

Deep Networks for Equalization in Communications

by

Laura Brink

A thesis submitted in partial satisfaction of the
requirements for the degree of
Master's of Science

in

Electrical Engineering and Computer Science

in the

Graduate Division

of the

University of California, Berkeley

Committee in charge:

Professor Anant Sahai, Chair
Professor John Wawrzynek

Fall 2018

The thesis of Laura Brink, titled Deep Networks for Equalization in Communications, is approved:

| | | | |
|-------|-------|------|-------|
| Chair | _____ | Date | _____ |
| | _____ | Date | _____ |
| | _____ | Date | _____ |

University of California, Berkeley

Deep Networks for Equalization in Communications

Copyright 2018
by
Laura Brink

Abstract

Deep Networks for Equalization in Communications

by

Laura Brink

Master's of Science in Electrical Engineering and Computer Science

University of California, Berkeley

Professor Anant Sahai, Chair

We apply the techniques from meta-learning and machine learning to the communications domain. Specifically, we explore how neural networks can learn to equalize new channel environments without training on them and how neural networks can learn to estimate and correct carrier frequency offset for new rates of rotation without training on them. We show that deep neural networks can learn to learn to estimate channel taps for two tap channels. We also explore how deep recursive neural networks learn to learn to equalize for any given channel. We demonstrate that neural networks can learn to learn to estimate and correct carrier frequency offset for new rates of rotation. Crucially, we do all of this without using backpropagation to re-train the networks for each new set of environmental conditions.

Contents

| | |
|---|-----------|
| Contents | i |
| List of Figures | ii |
| 1 Introduction | 1 |
| 1.1 Motivation | 1 |
| Context | 2 |
| 1.2 Background | 2 |
| Inter-symbol Interference and Equalization | 3 |
| Carrier Frequency Offset and Correction | 6 |
| 1.3 Related Works | 8 |
| 1.4 Contributions of this Report | 11 |
| 2 Deep Networks for Equalization | 12 |
| 2.1 Channel Estimation | 12 |
| 2.2 Channel Equalization | 13 |
| Learning an inverse | 14 |
| Learning to multiply two inputs | 14 |
| RNN for Channel Equalization | 15 |
| 2.3 Future Work | 19 |
| 3 Deep Networks for Carrier Frequency Offset | 20 |
| 3.1 Recurrent Neural Network Follows a Circle | 20 |
| Single Rate | 20 |
| Different Rates | 21 |
| 3.2 Deep Network Carrier Frequency Offset Estimation and Correction | 22 |
| 3.3 Future Work | 24 |
| 4 Conclusion | 26 |
| Bibliography | 27 |

List of Figures

| | | |
|-----|---|----|
| 1.1 | The effects of a two tap channel on the QPSK constellation. | 4 |
| 1.2 | The effects of a carrier frequency offset on the QPSK constellation. | 7 |
| 1.3 | The effects of a two tap channel and a carrier frequency offset on the QPSK constellation. | 9 |
| 1.4 | The impacts of multipath channels and initial phase offsets on bit error performance [17]. Note, the network is not re-trained for each new environment. As the number of paths and carrier frequency offset increase, the performance of the CNN decreases. | 10 |
| 2.1 | Comparison of a neural network, a K-Nearest-Neighbors, and a least-squares channel estimators for two tap, real channels. | 13 |
| 2.2 | Topographical surface representation of the division function; $z = \frac{x}{y}$ | 15 |
| 2.3 | Neural network division loss with respect to the lower bound of y | 16 |
| 2.4 | Topographical surface representation of the multiplication function; $z = xy$ | 17 |
| 2.5 | Neural network multiplication loss with respect to the upper bound of y | 18 |
| 2.6 | Comparison of RNN and MMSE equalizers loss over SNR. | 18 |
| 2.7 | What two tap channels does the equalizer get wrong? | 19 |
| 3.1 | Linear neural network: follow a circle for a constant CFO rate, $\omega = 0.02$. The figure on the left shows the original data of length 100. The figure in the middle shows the estimated data from the neural network given the starting point from the left figure. The figure on the right shows the log of the test cost of the neural network versus traing epoch. | 21 |
| 3.2 | Nonlinear neural network: follow a circle for different CFO rates. The figure on the left shows the original data of length 100 for $\omega = 0.01063284$. The figure in the middle shows the estimated data from the neural network given the starting point and rate of rotation from the left figure. The figure on the right shows the log of the test cost of the neural network versus traing epoch. | 22 |
| 3.3 | The log BER of received signals passed through a neural network CFO estimator with a rotation matrix CFO correction and classic demodulator. The log BER of the same signals passed through just the classic demodulator. | 23 |

| | | |
|-----|---|----|
| 3.4 | Nonlinear neural network estimating and correcting CFO for different ω . In this particular test example, $\omega = 0.00969545$ and $\text{SNR} = \infty$. The figure on the left shows the constellation of the 100 original data symbols. The figure in the middle shows the constellation of the original data symbols after CFO. The figure on the right shows the constellation of the data symbols after the neural network CFO correction. | 24 |
|-----|---|----|

Chapter 1

Introduction

Signal processing has long relied on well-defined, structured processes and protocols to function. However, in order to move towards more robust and adaptive systems, we will need to overhaul these tightly structured processes. We must design new robust and adaptive communications systems.

From an academic perspective, the rise of machine learning tools and processing has allowed us to tackle problems we have not yet been able to, like in image processing. However, we still do not fully understand the possibilities or the limitations of this technology. Most applications for machine learning techniques lack good baselines to compare performance. Applications such as image processing [13], videogame control [16], playing the game Go, [24], and natural language processing [2], are all missing good baselines to compare against. In order to study the limitations of machine learning techniques, we must apply them to spaces that we have studied extensively, like communications, and compare them to the well-known baselines that are already very good.

Most communications systems have three main processes at the receiver; equalization, demodulation, and error-correction. While we will need to design robust forms of all of these processes, we will focus on equalization for the remainder of this report.

1.1 Motivation

Meta-learning is the idea that algorithms need to be able to ‘learn to learn’ in order to generalize to different applications. For example, if a robot is trained to pick up coffee mugs, then we want that robot to be able to quickly learn how to pick up water bottles without having to re-train the robot. Meta-learning is inspired by humans ability to generalize how to learn [14]. Additionally, researchers in neurology have studied how the brain synapses change over time, suggesting that our algorithms will have to change over time to continue learning [1]. Our brains know how to update the synapses as we learn new things. This is like updating the weights of a neural network for new data points. We refer the reader to [15] for a survey of meta-learning technologies and to [8] for recent developments in meta-learning

algorithms.

In this report, we want to understand if our neural network based communication processes can learn to learn. As in, can they handle new environments with new data sequences without having to re-train for the new variables. The communications application is arguably an easier application of meta-learning than others have been exploring, like image classification [10]. In order to better understand the possibilities and limitations of meta-learning, we need to apply it to something simpler, like communications. In this report, we will explore whether or not we can learn to learn to communicate.

Context

A team at UC Berkeley and Nokia Bell Labs participated in the DARPA Spectrum Collaboration Challenge, a competition to build collaborative, intelligent radios¹. The challenge consists of putting multiple radio networks in the same space and bandwidth. Without having been codesigned, the networks must collaborate and share all of the resources. The environment was open to all of the teams so there are no rigid bands like we use today.

One of the potential scenarios for this challenge was that networks might have stranded radios where no radio in their own network can reach them. In this case, that network would need to rely on another network to relay messages to and from the stranded node. The challenge is that the relay network must forward along a signal that it cannot fully decode itself. This requires the network to equalize a signal and correct for CFO and then send it without using demodulation or error-correction. This challenge is one of the practical applications for robust and adaptive processes for equalization.

1.2 Background

While we will not be designing new modulation or demodulation processes in this report, we must introduce some terminology. A constellation is a mapping of bits to complex symbols. A classic quadrature phase shift keying (QPSK) modulator maps pairs of bits to complex symbols. For the remainder of the report, when we use a QPSK constellation, we use the following mapping.

| Bits | Complex Symbol |
|------|----------------|
| 00 | $-1 - 1j$ |
| 10 | $1 - 1j$ |
| 01 | $-1 + 1j$ |
| 11 | $1 + 1j$ |

We will also use signal to noise ratio (SNR) to evaluate our models. SNR is the ratio of the signal power to the noise power. If the SNR is greater than one, then this means

¹I was heavily involved with this project and this report grew out of my work on this challenge.

that there is more signal than noise. If the SNR is infinity, this means that there is no noise added. We use $\text{SNR} = \frac{E_b}{N_0}$ where E_b is the energy per bit and N_0 is the noise power spectral density. We set the SNR in our simulations by changing the noise power of a Gaussian random variable.

Inter-symbol Interference and Equalization

Inter-symbol interference occurs when we are transmitting over a channel that has some echos. These echos cause the receiver to hear a garbled signal instead of the original signal from the transmitter. This is called inter-symbol interference because the receiver is hearing a combination of symbols across time.

Let $\vec{x} = [x_0, x_1, \dots, x_n]$ be the set of n complex symbols that the transmitter sends over the channel that connects the transmitter to the receiver. Each channel will have different characteristics. Some channels may have echos, others may have delays, often channels will have both. When a channel has echos, this is called a multipath channel because there are multiple paths to reach the receiver. When a channel is modeled as a delay line, each multipath signal component at the receiver is modeled as a tap on the delay line. We can characterize a channel by characterizing the channel taps.

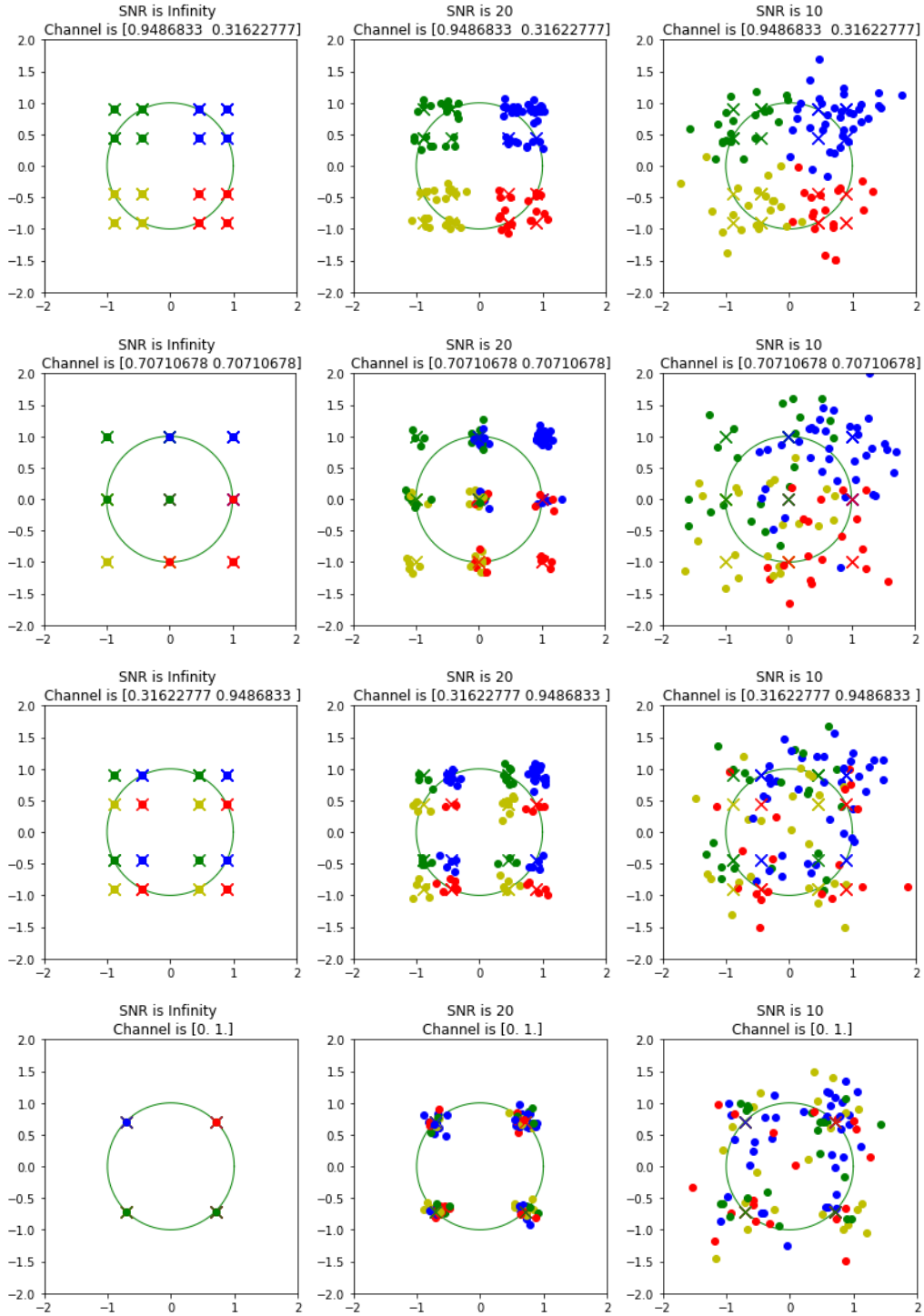
Let $\vec{a} = [a_0, a_1, \dots, a_\ell]$ be the characteristics for a multipath channel that has $\ell + 1$ taps. When a sequence of symbols like \vec{x} is transmitted over this channel, the channel taps are convolved over the sequence. Additionally, there is noise in the system denoted by v_m .

$$\tilde{x}_m = \sum_{i=0}^{\ell} a_i x_{m-i} + v_m \quad (1.1)$$

The receiver will hear a signal that is corrupted by inter-symbol interference and noise; $\vec{\tilde{x}} = [\tilde{x}_0, \tilde{x}_1, \dots, \tilde{x}_{n+\ell}]$. Receivers must be able to handle garbled signals in order to transmit data in the real world. The process of removing the inter-symbol interference is called equalization. The goal of equalization is to take in a garbled signal and output a signal with minimal inter-symbol interference.

Figure 1.1 demonstrates the effects of two tap channels on a QPSK modulation constellation. We show what the received signal constellations are from a sequence of 100 symbols modulated in QPSK through different two tap channels. The channel taps are normalized, $\|\vec{a}\|^2 = 1$. We assume that the channel taps are constant during the transmission of the signal. We see that under certain channel conditions, like when the two taps are equal, it is very difficult to distinguish between the four constellations. We also show how the constellations change with lower SNR. When the SNR is infinity, i.e. when no noise is added, the first and third channel constellations are distinct and clear. However, as we add noise (lower SNR), the symbol groupings get mixed and there are less clear distinctions between the symbols. Engineers have built processes to remove inter-symbol interference. First, let's go into the case when the channel characteristics are known.

Figure 1.1: The effects of a two tap channel on the QPSK constellation.



Equalization for a known channel

If the receiver knows the channel characteristics, \vec{a} , perfectly, then there are a few different methods that can be used. The zero-forcing equalizer applies the inverse of the channel response to the received signal. It is called zero-forcing because there will be zero inter-symbol interference if there is no noise. There are some limitations of the zero-forcing equalizer. First, the impulse response of the equalizer needs to be infinitely long. Second, if there is a weak signal at a frequency, then the inverse gain is going to be very large. This will amplify any noise in the system. Third, if there are any zeros in the frequency response, these cannot be inverted.

Another equalizer, the minimum mean squared error (MMSE) equalizer, handles noise much better than the zero-forcing equalizer. The MMSE equalizer minimizes the error between the equalized preamble and the original known preamble by choosing the optimal inverse of the channel response, W . Let H be the circulant matrix of the channel taps.

$$\vec{x}_{pre} = H\vec{x}_{pre} + \vec{v} \quad (1.2)$$

$$H = \begin{bmatrix} a_0 & 0 & & \cdots & 0 \\ a_1 & a_0 & 0 & \cdots & 0 \\ \vdots & & & & \vdots \\ a_\ell & a_{\ell-1} & \cdots & a_0 & 0 & \cdots & 0 \\ 0 & a_\ell & a_{\ell-1} & \cdots & a_0 & \cdots & 0 \\ \vdots & & & & & & \vdots \\ 0 & \cdots & 0 & a_\ell & a_{\ell-1} & \cdots & a_0 \end{bmatrix} \quad (1.3)$$

$$\min_W \|\vec{x}_{pre} - \vec{\hat{x}}_{pre}\|^2 \quad (1.4)$$

$$\vec{\hat{x}}_{pre} = W(H\vec{x}_{pre} + \vec{v}) \quad (1.5)$$

$$W^* = \vec{x}_{pre}\vec{x}_{pre}^T H^T (H\vec{x}_{pre}\vec{x}_{pre}^T H^T + \sigma_v^2 I)^{-1} \quad (1.6)$$

The MMSE equalizer works well for known and unknown channels and does not amplify noise like the zero-forcing equalizer. While it is important to consider how well a receiver can equalize with a known channel, this is rarely the case. Usually, we do not know the channel characteristics.

Equalization for an unknown channel

When the receiver does not know the channel characteristics, the process of equalization essentially has two jobs; first, identify the channel, second, remove the inter-symbol interference. If the receiver did not identify the channel first, there would be no way to remove the effects of it on the received signal.

In order to estimate a channel, most systems require that packets begin with a known sequence called a preamble. The signal sent is broken into two parts; $\vec{x} = [\vec{x}_{pre}, \vec{x}_{data}]$. The signal received on the transmitter is $\vec{\hat{x}} = [\vec{\hat{x}}_{pre}, \vec{\hat{x}}_{data}]$. The receiver knows what the original

preamble sequence was, \vec{x}_{pre} , and can use the received preamble sequence, $\vec{\tilde{x}}_{pre}$, to estimate the behavior of the channel. Once the channel is estimated, the receiver then equalizes the data, $\vec{\tilde{x}}_{data}$.

One common method to estimate the channel is to use the least-squares optimization framework. Let H be the circulant matrix of estimates of the channel taps. The least-squares channel estimator wants to find H that minimizes the mean squared error between the received signal, $\vec{\tilde{x}}_{pre}$, and what the predicted received signal would be if the the estimated channel H were applied to the original preamble, \vec{x}_{pre} .

$$\min_H \|\vec{\tilde{x}}_{pre} - H\vec{x}_{pre}\|^2 \quad (1.7)$$

The receiver needs to choose the number of non-zeros in H , representing how many taps (or echos) there might be in the channel. The number of taps, or non-zeros, in H will be set based on the environment of the transmitter and receiver. Once the channel response has been estimated, an equalizer uses that information to remove the inter-symbol interference from the received signal, like the MMSE equalizer defined above.

Carrier Frequency Offset and Correction

If we were to implement our minimum mean squared error equalizer on a real-world physical receiver, we would find some problems with our equalization process. Our equalizer will equalize the first symbols very well. However, as we equalize later parts of our sequence, we will encounter a physical phenomenon called carrier frequency offset, CFO. Carrier frequency offset can occur when the transmitter and receiver are at slightly different frequencies. It can also occur when the transmitter and receiver are moving, causing a sort of Doppler effect.

When there is a significant CFO present, the symbols will gradually rotate. The received signals will be the original signals rotated at a rate ω . Gaussian noise, v_m , is also added.

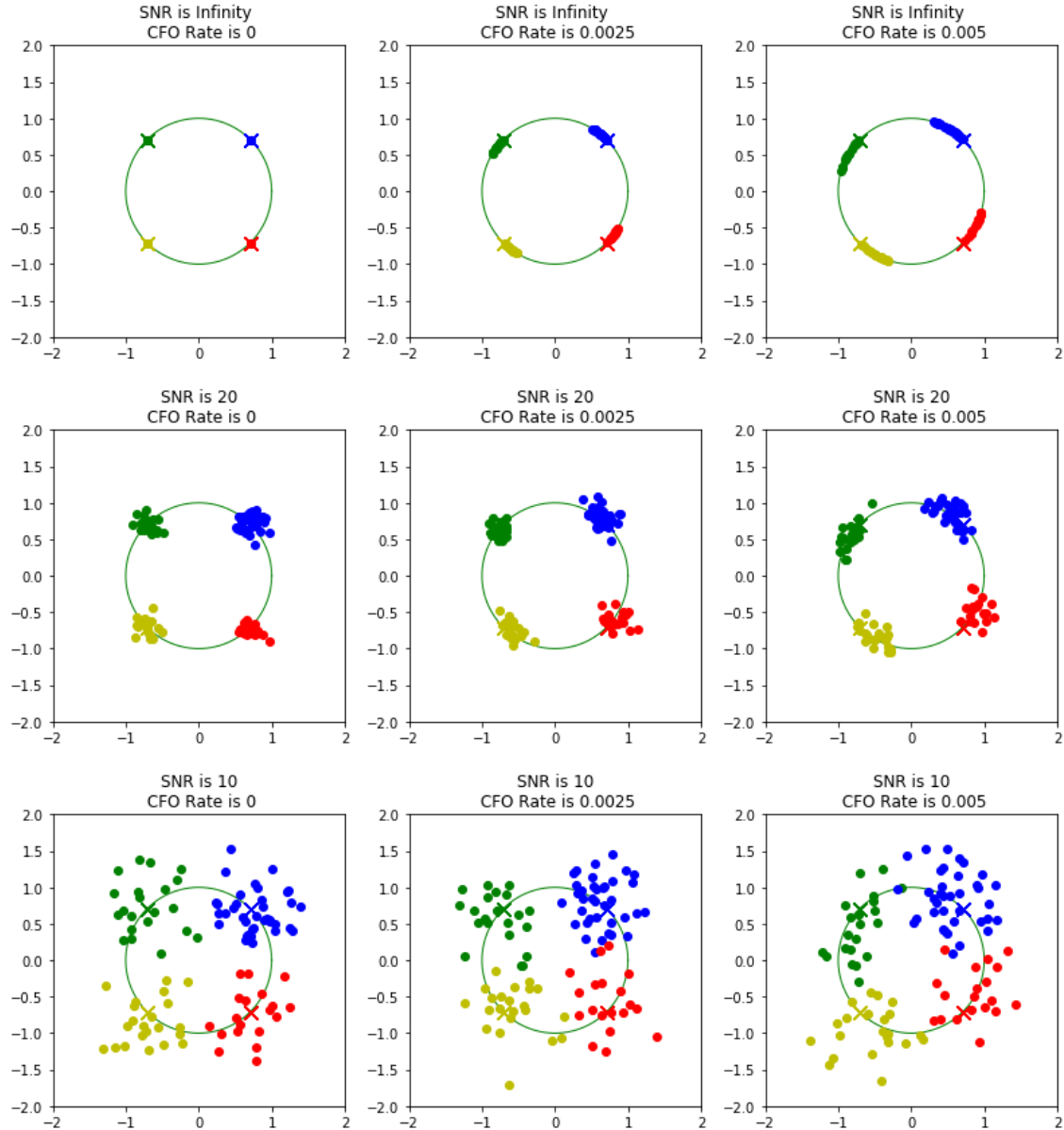
$$\tilde{x}_m = x_m e^{mj\omega} + v_m \quad (1.8)$$

Figure 1.2 demonstrates the effects of pure CFO on a QPSK modulation constellation (no multipath). We show what the received signal symbols constellations are from a sequence of 100 symbols modulated in QPSK with different CFO rates. We assume that the CFO rate, ω , is constant during the transmission of the signal. For small CFO rates, the data does not rotate that much. For larger CFO rates, the data turns more and as we add noise (lower SNR), the groupings of the symbols become mixed.

There are a few ways to handle CFO, some are more elegant than others. The first solution is to try to side-step the problem entirely. Since the effects of CFO depend on the length of a packet, one solution is to make packets so short that the symbols only move a little bit. In this case, the effects of CFO can be ignored.

Another solution is using a phase-locked loop, which is a control system that outputs a signal with a phase related to the input signal. A Costas loop is a circuit that implements a

Figure 1.2: The effects of a carrier frequency offset on the QPSK constellation.



phase-locked loop for CFO correction for continuous time signals by adapting the sampling rate [4]. However, this requires adjusting the sampling and is implemented in hardware. This does not allow the robustness and adaptability that we want.

What happens when there is both CFO and inter-symbol interference? The received signal will have the effects of the channel and CFO as well as the noise.

$$\tilde{x}_m = \left(\sum_{i=0}^{\ell} a_i x_{m-i} \right) e^{mj\omega} + v_m \quad (1.9)$$

Figure 1.3 demonstrates the effects of CFO and multi-tap channels on a QPSK modulation constellation. We show what the received signal symbols constellations are from a sequence of 100 symbols modulated in QPSK with different CFO rates and two tap channels. We assume that the CFO rate, ω , and the channel taps are constant during the transmission of the signal. The channel taps are normalized, $||\vec{a}||^2 = 1$ for each of the channels. We can see that for certain channel conditions, like for the channel in the third row, that the combination of carrier frequency offset and noise makes grouping these symbols quite difficult. Modern day receivers have to combine CFO correction, channel estimation, and equalization processes in order to handle these received signals.

1.3 Related Works

Recently, many researchers have become more interested in applying deep learning techniques to communications systems. We refer the reader to [3][5][25][9][21] for surveys and motivations of these kind of works. Most of the works that we have found typically only deal with additive white gaussian noise (AWGN) channels or Rayleigh fading channels.

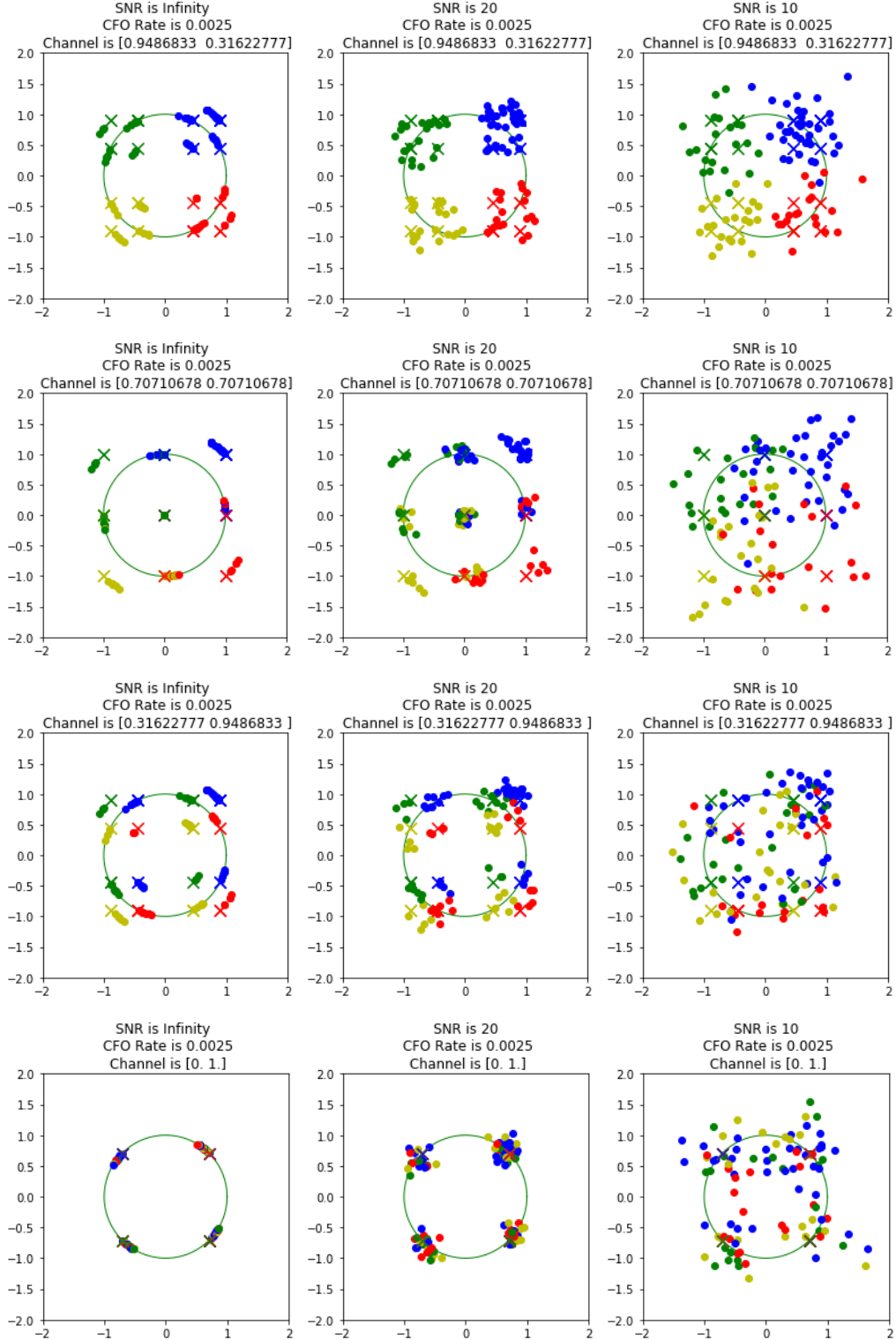
In [6], Dörner et. al. implemented an end to end transmitter and receiver with neural networks that allowed gradients to flow all the way back from the receiver during training. However, they restricted themselves to only sending a certain set of messages, only seeing AWGN channels, and they did not address CFO.

Other groups have been focusing on decoding with recurrent neural networks (RNNs) but also only deal with AWGN channels [11][12]. Others have been working with generative adversarial networks (GANs) to train an end-to-end communication system [26]. However, they also only consider AWGN channels or Rayleigh fading channels.

For those that do consider more complex channels, most re-train their models for each new channel seen. Ye et. al. consider OFDM systems where their feedforward neural networks estimate the channel state information and then train offline for that specific channel [27]. [23] considers nonlinear channels but re-trains the network for each new channel. Note, this work does not go into detail about the architecture used or how the networks are trained.

Goldsmith and Farsad train a detector for optical and molecular channels [7]. They assume a Poisson model for the channel and attempt to predict the probability mass function.

Figure 1.3: The effects of a two tap channel and a carrier frequency offset on the QPSK constellation.



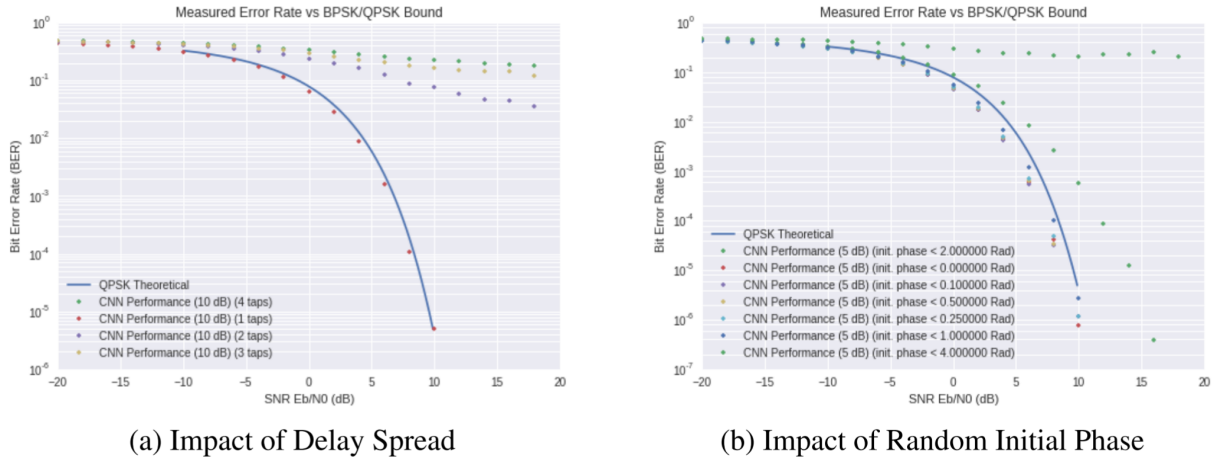


Figure 1.4: The impacts of multipath channels and initial phase offsets on bit error performance [17]. Note, the network is not re-trained for each new environment. As the number of paths and carrier frequency offset increase, the performance of the CNN decreases.

However, they assume that they retrain for each Poisson parameter and they do not address CFO.

Timothy O'Shea's group has been doing some excellent work in this area. O'Shea's first work in this area jointly optimized a transmitter and receiver for a given channel model (AWGN and Rayleigh fading) but the work does not go into detail about the performance for various channels [21].

In a subsequent paper, O'Shea et. al. explored how multipath channels and random initial phase affects the system performance [17]. This is the only work that we know of that does not re-train for each new multipath channel or carrier frequency offset. Figure 1.4 shows how the convolutional neural network performance drastically decreased with multipath channels and for non-zero phase offset. This particularly motivates our work as we hope to increase performance on learning to learn equalization and CFO.

In [19], they added an attention model to perform synchronization for time, frequency, phase, and sample timing offset. However, their work showed that the synchronization was still quite noisy and did not perform much better than without the attention model [17]. Additionally, this work did not consider multipath channels.

O'Shea et. al. have also attempted to use GANs to approximate the channel response model for AWGN channels with and without phase noise [18]. The GANs were unable to find accurate probability density functions to represent the channel response.

In [22], they use convolutional neural networks (CNNs) to estimate, but not correct, carrier frequency offset and timing offset. However, they only consider AWGN channels and Rayleigh fading channels. In [20], they explore how unsupervised learning can train autoencoders for multiple antenna communications. They do re-train for each new channel

and use a Rayleigh fading channel model.

1.4 Contributions of this Report

In this report, we show that a channel estimator can learn to estimate two tap, real-valued channels but that more work is needed for the network to perform comparably to the baseline, the least-squares channel estimator. We also design an equalizer recursive neural network that equalizes a stream of data for any two tap, real-valued channel. This equalizer does not need to re-train for each new channel and performs better than the baseline, the minimum mean squared error equalizer, for low signal to noise ratios. We also show that the baseline is much better than the neural network at handling difficult channels, specifically channels that have approximately equal taps.

This work is novel in this area because all of the other related works have re-trained the networks for each new channel they encounter. This is the first equalizer that learns to learn to handle each new channel without having to re-train. Table 1.1 compares our work with previous works.

Table 1.1: Comparison of our work and related works.

| | Goldsmith [7] | Ye [27] | O'Shea [17] | Our Work |
|--|---------------|-------------------|-------------------|-------------------|
| Types of channels considered | Poisson | AWGN Multipath | AWGN Multipath | AWGN Multipath |
| Do they address CFO? | No | No | Yes | Yes |
| Do they re-train for each new environment? | No | Yes | No | No |

Additionally, we design a neural network to estimate and correct carrier frequency offset that does not need to be retrained for each new carrier frequency offset. This work is novel in this area because the other related works have re-trained the networks for each new carrier frequency offset. This is the first carrier frequency offset correction network that learns to learn to handle each new carrier frequency offset without having to re-train. We also explore whether recursive neural networks can behave as clocks and show that this is a non-trivial problem that needs more exploration.

Chapter 2

Deep Networks for Equalization¹

The equalization process removes inter-symbol interference caused by the channel from a sequence of symbols. We will explore how neural networks can both estimate the channel characteristics and remove the inter-symbol interference. For the remainder of this chapter, we will assume that all channels only have real parts.

2.1 Channel Estimation

The process of channel estimation tries to identify the coefficients of the channel taps as best as possible. We want to estimate \vec{a} that minimizes the difference between the estimated channel taps, $\vec{\hat{a}}$, and the true channel taps, \vec{a} . We can choose $\vec{\hat{a}}$ to minimize this difference.

In Figure 2.1, we compare the performance of neural network channel estimator, a K-Nearest-Neighbors (KNN) channel estimator (where $K = 15$), and the least-squares channel estimator for two tap real channels. We train the neural network and KNN on $40k$ data sequences with random channels. Each data sequence has a preamble length of 100 bits modulated to 50 QPSK symbols. Each channel tap is a uniform random variable, $[-1, 1]$. We normalize the channel power such that $\|\vec{a}\|^2 = 1$.

The neural network channel estimator takes the received preamble, \vec{x}_{pre} , and the known preamble, \vec{x}_{pre} , as inputs. The output of the neural network is the estimate of the channel taps, $\vec{\hat{a}}$. We assume that the neural network knows there are only two channel taps and thus, limit the size of the output to two. The neural network has an architecture that consists of three dense layers. The first and second layers have 300 nodes with sigmoid activation functions. The final layer has two nodes with the tanh activation function. The network has a decaying learning rate that starts at 0.01. The loss function is mean squared error of the true channel and channel estimates; $\|\vec{a} - \vec{\hat{a}}\|^2$.

All three channel estimators are tested on $10k$ data sequences with random channels. The neural network and the KNN channel estimators are not re-trained for each new channel.

¹The contents of this chapter were produced in collaboration with Nipun Ramakrishnan, an undergraduate researcher that I am mentoring.

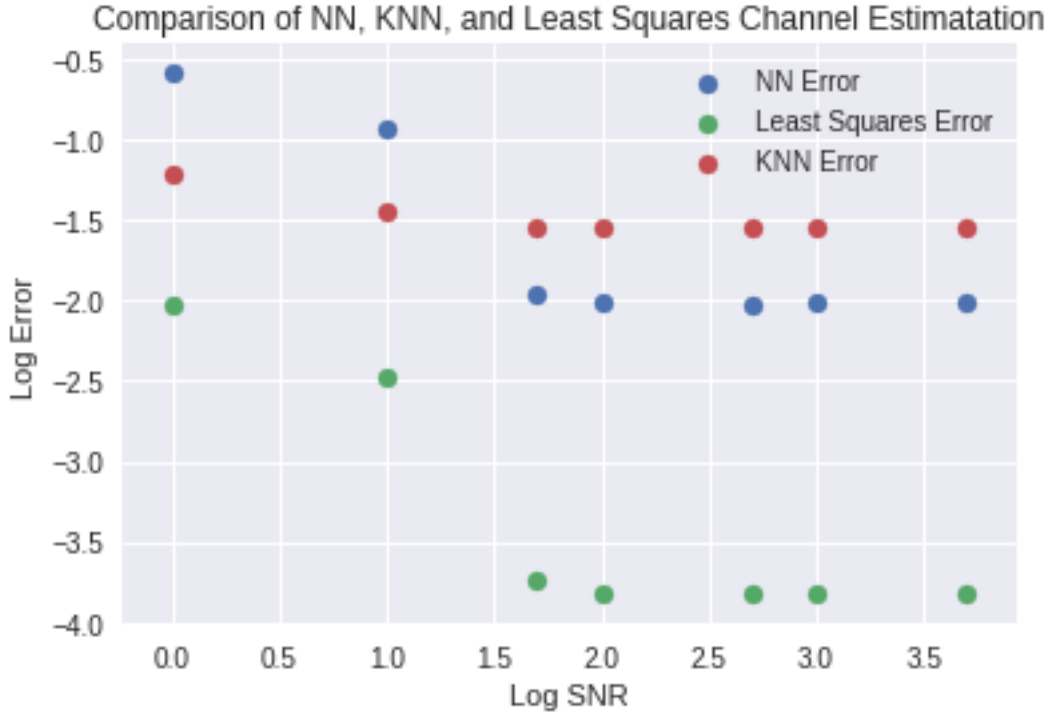


Figure 2.1: Comparison of a neural network, a K-Nearest-Neighbors, and a least-squares channel estimators for two tap, real channels.

The least-squares estimator consistently outperforms the neural network and KNN channel estimators, especially for high SNRs. For low SNRs, the least-squares estimator is better than the neural network by a factor of about five. For high SNRs, the least-squares estimator is better than the neural networks by a factor of about thirty. The neural network channel estimator performs better than the KNN channel estimator for high SNRs by a factor of about three but worse for low SNRs by a factor of about three.

2.2 Channel Equalization

The process of equalization tries to remove inter-symbol interference from a sequence of symbols. In a two tap channel, the equalization process is trying to remove the effect from the second channel tap. Given the received symbol, \tilde{x}_m , the channel taps, a_0, a_1 , and the previous symbol, x_{m-1} , we can solve for the best estimate of the original symbol.

$$\tilde{x}_m = a_0 x_m + a_1 x_{m-1} + v_m \quad (2.1)$$

$$\hat{x}_m = \frac{\tilde{x}_m - a_1 x_{m-1}}{a_0} \quad (2.2)$$

If we want neural networks to perform equalization, they might try to solve for \hat{x}_m in this way. If that is the case, then the neural networks will need to know how to take an input and divide by it and how to multiply two inputs.

Learning an inverse

We explore whether a neural network can learn to do division. We did not do a thorough architecture search or train for very long as we wanted an idea of the behavior of the neural network and will do this in future work. Figure 2.2 shows the topological representation of the division function, $z = \frac{x}{y}$. A neural network that is trained to do division will need to approximate the surface of this function. When y is close to zero, z goes to infinity or negative infinity. These large peaks will make it difficult for a neural network to do division.

Figure 2.3 shows the performance of a neural network learning division. The neural network's inputs are x and y and the output is \hat{z} . Both inputs are uniform random variables; x is drawn uniformly $[-1, 1]$ and y is drawn uniformly $[\beta, 1]$. The loss function is the mean squared error between the estimated division and the true division, $||\hat{z} - \frac{x}{y}||^2$. The architecture consists of two dense layers with 50 nodes each and sigmoid activation functions. This feeds into a linear layer that outputs the scalar estimate of \hat{z} . The network is re-trained with $10k$ data points for different values of β . As β gets close to zero, the error increases dramatically. This is expected because the closer β is to zero, the stronger the effects are of the infinite peaks of the division function.

Learning to multiply two inputs

We explore whether a neural network can learn to do multiplication. We did not do a thorough architecture search or train for very long as we wanted an idea of the behavior of the neural network and will do this in future work. Figure 2.4 shows the topological representation of the multiplication function, $z = xy$. A neural network that is trained to do multiplication will need to approximate the surface of this function. When x and y are close to zero, the function is stable. However, as x and y grow, the function becomes unstable and grows to infinity or negative infinity.

Figure 2.5 shows the performance of a neural network learning multiplication. The neural network's inputs are x and y and the output is \hat{z} . Both inputs are uniform random variables; x is drawn uniformly $[-1, 1]$ and y is drawn uniformly $[0, \gamma]$. The loss function is the mean squared error between the estimated multiplication and the true multiplication, $||\hat{z} - xy||^2$. The architecture consists of two dense layers with 50 nodes each and sigmoid activation functions. This feeds into a linear layer that outputs the scalar estimate of \hat{z} . The network is re-trained with $50k$ data points for different values of γ . As γ increases, the error increases drastically. This is expected because the larger γ is, the stronger the effects are of the infinite limits of the multiplication function.

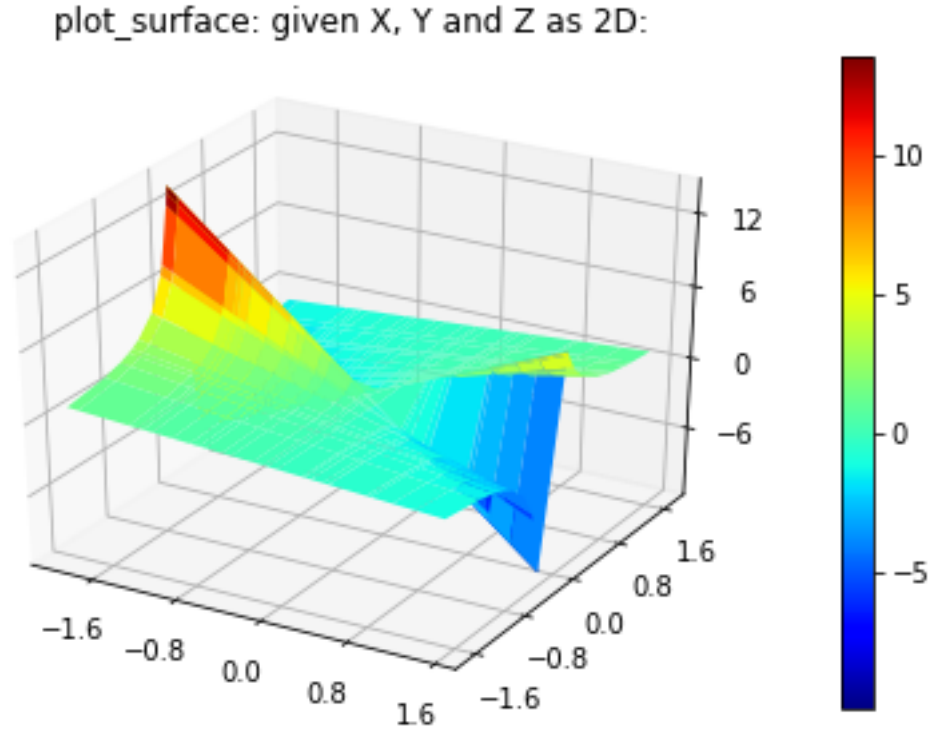


Figure 2.2: Topographical surface representation of the division function; $z = \frac{x}{y}$.

RNN for Channel Equalization

We design an RNN to perform channel equalization. The inputs of the RNN are the true channel taps. We assume a two tap channel, so a_0, a_1 are inputs to the RNN. Another input is the received data sequence, \vec{x}_{data} . The output of the RNN is an estimate of the original data sequence, $\vec{\hat{x}}_{data}$. The loss function that the RNN is trained on is the mean squared error between the original and estimated data sequence, $\|\vec{x}_{data} - \vec{\hat{x}}_{data}\|^2$.

The neural network architecture for our equalizer uses a special type of RNN called a bidirectional long-short term memory (LSTM) network. A bidirectional RNN connects the forward and backwards nodes of the RNN, allowing the output to depend on both future and past states. An RNN that has LSTM units allows the nodes to store memory for a short period of time. Our network also has time-distributed dense layers which are used after the LSTMs to flatten the output by applying a fully connected dense layer at each time step.

The inputs to our RNN are the channel taps concatenated with a sequence of 10 of the data symbols. The inputs are fed into four layers of bidirectional LSTM units. Each layer has a state size of 90 units for each forward and backward state. The output of the four bidirectional LSTM layers is 180 units for each of the 10 symbols. This is then fed into two time-distributed dense layers with 180 nodes and 100 nodes respectively and they each have

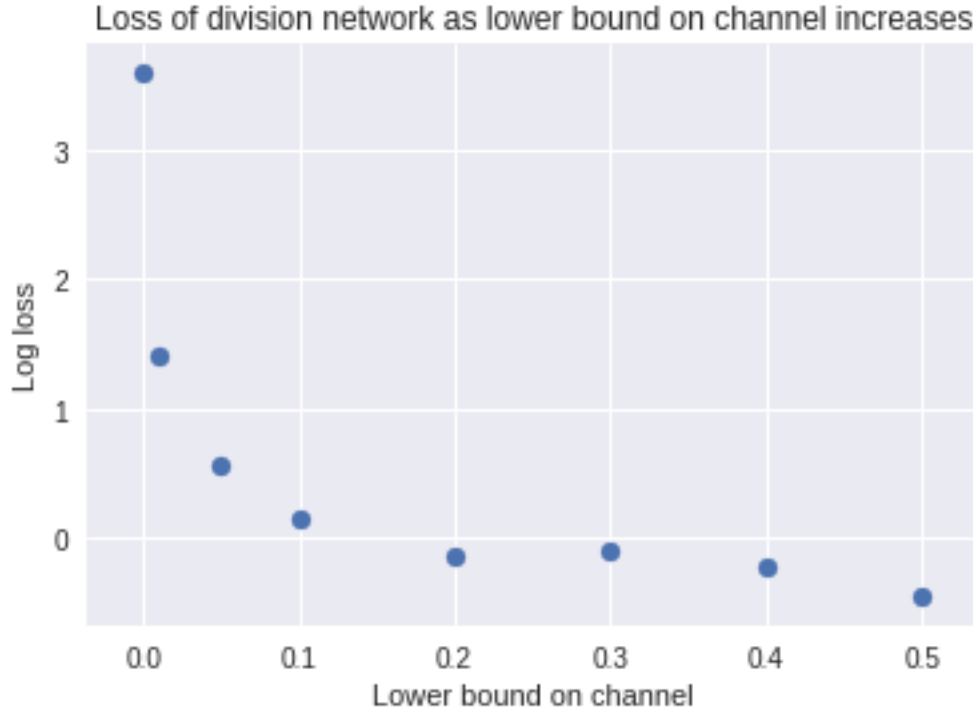


Figure 2.3: Neural network division loss with respect to the lower bound of y .

Relu activation functions. These time-distributed dense layers bring the output size down from 180 to 100 to 100. The output is then fed into one final time-distributed dense layer with 100 nodes and linear activation. The outputs of the RNN are 10 equalized symbols. For data sequences longer than 10 symbols, the process is repeated for each set of 10.

Figure 2.6 compares the mean squared error loss on test data for the RNN architecture defined above and the MMSE equalizer defined in the introduction. The RNN is trained on $40k$ data sequences, with a decaying learning rate starting at 0.01. The RNN and MMSE are tested on the same $10k$ data sequences. New data is generated for each SNR and the RNN is re-trained for each SNR. Each data sequence consists of 60 bits modulated to 30 complex symbols. We assume that the RNN and MMSE both have access to the true channel taps.

The RNN equalizer performs comparably to the MMSE equalizer for low SNRs and performs better than MMSE equalizer for medium SNRs. However, the RNN equalizer performs worse than the MMSE equalizer for high SNRs. We take a closer look at what kind of channels the RNN equalizer is getting wrong for these high SNRs.

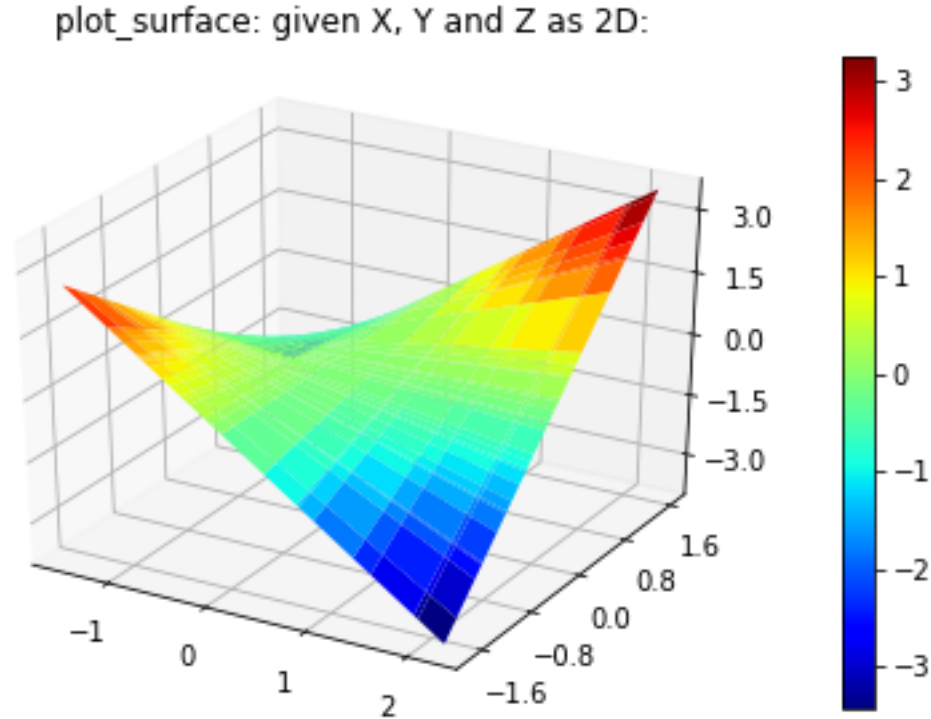


Figure 2.4: Topographical surface representation of the multiplication function; $z = xy$.

| Tap 1 | Tap 2 | Counts of bad estimates | Mean Squared Error |
|--------|--------|-------------------------|--------------------|
| 0.707 | 0.707 | 14 | 0.7162 |
| -0.707 | 0.707 | 1 | 0.5236 |
| 0.707 | -0.707 | 16 | 0.8172 |
| -0.707 | -0.707 | 12 | 0.6492 |

Figure 2.7 shows which channels the RNN equalizer paired with a classic demodulator get any bits wrong. There are a total of $50k$ random channels in the test set with SNR= 100. Of that set, only 43 channels result in incorrect bits after the RNN equalizer and classic demodulator. From the figure, it is clear that these difficult channels are clustered into four regions. All of the four regions are when the first and second tap of the channel are equal in magnitude; $|a_0| \approx |a_1|$. The mean squared error between the equalized data and the true data among just the bad channels is 0.7306. The mean squared error among the good channels is 0.000264. So when the RNN equalizes incorrectly, it gets it very wrong.

We expect these channels to be difficult. Refer to Figure 1.1 to see how the constellations of two bits fall right on top of each other. Additionally, if we think about these channels in terms of the infinite impulse response, these channels could result in highly variable inverses.

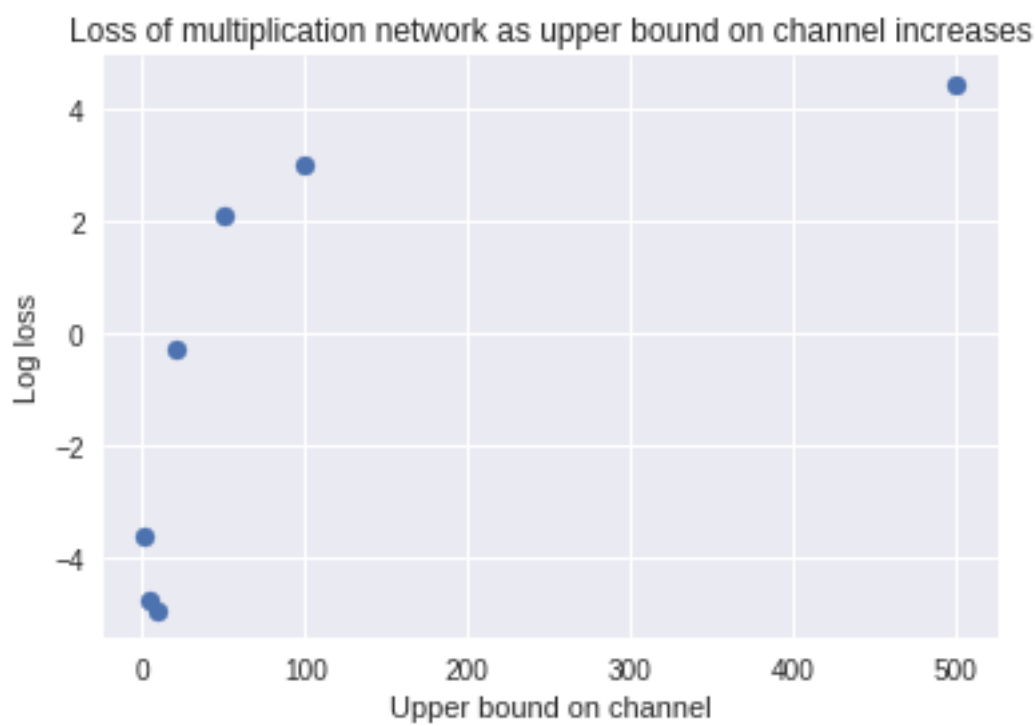


Figure 2.5: Neural network multiplication loss with respect to the upper bound of y .

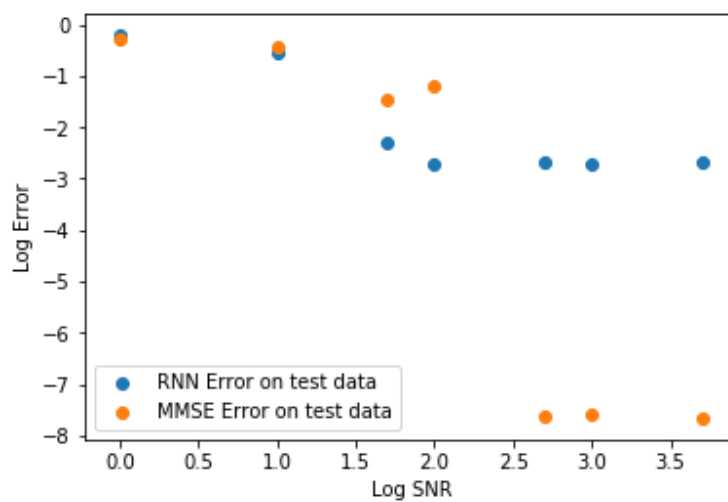


Figure 2.6: Comparison of RNN and MMSE equalizers loss over SNR.

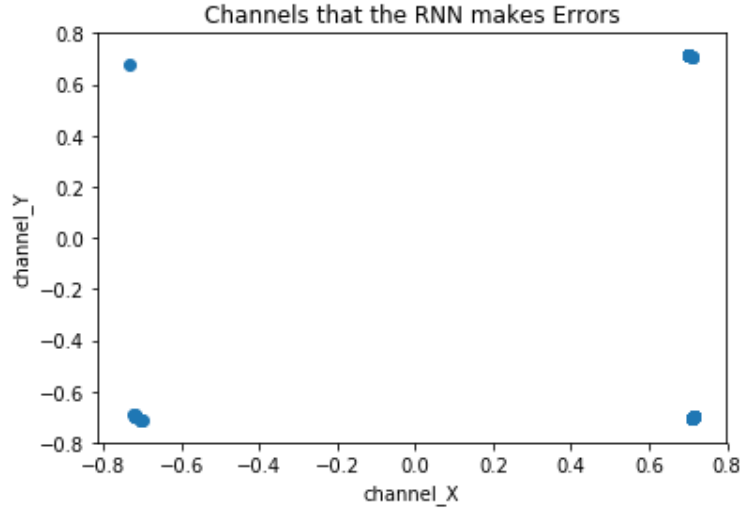


Figure 2.7: What two tap channels does the equalizer get wrong?

Because the two channel taps are equal, the infinite impulse response is constantly cancelling out and everything will have to balance out exactly. This representation means there will be a lot of variation of the inverse of those channels, causing the system to be almost unstable.

In this chapter, we designed a deep neural network to estimate two tap channels. We explored the difficulties in learning to divide and learning to multiply with neural networks. We showed that a deep recursive neural network, with bidirectional LSTM layers, can learn to learn to equalize for random two tap channels. We also examined when the RNN equalizer fails.

2.3 Future Work

In the future, we would like to expand our channel estimation networks to handle any number of channel taps and channel taps that also have imaginary parts. We would also like to explore how we might be able to learn to estimate the channel even without a shared preamble. We hope to explore and expand our research on equalization networks. One first step will be investigating how adding logarithmic layers into the networks affects the performance. We want to combine the channel estimator and equalizer networks into one deep network and have gradients flow all the way back. This will take the work one step further in the direction of learning to learn. We also want to do more thorough architecture searches and train the networks for longer.

Chapter 3

Deep Networks for Carrier Frequency Offset¹

The carrier frequency offset (CFO) correction process removes the symbol rotation caused by the non-ideal conditions of the transmitter and receiver environment. We will explore how neural networks can both estimate the CFO and correct for the rate of rotation without using backpropagation for each new rate.

3.1 Recurrent Neural Network Follows a Circle

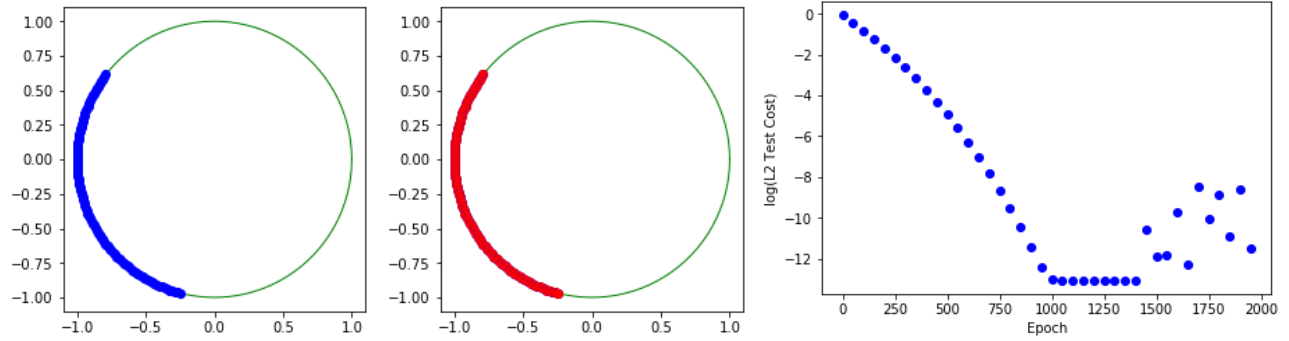
A carrier frequency offset affects the received symbols by slowly rotating them around the origin. If we want a neural network to be able to correct this, we essentially need it to act as a clock, following around a circle at a certain rate. We first explore what kind of architectures are necessary for a recurrent neural network (RNN) to learn how to follow a circle. We show that a simple, small linear RNN can learn to act as the rotation matrix for a single rate of rotation around a circle. We also show that in order to have an RNN follow a circle for different rates of rotation, we need a large non-linear layer. We will then design a deep neural network to estimate the rate of rotation caused by CFO and correct it.

Single Rate

Figure 3.1 shows a recurrent neural network (RNN) that follows a circle for a single rate of rotation. The network's input is the starting point, x_0 , and it must predict the next 99 points, $x_1 \dots x_{99}$. Point x_m is rotated by $e^{mj\omega}$ where ω is the rate of rotation. The RNN architecture is one linear layer that takes the current point as input, x_m , and outputs the next point, x_{m+1} . We are forcing the state of the RNN to be the estimate of the next point;

¹The contents of this chapter were produced in collaboration with Nikhil Shinde, an undergraduate researcher that I am mentoring.

Figure 3.1: Linear neural network: follow a circle for a constant CFO rate, $\omega = 0.02$. The figure on the left shows the original data of length 100. The figure in the middle shows the estimated data from the neural network given the starting point from the left figure. The figure on the right shows the log of the test cost of the neural network versus training epoch.



state = \hat{x}_{m+1} . We use a linear layer here because the network essentially has to learn how to become the rotation matrix which is linear with respect to the input.

$$R = \begin{bmatrix} \cos(\omega) & -\sin(\omega) \\ \sin(\omega) & \cos(\omega) \end{bmatrix} \quad (3.1)$$

The RNN is trained for a constant rate of rotation, ω , applied to sequences of 100 points. The network trained for $2k$ epochs, each with a batch size of $1k$ data sequences and with a learning rate of 0.006. The initial starting point, x_0 is a uniform random variable drawn from on the unit circle. The network is tested on $1k$ new data sequences but with the same ω . Figure 3.1 shows the results of an RNN trained and tested for $\omega = 0.02$. The network achieves a loss on the order of 10^{-13} . However, the minimum loss is not very sticky. Notice that the figure on the right in figure 3.1 shows the network jumping away from this minimum while training at around 1400 epochs. A more thorough search over learning rate decay is needed to prevent this jumping from happening.

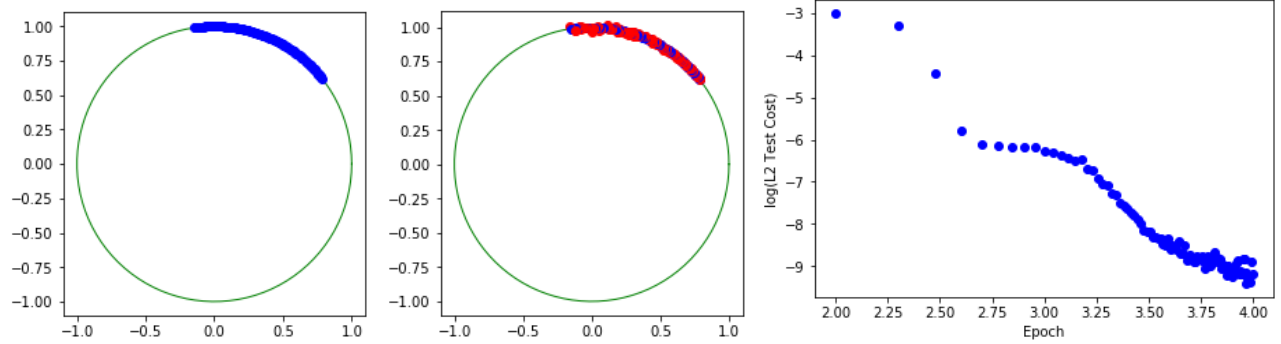
Different Rates

Figure 3.2 shows a recurrent neural network (RNN) that follows a circle for a given rate of rotation. The network's inputs are the starting point, x_0 and the rate of rotation, ω . The RNN must predict the next 99 points, $x_1 \dots x_{99}$. Point x_m is rotated by $e^{mj\omega}$.

We cannot use just linear layers in this case because the RNN must learn to find $\cos(\omega)$ and $\sin(\omega)$ for different ω . Essentially, this RNN needs to approximate the sin, cos functions and then apply them to the state.

The RNN architecture consists of one non-linear layer, with 100 nodes and Relu activation functions, and a linear layer at the output with 2 nodes. The inputs to the RNN are the

Figure 3.2: Nonlinear neural network: follow a circle for different CFO rates. The figure on the left shows the original data of length 100 for $\omega = 0.01063284$. The figure in the middle shows the estimated data from the neural network given the starting point and rate of rotation from the left figure. The figure on the right shows the log of the test cost of the neural network versus training epoch.



estimated current point, \hat{x}_m , and the rate of rotation, ω . The RNN outputs the estimate of the next point, \hat{x}_{m+1} . We are forcing the state of the RNN to be the estimate of the next point concatenated with the rate of rotation; state = $[\hat{x}_{m+1}, \omega]$ so that state then becomes the input to the next run of the RNN.

The RNN is trained for different rates of rotation, ω , applied to sequences of 100 points. The rate of rotation is a uniform random variable $[0, \frac{1}{50}]$. We use mean squared error of the true data sequence and the estimated data sequence for the loss function; $\|\vec{x} - \vec{\hat{x}}\|^2$. The network is trained for $10k$ epochs, each with a batch size of $1k$ data sequences and with a learning rate of 0.01. The initial starting point, x_0 is a uniform random variable drawn from on the unit circle. The network is tested on $100k$ new data sequences with random ω . Figure 3.2 shows the results of an RNN trained and tested for $10k$ epochs. The network achieves a loss on the order of 10^{-5} but is still data hungry and should be trained for more epochs.

3.2 Deep Network Carrier Frequency Offset Estimation and Correction

We explore estimating carrier frequency offset and correction with neural networks. One of the challenges with CFO is that it cannot be directly measured in the real world. Therefore, all of the loss functions will not depend on ω . Again, we want a neural network that does not need to update its weights every time it encounters a new rate of rotation.

The inputs of the neural network are the known preamble, \vec{x}_{pre} , and the received preamble, $\vec{\hat{x}}_{pre}$. The neural network outputs the estimated preamble symbols, $\vec{\hat{x}}_{pre}$. The loss

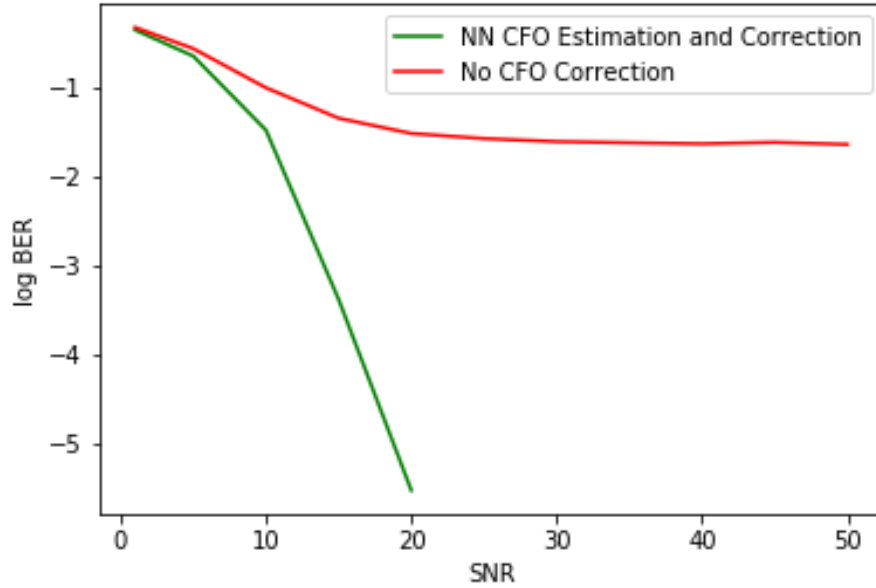


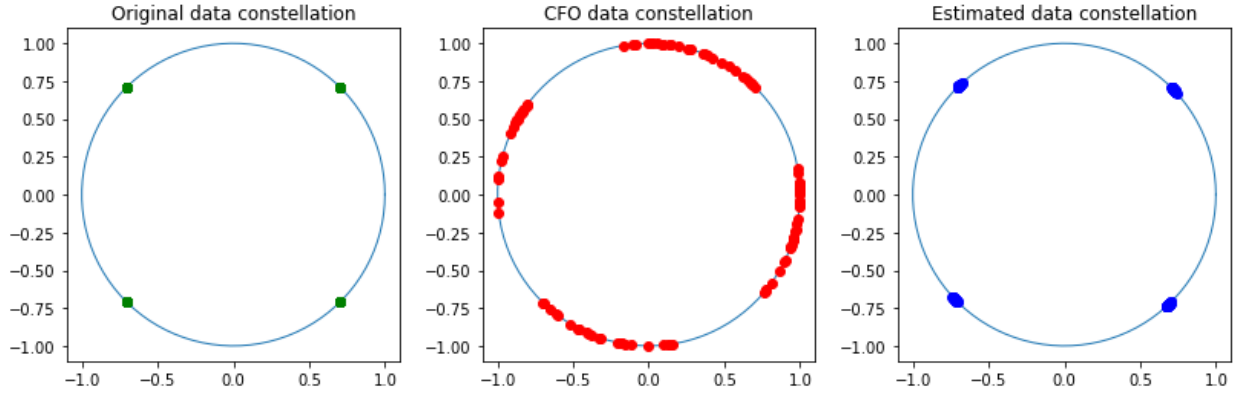
Figure 3.3: The log BER of received signals passed through a neural network CFO estimator with a rotation matrix CFO correction and classic demodulator. The log BER of the same signals passed through just the classic demodulator.

function is the mean squared error between the original preamble and the estimated preamble, $\|\vec{x}_{pre} - \hat{\vec{x}}_{pre}\|^2$. The neural network architecture consists of two non-linear layers, each with 10 nodes and sigmoid activations. This then feeds into a linear layer that outputs a scalar, treated as the estimate for the rate of rotation, $\hat{\omega}$. The received preamble is then rotated using a special layer with the function $e^{-mj\hat{\omega}}$ that is differentiable and can pass gradients.

Figure 3.3 compares the performance of the neural network correcting CFO versus no CFO correction. The training and testing data are solely for a one tap channel and have varying values of ω . The preamble length is kept constant at 100 symbols in QPSK. The neural network is trained for $15k$ epochs with batch size of 100 sequences and with a learning rate of 0.001. The neural network and no CFO correction are tested on $10k$ sequences with different rates of rotation, ω . The rate of rotation, ω , is a uniform random variable, $[0, \frac{1}{100}]$. The training and testing process is repeated for each SNR. Note, we re-train the network for each SNR. We compare the bit error rate of the neural network and classic demodulator versus the bit error rate of just the classic demodulator without any CFO correction.

For SNR above 25, the neural network does not have any incorrect bits in the $10k$ data sequences each with 100 symbols, or 200 bits. The neural network consistently performs better than just the classic demodulator for all SNR. Although, the performance is similar for SNR= 1. Figure 3.4 shows one example of a data sequence passed through the neural

Figure 3.4: Nonlinear neural network estimating and correcting CFO for different ω . In this particular test example, $\omega = 0.00969545$ and $\text{SNR} = \infty$. The figure on the left shows the constellation of the 100 original data symbols. The figure in the middle shows the constellation of the original data symbols after CFO. The figure on the right shows the constellation of the data symbols after the neural network CFO correction.



network described above. The neural network is trained for $\text{SNR} = \infty$ with the same training sequence as described above. The estimated data constellation plot has removed most of the effects of the CFO and looks quite similar to the original data constellation plot.

In this chapter, we demonstrated that recursive neural networks can act like clocks; they can follow a circle for a given rate of rotation. They can learn to learn to follow a circle for a given rate. We have also demonstrated that a neural network can learn to learn how to estimate CFO and correct for it. All of the simulations in this chapter were for single tap channels. We attempted to learn how to estimate CFO and correct for it when there were random two tap channels. However, the results were not promising and will need to be explored more in future work.

3.3 Future Work

In addition to the future work described in Chapter 2, we want to expand our CFO estimation and correction work. We want to estimate CFO even in the presence of multi-tap channels and when the channels have complex taps. We also want to explore replacing the rotation layer that does the CFO correction with the RNN that follows a circle. Then the entire process will be based in neural networks without much structure. We also want to combine all of the aspects in this report to create one deep equalizer network that corrects for CFO and equalizes the channel. We may not explicitly need networks that do the estimation but it is a good place to start.

In general, this report was more of an exploration into how neural networks can learn to learn to communicate for equalization and CFO correction. We would like to do a more

thorough architecture search and longer training runs. We also hope to verify the ideas presented in this report by training on real radio data instead of simulated data.

Chapter 4

Conclusion

In this report, we reviewed the current literature on machine learning systems for communications systems; specifically equalization and carrier frequency offset. We also discussed the value of robustness and adaptability in future communications systems requiring an investigation of meta-learning. We demonstrated that neural networks can learn to learn to estimate channel characteristics for two tap channels. We explored using deep recursive neural networks to learn to learn to equalize. We also used recursive neural networks to learn to learn to be a clock by rotating in a circle at different rates. Lastly, we used deep neural networks to estimate the rate of carrier frequency offset and correct it.

All of this was achieved without using backpropagation when faced with new environments. From here, we plan to continue developing our work and design an end-to-end receiver, including demodulation and error correction, in neural networks that can adapt to new multipath channels and carrier frequency offsets. Communications is a perfect application to explore the possibilities and limitations of the current meta-learning and machine learning frameworks.

Bibliography

- [1] Yoshua Bengio, Samy Bengio, and Jocelyn Cloutier. “Learning a synaptic learning rule”. In: *Technical Report 751, Département d’Informatique et de Recherche Opérationnelle, Université de Montréal, Canada* (1990).
- [2] Yoshua Bengio et al. “A Neural Probabilistic Language Model”. In: *Neural Information Processing Systems*. Vol. 13. 2001, pp. 932–938.
- [3] Corina Botoca, Georgeta Budura, and Miclau Nicolae. “Nonlinear Channel Equalization Using Complex Neural Networks”. In: *Seria Electronică Și Telecomunicații, Transactions on Electronics and Communications* (2004), pp. 227–231.
- [4] John P. Costas. “Synchronous Communications”. In: *Proceedings of the IRE*. Vol. 44. 12. IEEE, 1956, pp. 1713–1718.
- [5] Theo Diamandis. “Survey on Deep Learning Techniques for Wireless Communications”. In: *Departement of Electrical Engineering, Stanford University* (2017).
- [6] Sebastian Dörner et al. “Deep Learning-Based Communication Over the Air”. In: *arXiv preprint arxiv:1707.03384* (2017).
- [7] Nariman Farsad and Andrea Goldsmith. “Neural Network Detection of Data Sequences in Communication Systems”. In: *IEEE Transactions on Signal Processing* PP (2018).
- [8] Chelsea Finn, Pieter Abbeel, and Sergey Levine. “Model-agnostic meta-learning for fast adaptation of deep networks”. In: *International Conference on Machine Learning* (2017).
- [9] Hengtao He et al. “Model-Driven Deep Learning for Physical Layer Communications”. In: *arXiv preprint arxiv:1809.06059* (2018).
- [10] Siavash Khodadadeh, Ladislau Bölöni, and Mubarak Shah. “Unsupervised Meta-Learning For Few-Shot Image and Video Classification”. In: *arXiv preprint arxiv:1811.11819* (2018).
- [11] Hyeji Kim et al. “Communication Algorithms via Deep Learning”. In: *International Conference on Learning Representations* (2018).
- [12] Hyeji Kim et al. “Deepcode: Feedback Codes via Deep Learning”. In: 2018.

- [13] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. “ImageNet Classification with Deep Convolutional Neural Networks”. In: *Neural Information Processing Systems*. 2012, pp. 1106–1114.
- [14] Brenden M. Lake, Ruslan Salakhutdinov, and Joshua B. Tenenbaum. “Human-level concept learning through probabilistic program induction”. In: *Science*. Vol. 350. 6266. 2015, pp. 1332–1338.
- [15] Christiane Lemke, Marcin Budka, and Bogdan Gabrys. “Metalearning: a survey of trends and technologies”. In: *Artificial Intelligence Review*. Springer, 2015, pp. 117–130.
- [16] Volodymyr Mnih et al. “Human-level control through deep reinforcement learning”. In: *Nature*. Vol. 518. 2015, pp. 529–533.
- [17] Timothy J. O’Shea, Kiran Karra, and T. Charles Clancy. “Learning to Communicate: Channel Auto-encoders, Domain Specific Regularizers, and Attention”. In: *IEEE International Symposium on Signal Processing and Information Technology*. 2016, pp. 223–228.
- [18] Timothy J. O’Shea, Tamoghna Roy, and Nathan West. “Approximating the Void: Learning Stochastic Channel Models from Observation with Variational Generative Adversarial Networks”. In: *arXiv preprint arxiv:1805.06350* (2018).
- [19] Timothy J. O’Shea et al. “Radio Transformer Networks: Attention Models for Learning to Synchronize in Wireless Systems”. In: *arXiv preprint arxiv:1605.00716* (2016).
- [20] Timothy O’Shea, Tugba Erpek, and T. Charles Clancy. “Deep Learning-Based MIMO Communications”. In: *arXiv preprint arXiv:1707.07980* (2017).
- [21] Timothy O’Shea and Jakob Hoydis. “An Introduction to Deep Learning for the Physical Layer”. In: *IEEE Transactions on Cognitive Communications and Networking*. Vol. 3. 4. 2017, pp. 563–575.
- [22] Timothy O’Shea, Kiran Karra, and T. Charles Clancy. “Learning Approximate Neural Estimators for Wireless Channel State Information”. In: *arXiv preprint arXiv:1707.06260* (2017).
- [23] K. T. Raghavendra and Amiya K. Tripathy. “An Efficient Channel Equalizer Using Artificial Neural Networks”. In: *Neural Network World* 16 (2006), pp. 357–368.
- [24] David Silver et al. “Mastering the game of Go with deep neural networks and tree search”. In: *Nature*. Vol. 529. 2016, pp. 484–489.
- [25] Tianqi Wang et al. “Deep Learning for Wireless Physical Layer: Opportunities and Challenges”. In: *China Communications* 14.11 (2017), pp. 92–111.
- [26] Hao Ye, Geoffrey Ye Li, and Bing-Hwang Juang. “Channel Agnostic End-to-End Learning based Communication Systems with Conditional GAN”. In: *IEEE Global Communications Conference*. 2018.

- [27] Hao Ye, Geoffrey Ye Li, and Bing-Hwang Juang. “Power of Deep Learning for Channel Estimation and Signal Detection in OFDM Systems”. In: *IEEE Wireless Communications Letters*. 2018, pp. 114–117.