

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

THE COOPER UNION
ALBERT NERKEN SCHOOL OF ENGINEERING

A Deep Partitioned Autoencoder
for De-Noising Live Audio

by
Ethan Lusterman

A thesis submitted in partial fulfillment
of the requirements for the degree of
Master of Engineering

September 2016

Professor Sam Keene, Advisor

054 THE COOPER UNION FOR THE
055
056 ADVANCEMENT OF SCIENCE AND ART
057
058
059
060

061 ALBERT NERKEN SCHOOL OF ENGINEERING
062
063
064
065
066
067
068
069

070 This thesis was prepared under the direction of the Can-
071 didate's Thesis Advisor and has received approval. It was
072 submitted to the Dean of the School of Engineering and
073 the full Faculty, and was approved as partial fulfillment of
074 the requirements for the degree of Master of Engineering.
075
076
077
078
079
080
081
082
083
084
085
086

087 _____
088 Dean, School of Engineering Date
089
090
091
092
093
094

095 _____
096 Prof. Sam Keene, Thesis Advisor Date
097
098
099
100
101
102
103
104
105
106
107

108	Ack
109	
110	
111	
112	
113	
114	
115	
116	
117	
118	
119	
120	
121	
122	
123	
124	
125	
126	
127	
128	
129	
130	
131	
132	
133	
134	
135	
136	
137	
138	
139	
140	
141	
142	
143	
144	
145	
146	
147	
148	
149	
150	
151	
152	
153	
154	
155	
156	
157	
158	
159	
160	
161	

162	Abstract
163	
164	
165	
166	
167	
168	
169	
170	
171	
172	
173	
174	
175	
176	
177	
178	
179	
180	
181	
182	
183	
184	
185	
186	
187	
188	
189	
190	
191	
192	
193	
194	
195	
196	
197	
198	
199	
200	
201	
202	
203	
204	
205	
206	
207	
208	
209	
210	
211	
212	
213	
214	
215	

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

Contents

1	Introduction	1
2	Background	2
2.1	Machine Learning at Large	2
2.1.1	Differentiable parametric modeling	3
2.1.2	Basis expansions	3
2.1.3	Kernel methods	5
2.1.4	Curse of dimensionality	5
2.1.5	Over-fitting	5
2.2	Neural Networks	5
2.2.1	Neural networks as computational graphs	5
2.2.2	Fully connected neural networks	6
2.2.3	Convolutional neural networks	9
2.2.4	Convolutions with holes	9
2.2.5	Tranposed convolutions	9
2.2.6	Additional techniques	9
2.2.6.1	Dropout	9
2.2.6.2	Batch normalization	10
2.2.6.3	Rectified linear units	10
2.2.6.4	Exponential linear units	10
2.2.6.5	Residual networks	10
2.2.6.6	Minibatch training	10
2.2.6.7	Intelligent parameter initialization	11
2.2.6.8	ADAM optimizer	12
2.2.6.9	Style loss	12
2.3	Generative Models	12
2.3.1	Generative models at large	12
2.3.2	Generative adversarial networks	12
2.3.3	Deep convolutional generative adversarial newtorks . .	12

270	2.3.4	Additional techniques for training generative adversarial	
271		networks	12
272			
273			
274	3	System Description	12
275			
276	3.1	Energy distance matching	12
277			
278	3.2	Kernel based moment matching	12
279			
280	3.3	Convolutions with holes	12
281			
282	3.4	Style Loss term	12
283			
284			
285			
286			
287			
288			
289			
290			
291			
292			
293			
294			
295			
296			
297			
298			
299			
300			
301			
302			
303			
304			
305			
306			
307			
308			
309			
310			
311			
312			
313			
314			
315			
316			
317			
318			
319			
320			
321			
322			
323			

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

List of Figures

1	Linear regression of a noisy sine-wave	4
2	Linear regression of a noisy sine-wave with basis expansion . .	5
3	Polynomial Fit	5
4	Feed-forward fully-connected network	9

378	Table of Nomenclature
379	
380	
381	
382	
383	
384	
385	
386	
387	
388	
389	
390	
391	
392	
393	
394	
395	
396	
397	
398	
399	
400	
401	
402	
403	
404	
405	
406	
407	
408	
409	
410	
411	
412	
413	
414	
415	
416	
417	
418	
419	
420	
421	
422	
423	
424	
425	
426	
427	
428	
429	
430	
431	

1 Introduction

Advances in smartphone technology have led to smaller devices with more powerful audio hardware, allowing for common consumers to make higher quality recordings. However, recorded speech and music are subject to noisy conditions, often hampering intelligibility and listenability. The goal of denoising audio recordings is to improve intelligibility and perceived quality. A variety of applications of audio denoising exist, including listening to a recording of a band or an artist’s live performance in a noisy crowd, or listening to a recorded conversation or speech under noisy conditions.

A common technique for denoising involves the use of deep neural networks (DNN). [PARIS] Advances in parallel graphics processing units (GPU) and in machine learning algorithms have allowed for training deeper networks faster, utilizing more hidden layers with more neurons.

Prior work in denoising audio has involved access to noise-free training data. Since common consumers do not often have access to clean audio, we seek to denoise without the use of clean audio.

In this thesis, we compare several neural network architectures and problem scenarios, ranging from data input types, level of noise, depth of network, training objectives, and more. In Chapter 2, we present background information on machine learning and neural networks as well as prior work in audio denoising. In Chapter 3, we detail all considered network architectures. In Chapter 4, we compare results from different data inputs, levels of noise, network architectures, and training objectives and discuss methods of evaluation. Finally, we make conclusions and recommendations for future work in Chapter 5.

2 Background

2.1 Machine Learning at Large

Machine learning in general encompasses all tasks where we utilize data to make some kind of decision. Classical examples include handwritten digit recognition, and flower species classification based on sepal width and petal length. As rough mathematical example, consider a sample of N pairs of examples drawn from larger population \mathcal{P} :

$$\mathcal{D} = \{(x_i, y_i) : i = 1, 2, \dots, N\} \quad (1)$$

$$\mathcal{D} \subset \mathcal{P} \quad (2)$$

Our machine learning task may be to find some function $f(x)$ that estimates y for any (x, y) sample drawn from the population at large. To learn the function we might setup a loss function $l(y, \hat{y})$ to quantify how well our model performs on our sample \mathcal{D} . As a general rule, the loss function satisfies the following limit:

$$\lim_{\hat{y}_i \rightarrow y_i} l(y_i, \hat{y}_i) = 0 \quad (3)$$

In other words our goal is to find some $f(x)$ that minimizes the loss function $l(y, \hat{y})$ for all (x_i, y_i) pairs in \mathcal{D} . This type of task is referred to as supervised learning, because we have a domain of inputs and a specified codomain of targets; there are other applications with unknown targets, these tasks are referred to as unsupervised learning tasks. With a few extra constraints we are able use a number of techniques to accomplish this supervised learning goal.

2.1.1 Differentiable parametric modeling

Let us restrict our study of functions f specifically to functions that are differentiable and parametric in form, so from now on we refer to $f(x | \theta)$ where θ refers to the function's set of parameters. Using these new terms our goal is to find some $\hat{\theta}$ that satisfies the following:

$$\hat{\theta} = \arg \min_{\theta} l(y, f(x | \theta)), \forall x, y \quad (4)$$

For any problem of this type we would like to calculate gradient of l with respect to each θ_i , set them all to zero and solve for each θ_i :

$$\frac{\partial l(x, f(x | \theta))}{\partial \theta_i} = 0, \forall i \quad (5)$$

However this is not always possible analytically because without convexity constraints we are not able to guarantee the existence of a unique minimum in the loss function.

In lieu of an analytic solution we can use a numerical optimization approach, and jointly optimize the θ_i 's via a gradient descent algorithm. For example we could use this update rule, which we refer to as steepest descent:

$$\theta_i \leftarrow \theta_i - \eta \frac{\partial l(x, f(x | \theta))}{\partial \theta_i} \quad (6)$$

where η is a the learning rate parameter, which controls how large our steps are towards the minimum.

2.1.2 Basis expansions

Let us consider our first supervised learning example. Consider a noisy sine-wave consisting of k samples:

$$\mathbf{y} = \sin \mathbf{x} + \mathbf{n}, \mathbf{y}, \mathbf{x}, \mathbf{n} \in \mathbb{R}^k \quad (7)$$

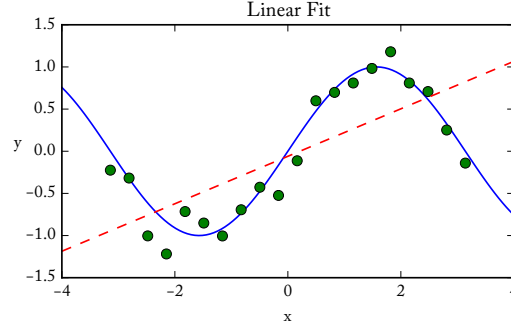


Figure 1: Linear regression of a noisy sine-wave.

We would like to find some $f(\mathbf{x} \mid \theta)$ that produces an estimate $\hat{\mathbf{y}}$. Let us consider the following functional form for f :

$$f(\mathbf{x}) = m\mathbf{x} + b, \quad m, b \in \mathbb{R} \quad (8)$$

Or we can rewrite this as a matrix multiplication by padding \mathbf{x} in with a ones vector to create a $2 \times k$ matrix we refer to as Φ , and combining m and b into a vector \mathbf{w} :

$$\hat{\mathbf{y}} = \mathbf{w}^T \Phi \quad (9)$$

We are able to compute \hat{m} and \hat{b} analytically, with a mean squared error loss function:

$$l(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{2} \sum_{i=1}^k (\hat{y}_i - y_i)^2 \quad (10)$$

Now let's substitute our model in for $\hat{\mathbf{y}}$ and set the gradient to zero:

$$\frac{\partial l(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{w}} \left(\frac{1}{2} \sum_{i=1}^k (\mathbf{w}^T \Phi - y_i)^2 \right) = 0 \quad (11)$$

Carrying out the derivative we see:

$$\frac{1}{2} \sum_{i=1}^k (\mathbf{w}^T \Phi - \mathbf{y}) \Phi^T = 0 \quad (12)$$

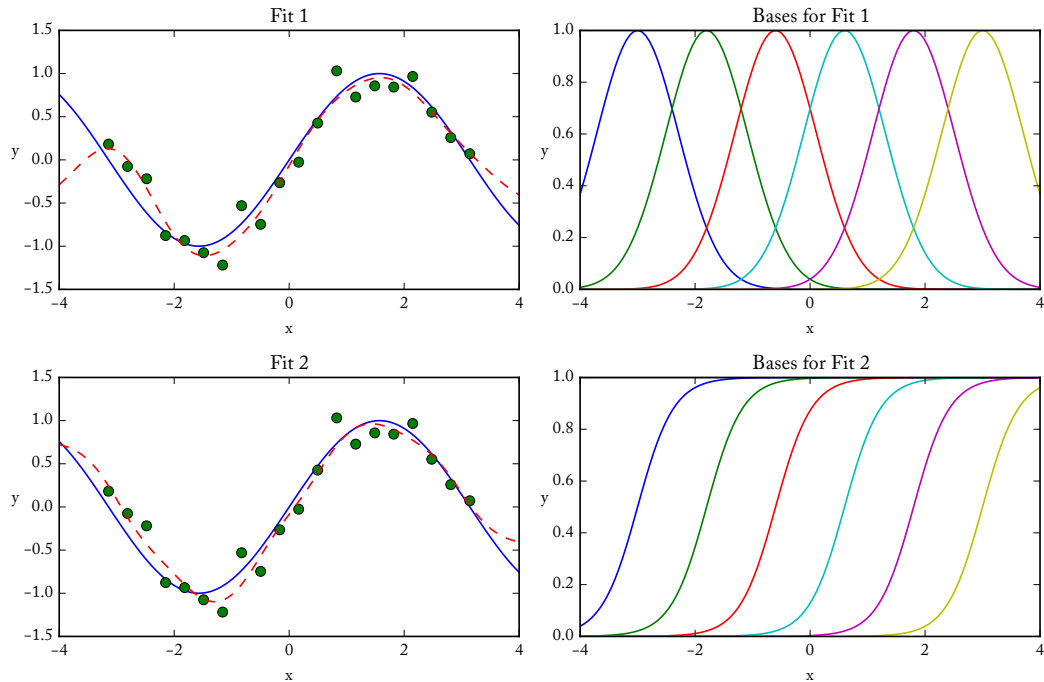


Figure 2: Linear regression of a noisy sine-wave with basis expansion

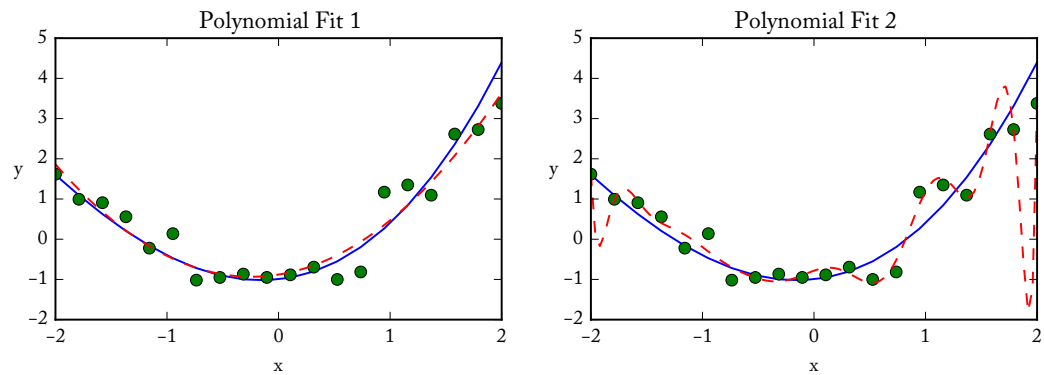


Figure 3: Polynomial Fit

2.1.3 Kernel methods

2.1.4 Curse of dimensionality

2.1.5 Over-fitting

2.2 Neural Networks

2.2.1 Neural networks as computational graphs

A neural network generally is a computational graph $G = (V, E)$ where the vertices V correspond to nodes of computation, and the edges E correspond to

data flow paths connecting the computational nodes. We limit our discussion to feed-forward or directed acyclic graphs, in other words graphs where given a starting vertex v we are unable to follow the directed edges away from it and return to v . With this limitation in place we are able to focus on stateless graphs, which treat input data examples independently of one another.

With the additional limitation that all computational nodes in the graph perform differentiable operations, and that some nodes are parametric, we are able to use automatic differentiation and a gradient descent like optimization algorithm to optimize the parameters of the graph against a differentiable loss function [?]. Automatic differentiation is an alternative to numeric and symbolic differentiation, that is efficient, accurate to machine precision, and scalable to arbitrary graphs consisting of elementary operations. The core tenet of automatic differentiation is that since graph vertices are differentiable we can perform the chain rule at each node to build up gradients of the full graph. Several machine learning frameworks, such as Theano and TensorFlow implement automatic differentiation, allowing us to specify the functional form of the graph and getting gradients at minimal cost [?,?].

2.2.2 Fully connected neural networks

While we have introduced the idea that neural networks can be made up of arbitrary computational nodes, there are several common types of nodes that are used to build up a toolkit of functional forms at our disposal. First, let us consider a simple inner product with signature:

$$\langle \cdot, \cdot \rangle: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R} \quad (13)$$

And functional form:

$$\langle \mathbf{w}, \mathbf{x} \rangle = \sum_{i=1}^n w_i x_i \quad (14)$$

We can interpret this inner product as a single node in our graph, with parameters \mathbf{w} and input vector \mathbf{x} . While this functional form is nice, in that it's linear, and easy to define, it's limited by the fact that it brings the rich space $\mathbb{R}^n \times \mathbb{R}^n$ down to \mathbb{R} which has limited representational power for problems of interest in machine learning. Instead we'll follow a common design pattern, and create an array of inner product nodes at a constant number of hops from the source node in the graph, or as we'll refer to from here on, at a constant depth. In this case the signature would be:

$$g: \mathbb{R}^{n \times m} \times \mathbb{R}^n \rightarrow \mathbb{R}^m \quad (15)$$

And functional form:

$$g(W, \mathbf{x})_j = \sum_{i=1}^n w_{ij} x_i \quad (16)$$

which we recognize as ordinary matrix multiplication $W\mathbf{x}$, where $W \in \mathbb{R}^{n \times m}$. We call this collection of nodes at a common depth a layer. This type of matrix multiplication layer is generally referred to as a linear, dense, or fully connected layer.

Layers can be connected together with the goal of creating a more powerful model. If we feed one dense layer into another directly we have a functional form:

$$\mathbf{y} = W_n \cdots W_3 W_2 W_1 \mathbf{x} \quad (17)$$

where the W_i 's are all compatible. However this approach is actually equivalent to a single matrix multiply $W\mathbf{x}$ where W is the matrix product of the W_i 's, so we did not actually achieve our goal of higher representational power. Our model is still merely linear.

To address this issue we introduce the addition of a non-linear transformation $f(\cdot)$ (known as an activation function) at the output of the matrix

multiply. So instead our functional form would be:

$$\mathbf{y} = f(W_n \cdots f(W_3 f(W_2 f(W_1 \mathbf{x}))) \cdots) \quad (18)$$

As long as $f(\cdot)$ is differentiable we are able to train our now highly non-linear model with our methods of automatic differentiation and gradient based optimization.

Finally we introduce a bias \mathbf{b} so that each node in a layer is not constrained to intercept zero. Therefore our full functional form for a single fully connected layer with a non-linearity is:

$$\mathbf{y} = f(W\mathbf{x} + \mathbf{b}) \quad (19)$$

with $W \in \mathbb{R}^{n \times m}$ and $\mathbf{b} \in \mathbb{R}^m$.

A traditional activation function is the sigmoid function:

$$S(t) = \frac{1}{1 + e^{-t}} \quad (20)$$

which maps $t \in \mathbb{R}$ (i.e the interval $[-\infty, \infty]$) to $S(t)$ on the interval $[0, 1]$.

This can allow the interpretation of $S(t)$ as a probability. Using this activation function Cybenko proved the universal approximation theorem [?] which shows that neural networks of the form in Figure 4 can approximate any function with appropriate domain and codomains given a large of enough number nodes in the middle layer, generally referred to as a hidden layer.

Figure 4 shows a feed-forward network with one hidden-layer. The outputs of the hidden layer are referred to as latent variables and are a new type of representation for the input data \mathbf{x} . With each hidden-layer added to the network increasingly complex data representations may be available, but the networks will also tend to over-fit their training data as the number of parameters increase. In order to combat over-fitting, and increase the generalization performance of the network new types of architectures have been proposed.

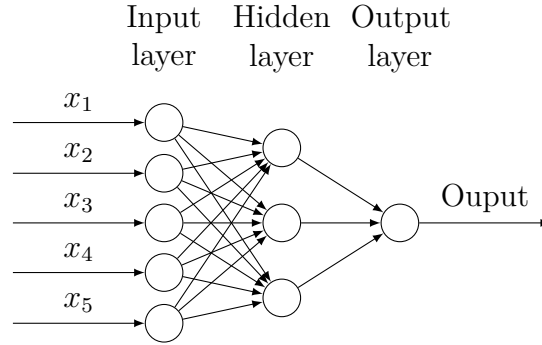


Figure 4: A single hidden layer feed-forward neural network. Each directed connection has a scaling factor w_{ij} associated with it; these entries make up the W matrix. Each circle, now referred to as a neuron, performs the sum and activation functions corresponding to the matrix multiply, and mapping by $f(\cdot)$ in Equation 19. Layers in the middle of the network are referred to as hidden layers because they are not directly observable — instead these are latent variables.

2.2.3 Convolutional neural networks

2.2.4 Convolutions with holes

2.2.5 Tranposed convolutions

2.2.6 Additional techniques

2.2.6.1 Dropout

Various teams have demonstrated that in deep networks neurons can learn complex co-adaptation schemes; that is deep neurons respond to “mistakes” in shallower neurons. In order to prevent co-adaptation, so that each neuron learns a meaningful representation, the dropout scheme has been proposed [?, ?, ?]. Dropout is a masking process. In order to apply dropout to a layer during training, a binary mask is applied across neurons. The mask is determined by sampling a Bernoulli distribution with a parameter p corresponding to the probability that a neuron will be masked. When a neuron is masked its output

is considered fixed at zero. When the network is used for evaluation no masks are applied and all neurons are connected.

2.2.6.2 Batch normalization

2.2.6.3 Rectified linear units

Various activations have been proposed on the basis of similarity to the potential activations in actual biological neurons in human eyes. Recently the rectified linear unit (ReLU) has demonstrated significant performance increases for network generalization and increased training speed [?, ?]. The ReLU activation is defined as $\max(0, x)$. In other words a ReLU activation forwards along any positive inputs and sets and negative inputs to zero.

2.2.6.4 Exponential linear units

2.2.6.5 Residual networks

2.2.6.6 Minibatch training

Traditionally there were two different approaches to optimizing neural network parameters, the online approach and the batch approach. In the online approach, the parameters of the network are updated after each exposure to a training example and gradient calculation. In the batch approach, the network is exposed to all of the training examples, the gradients are accumulated and then the network’s parameters are updated. This was long believed to be the better approach, as the accumulated and averaged gradient was more likely to be an estimate of the “true” gradient of the network towards an optimized solution. It was later shown that, while the batch mode may give a better estimate of the gradient, due to the noise and its stochastic nature, the online approach actually leads to a faster convergence time, in terms of number of

examples [?]. This is because in the stochastic approach the optimizer is less likely to get stuck in local minima. An alternative approach that recently has become popular is the mini- batch approach. In this approach some number of examples, say N , are exposed to the network, the gradients are accumulated, and the parameters are updated. In this case N is much less than the total number of training examples available. This approach, with serial computational resources really only represents a decrease in training speed, because it is fairly similar to the batch approach. The reason this approach has become popular lately is that with the parallel resources afforded by modern high performance graphics processing units (GPUs) the increase in example-wise training time becomes a decrease in wall-time to train.

2.2.6.7 Intelligent parameter initialization

Kaiming et al. have demonstrated an improved parameter initialization technique specifically designed for neurons with ReLU activations [?]. This approach derives a method that enables extremely deep models comprised of ReLU neurons to converge rather than stall, as they would with other initialization schemes. The initial parameters for the convolutional layers are drawn from a zero mean normal distribution with a standard deviation of $\sqrt{2/n_l}$. Where $n_l = k^2c$. This corresponds to the number of connections in the response for $k \times k$ kernels processing c input channels.

1026	2.2.6.8	ADAM optimizer
1027		
1028	2.2.6.9	Style loss
1029		
1030		
1031	2.3	Generative Models
1032		
1033		
1034	2.3.1	Generative models at large
1035		
1036	2.3.2	Generative adversarial networks
1037		
1038	2.3.3	Deep convolutional generative adversarial networks
1039		
1040	2.3.4	Additional techniques for training generative adversarial networks
1041		
1042		
1043		
1044	3	System Description
1045		
1046		
1047		
1048	3.1	Energy distance matching
1049		
1050	3.2	Kernel based moment matching
1051		
1052		
1053	3.3	Convolutions with holes
1054		
1055		
1056	3.4	Style Loss term
1057		
1058		
1059		
1060		
1061		
1062		
1063		
1064		
1065		
1066		
1067		
1068		
1069		
1070		
1071		
1072		
1073		
1074		
1075		
1076		
1077		
1078		
1079		