. 64, no. 3, pp. 1845–1866, 2017.
-  Z. Zhang, "Derivation of backpropagation in convolutional neural network (cnn)," *University of Tennessee, Knoxville, TN*, 2016.

-  X. Peng, H. Cao, and P. Natarajan, "Using convolutional encoder-decoder for document image binarization," in *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, vol. 1, pp. 708–713, IEEE, 2017.
-  P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, and G. Cottrell, "Understanding convolution for semantic segmentation," in *2018 IEEE winter conference on applications of computer vision (WACV)*, pp. 1451–1460, IEEE, 2018.
-  V. Podlozhnyuk, "Image convolution with cuda," *NVIDIA Corporation white paper*, June, vol. 2097, no. 3, 2007.
-  X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 315–323, 2011.
-  I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. The MIT Press, 2016.

- K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” *arXiv preprint arXiv:1502.03167*, 2015.
- S. Santurkar, D. Tsipras, A. Ilyas, and A. Madry, “How does batch normalization help optimization?,” in *Advances in Neural Information Processing Systems*, pp. 2483–2493, 2018.
- C. Gulcehre, M. Moczulski, M. Denil, and Y. Bengio, “Noisy activation functions,” in *International conference on machine learning*, pp. 3059–3068, 2016.
- S. Sharma, “Activation functions in neural networks,” *Towards Data Science*, vol. 6, 2017.

-  L. McInnes, J. Healy, and J. Melville, "Umap: Uniform manifold approximation and projection for dimension reduction," *arXiv preprint arXiv:1802.03426*, 2018.
-  Y. Li, S. Liu, J. Yang, and M.-H. Yang, "Generative face completion," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3911–3919, 2017.
-  J. Pennington, S. S. Schoenholz, and S. Ganguli, "The emergence of spectral universality in deep networks," *arXiv preprint arXiv:1802.09979*, 2018.
-  J. Zhang, T. He, S. Sra, and A. Jadbabaie, "Analysis of gradient clipping and adaptive scaling with a relaxed smoothness condition," *arXiv preprint arXiv:1905.11881*, 2019.
-  Y. Bengio and Y. Le Cun, "Word normalization for on-line handwritten word recognition," in *International Conference on Pattern Recognition*, pp. 409–409, IEEE COMPUTER SOCIETY PRESS, 1994.