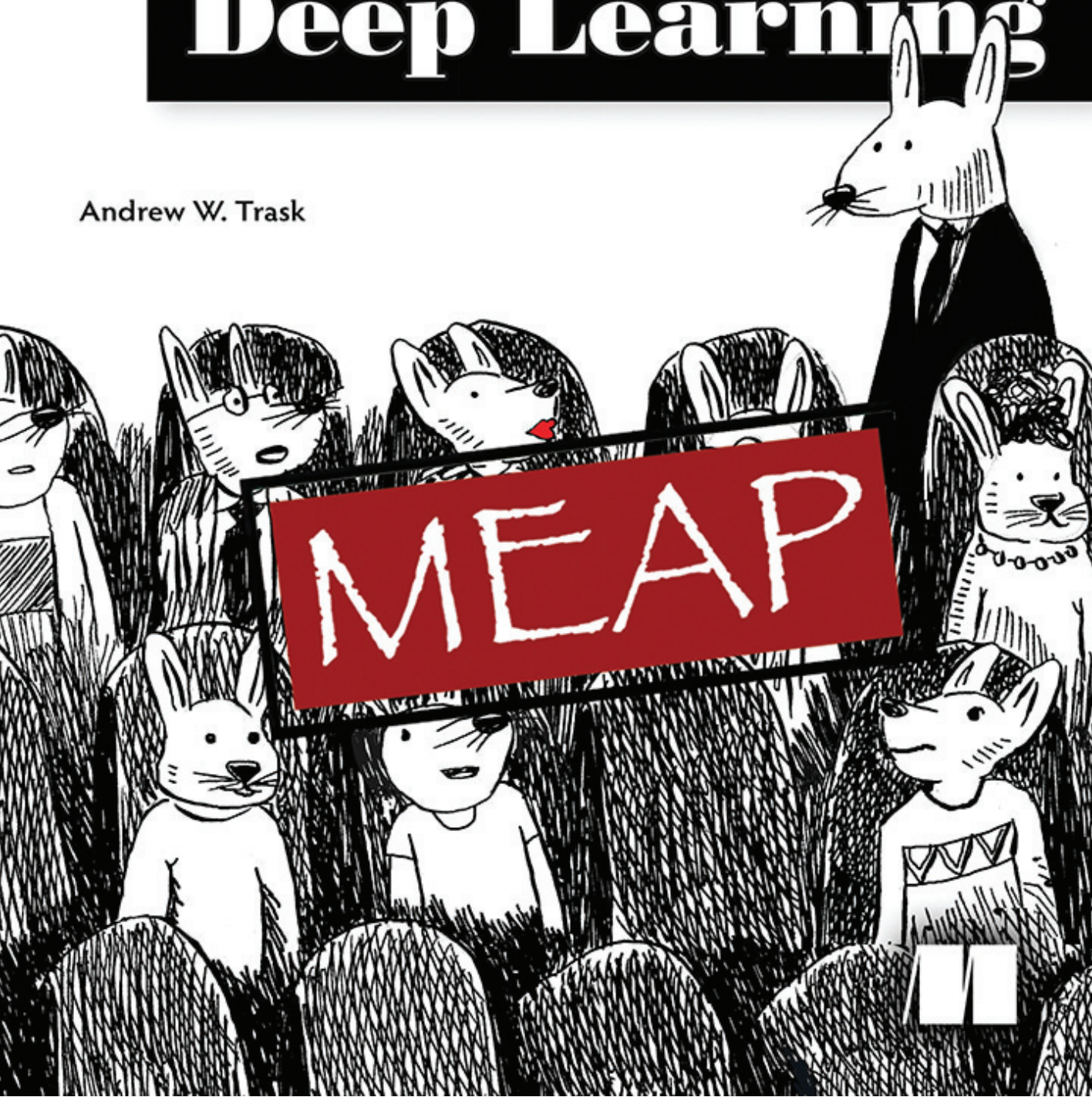


grokking

Deep Learning

Andrew W. Trask





**MEAP Edition
Manning Early Access Program
Grokking Deep Learning
Version 11**

Copyright 2018 Manning Publications

For more information on this and other Manning titles go to
www.manning.com

welcome

Thank you so much for purchasing *Grokking Deep Learning*. This book will teach you the fundamentals of Deep Learning from an intuitive perspective, so that you can understand *how machines learn using Deep Learning*. This book is not focused on learning a framework such as Torch, TensorFlow, or Keras. Instead, it is focused on teaching you the Deep Learning *methods* behind well known frameworks. Everything will be built from scratch using only Python and numpy (a matrix library). In this way, you will understand every detail that goes into training a neural network, not just how to use a code library. You should consider this book a prerequisite to mastering one of the major frameworks.

There are many other resources for learning Deep Learning. I'm glad that you came to this one, as I have intentionally written it with what I believe is the lowest barrier to entry possible. No knowledge of Linear Algebra, Calculus, Convex Optimization, or even Machine Learning is assumed. Everything from these subjects that is necessary to understand Deep Learning will be explained as we go. If you have passed high school mathematics and hacked around in Python, you're ready for this book, and when you complete this book, you will be ready to master a major deep learning framework like Torch, TensorFlow, or Keras.

Finally, as this is the MEAP, if there is any point in these first few chapters that something does not make sense, it is my hope that you would tweet your questions to me @iamtrask. I would be happy to help, and more importantly, I want to know if any section of the book is not fulfilling my personal commitment to the lowest barrier to entry possible so that I can adjust it for the final published work. Please, don't hesitate to reach out if you have questions.

These first few chapters will be walking you from a general introduction to Deep Learning all the way through to building your first working neural network. In these chapters, you will get a firm grasp on the philosophy behind how machines can learn the world you present to them. It's an exciting thing to see happen, and perhaps even more exciting, you will understand every nook and cranny of what makes this learning possible.

It is an honor to have your time and attention.

—Andrew Trask

brief contents

PART 1: NEURAL NETWORK BASICS

- 1 Introducing Deep Learning: Why you should learn it*
- 2 Fundamental Concepts: How do machines learn?*
- 3 Introduction to Neural Prediction: Forward Propagation*
- 4 Introduction to Neural Learning: Gradient Descent*
- 5 Learning Multiple Weights at a Time: Generalizing Gradient Descent*
- 6 Building Your First "Deep" Neural Network: Introduction to Backpropagation*
- 7 How to Picture Neural Networks: In Your Head and on Paper*
- 8 Learning Signal and Ignoring Noise: Introduction to Regularization & Batching*
- 9 Modeling Probabilities and Non-Linearities: Activation Functions*

PART 2: ADVANCED LAYERS AND ARCHITECTURES

- 10 Neural Networks that Understand Edges and Corners*
- 11 Neural Networks that Understand Language: King - Man + Woman == ?*
- 12 Neural Networks that Write like Shakespeare: Recurrent Layers for Variable Length Data*
- 13 Neural Networks that Read & Answer Questions*
- 14 Building a Neural Network that Destroys You in Pong*
- 15 Where to go from here*

IN THIS CHAPTER •

Why you should learn deep learning

Why you should read this book

What you need to get started

“ Do not worry about your difficulties in mathematics.
I can assure you mine are still greater.
— ALBERT EINSTEIN ”

Welcome to Grokking Deep Learning

You're about to learn some of the most valuable skills of the century!

I'm very excited that you're here! You should be too! Deep Learning represents an exciting intersection of Machine Learning and Artificial Intelligence, and a very significant disruption to society and industry. The methods discussed in this book are changing the world all around you. From optimizing the engine of your car to deciding which content you view on social media, it's everywhere, it's powerful, and fortunately, it's fun!

Why you should learn Deep Learning

It is a powerful tool for the *incremental automation of intelligence*.

From the beginning of time, humans have been building better and better tools to understand and control the environment around us. Deep Learning is today's chapter in this story of innovation. Perhaps what makes this chapter so compelling is that this field is more of a *mental* innovation than a *mechanical* one. Much like its sister fields in Machine Learning, Deep Learning seeks to *automate intelligence* bit by bit, and in the past few years it has achieved enormous success and progress in this endeavor, exceeding previous records in Computer Vision, Speech Recognition, Machine Translation, and many other tasks. This is particularly extraordinary given that Deep Learning seems to use *largely the same brain-inspired algorithm* (Neural Networks) for achieving these accomplishments across a vast number of fields. Even though Deep Learning is still an actively developing field with many challenges, recent developments have led to tremendous excitement that perhaps we have in fact discovered more than just a great tool, but a window into our own minds as well.

Deep Learning has the potential for significant automation of *skilled labor*.

There is a substantial amount of hype around the potential impacts of Deep Learning if the current trend of progress is extrapolated at varying speeds. While many of these predictions are over-zealous, there is one that I think merits your consideration: job displacement. I think that this claim stands out from the rest for no other reason than if Deep Learning's innovations stopped *today*, there would already be an incredible impact on skilled labor around the globe. Call center operators, taxi drivers, and low-level business analysts are compelling examples where Deep Learning can provide a low-cost alternative. Fortunately, the economy doesn't turn on a dime, but in many ways we are already past the point of concern with the current power of the technology. It is my hope that you (and people you know) will be enabled by this book to transition from perhaps one of the industries facing disruption into an industry ripe with growth and prosperity: Deep Learning.

It's fun and incredibly creative. You will discover much about what it is to be human by trying to simulate intelligence and creativity.

Personally, I got into Deep Learning because it's fascinating. It's an amazing intersection between man and machine. Unpacking exactly what it means to think, to reason, and to create is enlightening, engaging, and for me it's quite inspiring. Consider having a dataset filled with every painting ever painted, and then using that to teach a machine how to paint like Monet. Insanely, it's possible, and it's mind-bogglingly cool to see how it works.

Will this be hard to learn?

How hard will you have to work before there is a "fun" payoff?

This is my favorite question. My definition of a "fun" payoff is the experience of witnessing something that I built *learning*. There's just something amazing about seeing a creation of your hands do something like that. If you also feel this way, then the answer is simple. A few pages into Chapter 3, you will create your first neural network. The only work involved between now and then is reading the pages between here and there.

After Chapter 3, you may be interested to know that the *next* fun payoff occurs after you have memorized a small snippet of code and proceeded to read to the midway of Chapter 4. Each chapter will continue to work this way. Memorize a small code segment from the previous chapter, read the next chapter, and then experience the payoff of a new learning neural network.

Why you should read this book

It has a uniquely low barrier to entry.

The reason you should read this book is the same reason I'm writing it. I don't know of another resource (book, course, large blog series) that teaches Deep Learning **without assuming advanced knowledge of mathematics** (i.e. college degree in a mathy field). Don't get me wrong, there are really good reasons for teaching it using math. Math is, after all, a language. It is certainly more **efficient** to teach Deep Learning using this language, but I don't think it's absolutely necessary to assume advanced knowledge of math in order to become a skilled, knowledgeable practitioner who has a firm understanding of the "how" behind Deep Learning. So, why should you learn Deep Learning using this book? I'm going to assume you have a High School level background in math (and that it's rusty), and *explain everything else you need to know as we go along*. Remember multiplication? Remember x-y graphs (the square with lines on it)? Awesome! You'll be fine.

Why you should read this book (cont.)

To help you understand what's *inside* a framework (Torch, TensorFlow, etc.).

There are two major groups of Deep Learning educational material (books, courses, etc.). One group is focused around how to use popular frameworks and code libraries such as Torch, Tensorflow, Keras, and others. The other group is focused around teaching Deep Learning itself, otherwise known as the *science under the hood* of these major frameworks. Ultimately, learning about *both* is important. It's like if you want to be a NASCAR driver. You need to learn both about the particular model of car you're driving (the framework), and about driving itself (the science/skill). However, just learning about a framework is like learning about the pros and cons of a Generation-6 Chevrolet SS before you know what a stick shift is. This book is about teaching you what *Deep Learning* is so that you can then be prepared to learn a framework.

All math related material will be backed by intuitive *analogies*.

Whenever I encounter a math formula in the wild, I take a two-step approach. The first is to translate its methods into an intuitive *analogy* to the real world. I almost never just take a formula at face value. I break it into *parts*, each with a story of its own. That will be the approach of this book as well. Anytime we encounter a math concept, I'll offer an alternative *analogy* for what the formula is actually doing.

{ "Everything should be made as simple as possible, but no simpler"
- Albert Einstein }

Everything after the introduction chapters is "project" based.

If there is one thing I hate when learning something new, it is when I have to question whether or not what I'm learning is useful/relevant. If someone is teaching me everything there is to know about a hammer without actually taking my hand and helping me drive in a nail, then they're not really teaching me how to use a hammer. I know that there are going to be dots that weren't connected, and if I was thrown out into the real world with a hammer, a box of nails, and a bunch of 2x4s, I'm going to have to do some guesswork.

This book is about giving you the wood, nails, and a hammer *before* telling you about what they do. Each lesson is about picking up the tools and building stuff with them, explaining how stuff works along the way. In this way, you don't leave with a list of facts about the various Deep Learning tools we'll work with, you leave with the ability to use them to solve problems. Furthermore, you will understand the most important part: when and

why each tool is appropriate for each problem you want to solve. It is with this knowledge that you will be empowered to pursue a career in research and/or industry.

What you need to get started

Install Jupyter Notebook and the NumPy python library

My absolute favorite place to work is a Jupyter Notebook. One of the most important parts of learning Deep Learning (for me), is the ability to stop a network while it's training and tear apart absolutely every piece to see what it looks like. This is something that Jupyter Notebook is incredibly useful for. As for Numpy, perhaps the most compelling case for why this book leaves nothing out is that we'll only be using a single matrix library. In this way, you will understand **how** everything works, not just how to call a framework. This book teaches Deep Learning from absolute scratch, soup to nuts. Installation instructions for these two tools can be found at (<http://jupyter.org/>) for Jupyter and (<http://numpy.org>) for NumPy. I will be building these examples in Python 2.7, but will test them for Python 3 as well. For easy installation, I also recommend the Anaconda framework (<https://docs.continuum.io/anaconda/install>).

Pass High School Mathematics

There are some mathematical assumptions that are simply out of depth for this book, but the goal of this book is to teach Deep Learning only assuming you understand basic algebra.

Find a personal problem you are interested in

This might seem like an optional "need" to get started. I guess it could be, but seriously, I highly, highly recommend finding one. Everyone I know who has become successful at this stuff had some sort of problem they were trying to solve. Learning Deep Learning was just a "dependency" to solving some other interesting task. For me, it was using Twitter to predict the stock market. It's just something that I thought was really fascinating. It's what drove me to sit down and read the next chapter and build the next prototype. And as it turns out, this field is **so new**, and is changing **so fast**, that if you spend the next couple of years chasing one project with these tools, you will find yourself becoming one of the leading experts in that *particular problem* faster than you might think. For me, chasing this idea took me from barely knowing anything about programming to a research grant at a hedge fund applying what I learned in around 18 months! For Deep Learning, having a problem you're fascinated with that involves using one dataset to predict another is the key catalyst! Go find one!

You'll probably need some Python knowledge

Python is my teaching library of choice, but I'll provide a few others online.

Python is an amazingly intuitive language. I think it just might be the most widely adopted and intuitively readable language yet constructed. Furthermore, the Python community has a passion for simplicity that can't be beat. For these reasons, I want to stick with Python for all of the examples (Python 2.7 is what I'm working in). On this book's Github, I'll provide all of the examples in a variety of other languages as well, but for the in-page explanations, we're going to use Python.

How much coding experience should you have?

Scan through the Python Codecademy course (<https://www.codecademy.com/learn/python>). If you can read through the table of contents and feel comfortable with the terms mentioned, you're all set! If not, then just take the course and come back when you're done! It's designed to be a beginner course and it's very well crafted.

Conclusion and Primer for Chapter 2

If you've got your Jupyter Notebook in-hand and feel comfortable with the basics of Python, you're ready for the next chapter! As a heads up, Chapter 2 is the last chapter that will be mostly dialogue based (without building something). It's just designed to give you an awareness of the high level vocabulary, concepts, and fields in Artificial Intelligence, Machine Learning, and most importantly, Deep Learning.

IN THIS CHAPTER •

- What are Deep Learning, Machine Learning, and Artificial Intelligence?
- What is a Parametric Model?
- What is a Non-Parametric Model?
- What is Supervised Learning?
- What is Unsupervised Learning?
- How can machines learn?

“

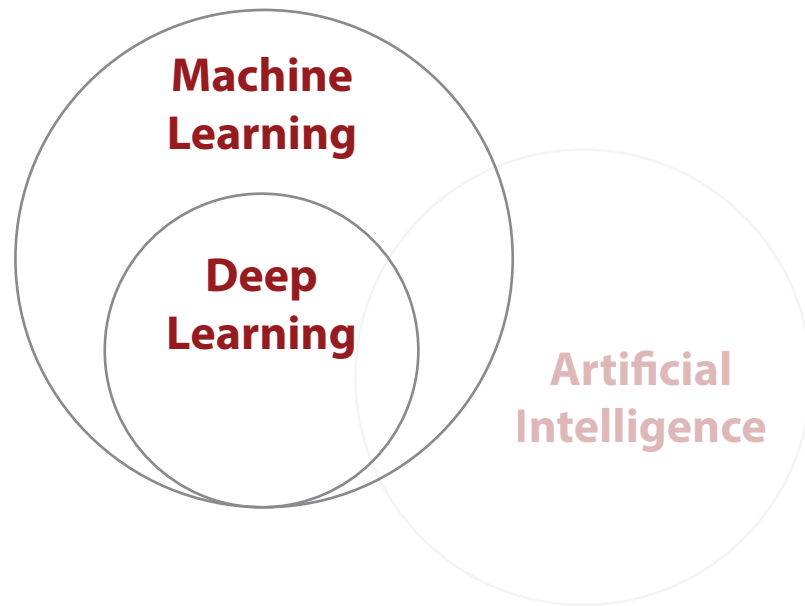
Machine Learning will cause every successful IPO win in 5 years.
— ERIC SCHMIDT (Google Chairman, 2016)

”

What is Deep Learning?

Deep Learning is a subfield of methods for Machine Learning

Deep Learning is a subset of Machine Learning, which is a field dedicated to the study and development of machines that can learn (sometimes with the goal of eventually attaining General Artificial Intelligence). In industry, Deep Learning is used to solve practical tasks in a variety of fields such as Computer Vision (Image), Natural Language Processing (Text), and Automatic Speech Recognition (Audio). In short, Deep Learning is a subset of *methods* in the Machine Learning toolbox, primarily leveraging **Artificial Neural Networks**, which are a class of algorithm loosely inspired by the human brain.



Notice in the figure above that not all of Deep Learning is focused around pursuing Generalized Artificial Intelligence (i.e. sentient machines like in the movies). In fact, many applications of this technology are applied to solve a wide variety of problems in industry. This book seeks to focus on learning the fundamentals of Deep Learning behind both cutting edge research and industry, helping to prepare you for either.

What is Machine Learning?

"A field of study that gives computers the ability to learn without being explicitly programmed"

- Arthur Samuel

Given that Deep Learning is a subset of Machine Learning, what is Machine Learning? Most generally, it is what its name implies. Machine Learning is a subfield of Computer Science wherein *machines learn* to perform tasks for which they were *not explicitly programmed*. In short, machines observe a pattern and attempt to imitate it in some way which can be either direct or indirect **imitation**.

machine learning \sim monkey see monkey do

I mention direct and indirect imitation as a parallel to the two main types of machine learning, **supervised** machine learning and **unsupervised** machine learning. Supervised machine learning is the direct imitation of a pattern between two datasets. It is always attempting to take an input dataset and transform it into an output dataset. This can be an incredibly powerful and useful capability. Consider the following examples: (**input** datasets in bold and *output* datasets in italic)

- Using the **pixels** of an image to detect the *presence or absence of a cat*.
- Using the **movies you've liked** to predict more *movies you may like*.
- Using someone's **words** to predict whether they are *happy* or *sad*.
- Using weather sensor **data** to predict the *probability of rain*.
- Using car engine **sensors** to predict the optimal tuning *settings*.
- Using news **data** to predict tomorrow's stock *price*.
- Using an input **number** to predict a *number* double its size.
- Using a raw **audio file** to predict a *transcript* of the audio.

These are all supervised machine learning tasks. In all cases the machine learning algorithm is attempting to imitate the pattern between the two datasets in such a way that it can **use one dataset to predict the other**. For any example above, imagine if you had the power to predict the *output* dataset given only the **input** dataset. This would be profound.

Supervised Machine Learning

Supervised Learning transforms one dataset into another.

Supervised Learning is a method for transforming one dataset into another. For example, if we had a dataset of "Monday Stock Prices" which recorded the price of every stock on every Monday for the past 10 years, and a second dataset of "Tuesday Stock Prices" recorded over the same time period, a supervised learning algorithm might try to use one to predict the other.



If we successfully trained our supervised machine learning algorithm on 10 years of Mondays and Tuesdays, then we could predict the stock price of any Tuesday in the future given the stock price on the immediately preceeding Monday. I encourage you to stop and consider this for a moment.

Supervised Machine Learning is the bread and butter of applied Artificial Intelligence (i.e. "Narrow AI"). It is useful for taking *what we do know* as input and quickly transforming it into **what we want to know**. This allows supervised machine learning algorithms to extend human intelligence and capabilities in a seemingly endless number of ways.

The majority of work leveraging machine learning results in the training of a supervised classifier of some kind. Even unsupervised machine learning (which we will learn more about in a moment) is typically done to aid in the development of an accurate supervised machine learning algorithm.



For the rest of this book, we will be creating algorithms that can take input data that is observable, recordable, and by extension **knowable** and transform it into valuable output data that requires logical analysis. This is the power of supervised machine learning.

Unsupervised Machine Learning

Unsupervised Learning groups your data.

Unsupervised learning shares a property in common with supervised learning. It transforms one dataset into another. However, the dataset that it transforms into is **not previously known or understood**. Unlike supervised learning, there is no "right answer" that we're trying to get the model to duplicate. We just tell an unsupervised algorithm to "find patterns in this data and tell me about them".

For example, *clustering a dataset into groups* is a type of unsupervised learning. "Clustering" transforms your sequence of *datapoints* into a sequence of *cluster labels*. If it learns 10 clusters, it's common for these labels to be the numbers 1-10. Each datapoint will get assigned to a number based on which cluster it is in. Thus, your dataset turns from a bunch of datapoints into a bunch of labels. Why are the labels numbers? The algorithm doesn't tell us what the clusters are. How could it know? It just says "Hey scientist!... I found some structure. It looks like there are some groups in your data. Here they are!".



I have good news! This idea of clustering is something you can reliably hold onto in your mind as the definition of unsupervised learning. Even though there are many forms of unsupervised learning, *all forms of unsupervised learning can be viewed as a form of clustering*. We will discover more on this later in the book.

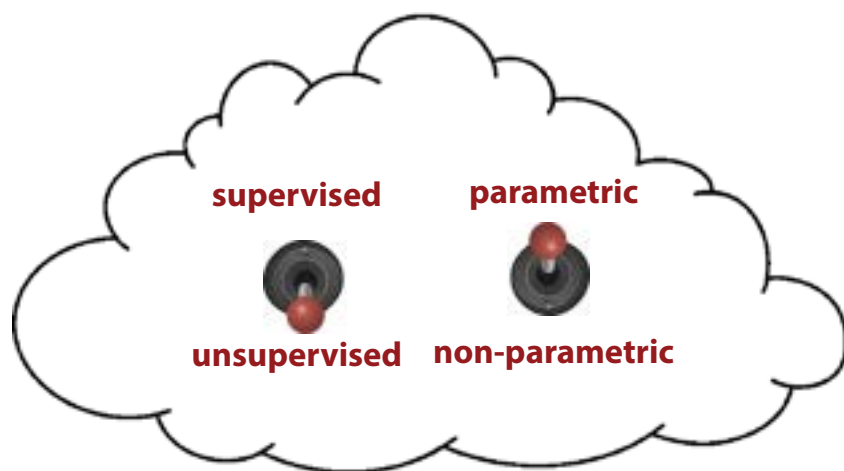


Check out the example above. Even though the algorithm didn't tell us what the clusters are named, can you figure out how it clustered the words? (1 == cute & 2 == delicious) Later, we will unpack how other forms of unsupervised learning are also just a form of clustering and why these clusters are useful for supervised learning.

Parametric vs Non-Parametric Learning

Oversimplified: Trial and error learning versus counting and probability

The last two pages divided all of our machine learning algorithms into two groups, supervised and unsupervised. Now, we're going to discuss another way to divide the same machine learning algorithms into two groups, parametric and non-parametric. So, if we think about our little machine learning cloud, it has two settings:



As you can see, we really have four different types of algorithm to choose from. An algorithm is either unsupervised or supervised and it is either parametric or non-parametric. Whereas the previous section on supervision is really about the **type of pattern** being learned, parametricism is about the way the learning is **stored** and often by extension, the **method for learning**. First, let's look at the formal definition for parametricism vs non-parametricism. For the record, there is still some debate around the exact difference.

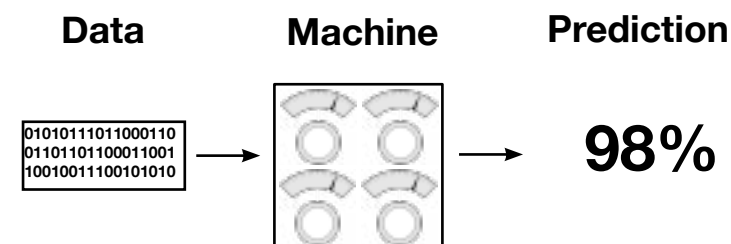
A parametric model is characterized by having a *fixed* number of parameters whereas a non-parametric model's number of parameters is *infinite* (determined by data).

As an example, let's say the problem was to fit a square peg into the correct (square) hole. Some humans (such as babies) just jam it into all the holes until it fits somewhere (parametric). A teenager, however, might just count the number of sides (4) and then search for the hole with an equal number (non-parametric). Parametric models tend to use "trial and error", where non-parametric models tend to "count". Let's look closer.

Supervised Parametric Learning

Oversimplified: Trial and error learning using knobs

Supervised parametric learning machines are machines with a fixed number of knobs (that's the parametric part), wherein learning occurs by turning the knobs. **Input data** comes in, is processed based on the angle of the knobs, and is transformed into a *prediction*.



Learning is accomplished by turning the knobs to different angles. If we're trying to predict the probability that the Red Socks will win the World Series, then this model would first take data (such as sports stats like win/loss record or average number of toes) and make a prediction (such as 98% chance). Next, the model would observe whether or not the Red Socks actually won. After it knew whether they won, our learning algorithm would **update the knobs** to make a more accurate prediction the next time it sees the **same/similar input data**.

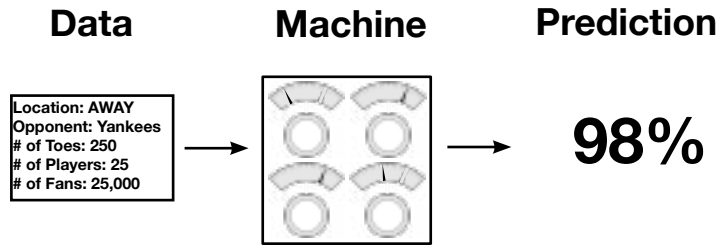
Perhaps it would "turn up" the "win/loss record" knob if the team's win/loss record was a good predictor. Inversely, it might turn down the "average number of toes" knob if that datapoint wasn't a good predictor. This is how parametric models learn!

Note that the entirety of what the model has learned can be captured in the positions of the knobs at any given time. One can also think of this type of learning model as a search algorithm. We are "searching" for the appropriate knob configuration by trying configurations, adjusting them, and retrying.

Note further that the notion of trial and error isn't the formal definition, but it is a very common (with exceptions) property to parametric models. When there is an arbitrary (but fixed) number of knobs to turn, then it requires some level of searching to find the optimal configuration. This is in contrast to non-parametric learning, which is often "count" based and (more or less) "adds new knobs" when it finds something new to count. Let's break down supervised parametric learning into its three steps.

Step 1: Predict

To illustrate supervised parametric learning, let's continue with our sports analogy where we're trying to predict whether or not the Red Socks will win the World Series. The first step, as mentioned, is to gather sports statistics, send them through our machine, and make a prediction on the probability that the Red Socks will win.



Step 2: Compare to Truth Pattern

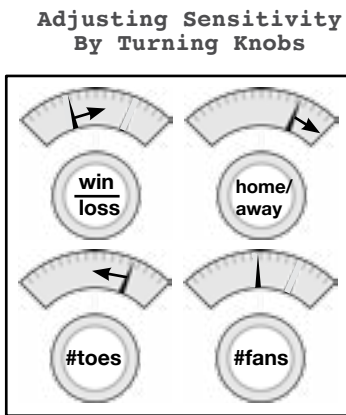
The second step is to compare the prediction (98%) with the pattern that we care about (whether the Red Socks won). Sadly, they lost, so our comparison is:

Pred: 98% > Truth: 0%

This step simply recognizes that if our model had predicted 0%, it would have perfectly predicted the upcoming loss of the team. We want our machine to be accurate, which takes us to Step 3.

Step 3: Learn the Pattern

This step adjusts the knobs by studying both how **much** the model missed (98%) and what the input data **was** (sports stats) at the time of prediction. It then turns the knobs to make a more accurate prediction given the input data. In theory, the next time it saw the same sports stats, the prediction would be lower than 98%. Note that each knob represents the *prediction's sensitivity to different types of input data*. That's what we're changing when we "learn".



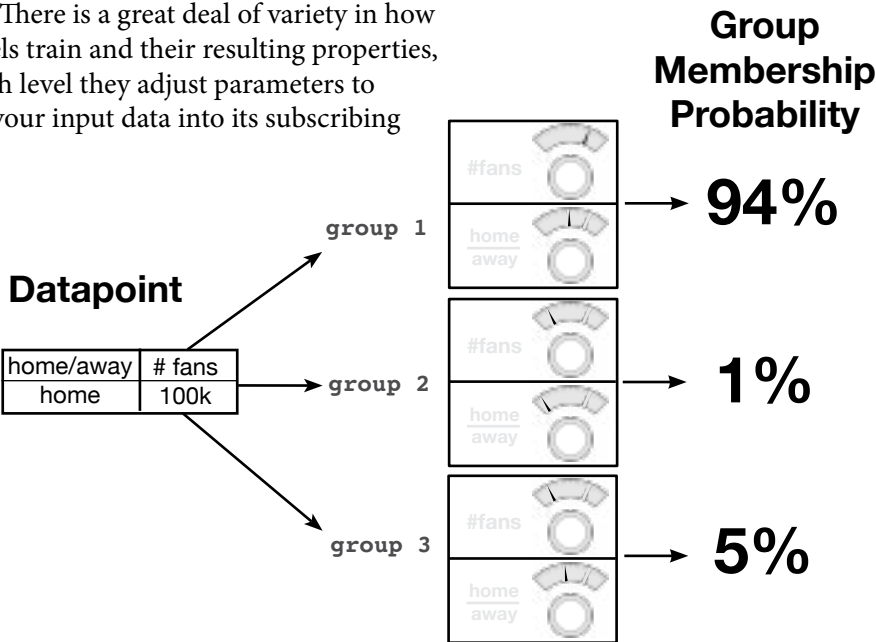
Unsupervised Parametric Learning

Unsupervised parametric learning leverages a very similar approach. Let's walk through the steps at a high level. Remember that unsupervised learning is all about grouping your data. Unsupervised *parametric* learning uses knobs to group your data. However, in this case, it usually has several knobs for each group, each that maps your input data's affinity to that particular group (with exception and nuance, this is a high level description). Let's look at an example where we assume we want to divide our data into three groups.

In the dataset on the right, I have identified three clusters in the data that we might want our parametric model to find. I identified them via formatting as **group 1**, group 2, and *group 3*. Let's propagate our first datapoint through a trained unsupervised model below. Notice that it maps most strongly to **group one**.

home/away	# fans
home	100k
away	50k
home	100k
home	99k
away	50k
away	10k
away	11k

Each group's machine attempts to transform the input data to a number between 0 and 1, telling us the *probability that the input data is a member of that group*. There is a great deal of variety in how these models train and their resulting properties, but at a high level they adjust parameters to transform your input data into its subscribing group(s).



Non-Parametric Learning

Oversimplified: Counting based methods

Non-Parametric learning is a class of algorithm wherein the number of parameters is based on data (instead of pre-defined). This lends itself to methods that generally count in one way or another, thus increasing the number of parameters based on the number of items being counted within the data. In the supervised setting, for example, a non-parametric model might count the number of times a particular color of streetlight causes cars to "go". After counting only a few examples, this model would then be able to predict that *middle* lights always (100%) cause cars to "go" and *right* lights only sometimes (50%) cause cars to "go".



Notice that this model would have 3 parameters, 3 counts indicating the number of times each colored light turned on and cars "go" (perhaps divided by the number of total observations). If there had been 5 lights, there would have been 5 counts (5 parameters). What makes this simple model *non-parametric* is this trait wherein the number of parameters changes based on the data (in this case, the number of lights). This is in contrast to parametric models, which start with a set number of parameters and, more importantly, can have more or less parameters purely at the discretion of the scientist training the model (irrespective of data).

A close eye might question this idea. Our parametric model from before seemed to have a knob for each input datapoint. In fact, most parametric models still have to have some sort of *input* that is based on the number of classes in the data. Thus you can see that there is a bit of *grayness* between parametric and non-parametric algorithms. Even parametric algorithms still are somewhat influenced by the number of classes in the data, even if they are not explicitly counting patterns.

Perhaps this also illuminates that *parameters* is actually a very generic term, referring only to the set of numbers used to model a pattern (without any limitation on how those numbers are used). Counts are parameters. Weights are parameters. Normalized variants of counts or weights are parameters. Correlation coefficients can be parameters. It's simply referring to the set of numbers used to model a pattern. As it happens, Deep Learning is a class of parametric models. We won't be discussing non-parametric models further in this book, but they are a very interesting and powerful class of algorithm.

Conclusion

In this chapter, we have gone a level deeper into the various flavors of Machine Learning. We have learned that a Machine Learning algorithm is either Supervised or Unsupervised and either Parametric or Non-Parametric. Furthermore, we have explored exactly what makes these 4 different groups of algorithms distinct. We have learned that Supervised Machine Learning is a class of algorithm where we learn to predict one dataset given another and that Unsupervised Learning generally groups a single dataset into various kinds of clusters. We learned that Parametric algorithms have a fixed number of *parameters* and that Non-Parametric algorithms adjust the number of parameters they have based on the dataset.

Deep Learning leverages Neural Networks to perform both supervised and unsupervised prediction. Up until now, we have stayed at mostly a conceptual level, gaining our bearings on the field as a whole and our place in it. In the next chapter, we will be building our first neural network, and all subsequent chapters will be *project based*. So, pull out your Jupyter notebook and let's jump right in!

IN THIS CHAPTER •

- A Simple Network Making a Prediction
- What is a Neural Network and what does it do?
- Making a Prediction with Multiple Inputs
- Making a Prediction with Multiple Outputs
- Making a Prediction with Multiple Inputs and Outputs
- Predicting on Predictions

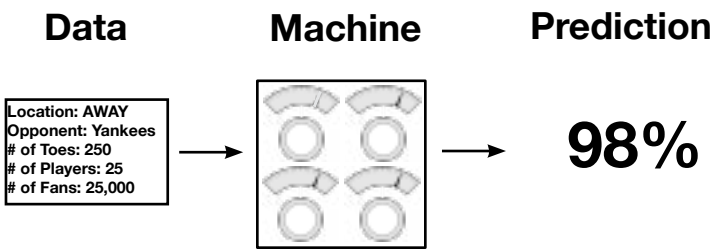
“ I try not to get involved in the business of prediction. It's a quick way to look like an idiot. ”

— WARREN ELLIS

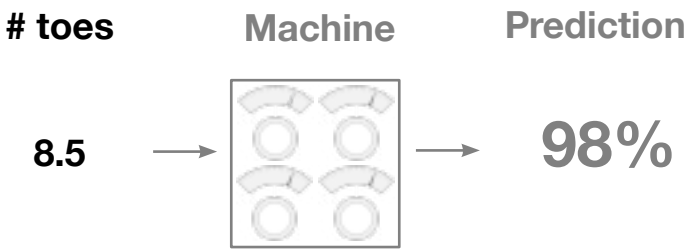
Step 1: Predict

This chapter is about "Prediction"

In the previous chapter, we learned about the paradigm: "Predict, Compare, Learn". In this chapter, we will dive deep into the first step: "Predict". You may remember that the Predict step looks a lot like this.

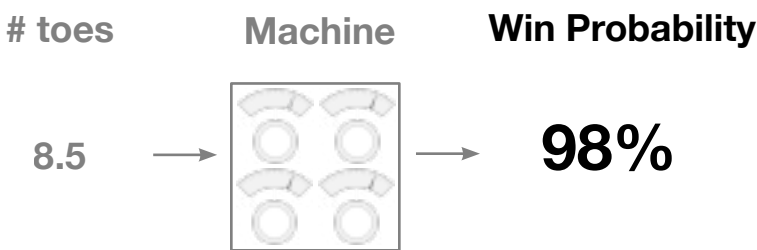


In this chapter, we're going to learn more about what these 3 different parts of a neural network prediction really look like under the hood. Let's start with the first one, the Data. In our first neural network, we're going to predict one datapoint at a time, like so:

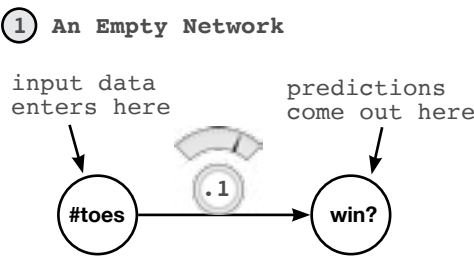


Later on, we will find that the "number of datapoints at a time" that we want to process will have a significant impact on what our network looks like. You might be wondering, "how do I choose how many datapoints to propagate at a time?" The answer to this question is based on whether or not you think the neural network can be accurate with the data you give it. For example, if I'm trying to predict whether or not there's a cat in a photo, I definitely need to show my network all the pixels of an image at once. Why? Well, if I only sent you one pixel of an image, could you classify whether the image contained a cat? Me neither! (That's a general rule of thumb by the way. Always present enough information to the network, where "enough information" is defined loosely as how much a human might need to make the same prediction.)

Let's skip over the network for now. As it turns out, we can only create our network once we understand the shape of our input and output datasets (for now, shape means "number of columns" or "number of datapoints we're processing at once"). For now, we're going to stick with the "single-prediction" of "likelihood that the baseball team will win".



Ok, so now that we know that we want to take one input datapoint and output one prediction, we can create our neural network. Since we only have one input datapoint and one output datapoint, we're going to build a network with a single knob mapping from the input point to the output. Abstractly these "knob"s are actually called "weight"s, and we will refer to them as such from here on out. So, without further ado, here's our first neural network with a single weight mapping from our input "#toes" to output "win?"



As you can see, with one weight, this network takes in one datapoint at a time (average number of toes on the baseball team) and outputs a single prediction (whether or not it thinks the team will win).

A Simple Neural Network Making a Prediction

Let's start with the simplest neural network possible.

1

An Empty Network

input data enters here

predictions come out here

#toes

.1

win?

```
weight = 0.1
def neural_network(input, weight):
    prediction = input * weight
    return prediction
```

2

Inserting One Input Datapoint

input data (#toes)

8.5

.1

```
number_of_toes = [8.5, 9.5, 10, 9]
input = number_of_toes[0]
pred = neural_network(input,weight)
print(pred)
```

3

Multiplying Input By Weight

(8.5 * 0.1 = 0.85)

8.5

.1

```
def neural_network(input, weight):
    prediction = input * weight
    return prediction
```

4

Depositing Prediction

prediction

8.5

.1

0.85

```
number_of_toes = [8.5, 9.5, 10, 9]
input = number_of_toes[0]
pred = neural_network(input,weight)
```

What is a Neural Network?

This is a neural network.
Open up a Jupyter Notebook and run the following:

the network

weight = 0.1

def neural_network(input, weight):

prediction = input * weight

return prediction

how we use the network to predict something

number_of_toes = [8.5, 9.5, 10, 9]

input = number_of_toes[0]

pred = neural_network(input,weight)

print(pred)

You just made your first neural network and used it to predict! Congratulations! The last line prints the prediction (`pred`). It should be 0.85. So what is a neural network? For now, it's one or more *weights* which we can multiply by our input data to make a prediction.

What is input data?
It's a number that we recorded in the real world somewhere. It's usually something that is easily knowable, like today's temperature, a baseball player's batting average, or yesterday's stock price.

What is a prediction?
A prediction is what the neural network tells us *given our input data* such as "given the temperature, it is **0%** likely that people will wear sweatsuits today" or "given a baseball player's batting average, he is **30%** likely to hit a home run" or "given yesterday's stock price, today's stock price will be **101.52**".

Is this prediction always right?
No. Sometimes our neural network will make mistakes, but it can learn from them. For example, if it predicts too high, it will adjust it's *weight* to predict lower next time and vice versa.

How does the network learn?
Trial and error! First, it tries to make a *prediction*. Then, it sees whether it was too high or too low. Finally, it changes the *weight* (up or down) to predict more accurately the next time it sees the same input.

What does this Neural Network do?

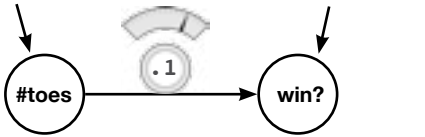
It multiplies the input by a weight. It "scales" the input by a certain amount.

On the previous page, we made our first prediction with a neural network. A neural network, in it's simplest form, uses the power of *multiplication*. It takes our input datapoint (in this case, 8.5) and *multiplies* it by our *weight*. If the *weight* is 2, then it would *double our input*. If the *weight* is 0.01, then it would *divide* the input by 100. As you can see, some *weight* values make the input *bigger* and other values make it *smaller*.

1 An Empty Network

input data enters here

predictions come out here




```
weight = 0.1

def neural_network(input, weight):
    prediction = input * weight
    return prediction
```

The interface for our neural network is really quite simple. It accepts an input variable as *information*, and a weight variable as *knowledge* and outputs a prediction. Every neural network you will ever see works this way. It uses the *knowledge* in the weights to interpret the *information* in the input data. Later neural networks will accept larger, more complicated input and weight values, but this same underlying premise will always ring true.

2 Inserting One Input Datapoint

input data (#toes)

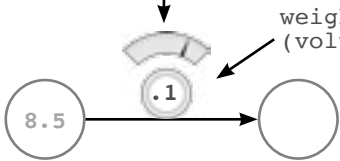


```
number_of_toes = [8.5, 9.5, 10, 9]
input = number_of_toes[0]
pred = neural_network(input,weight)
```

In this case, the "information" is the average number of toes on a baseball team before a game. Notice several things: The neural network does NOT have access to any information except *one* instance. If, after this prediction, we were to feed in `number_of_toes[1]`, it would not remember the prediction it made in the last timestep. A neural network only knows what you feed it as input. It forgets everything else. Later, we will learn how to give neural networks "short term memories" by feeding in multiple inputs at once.

3 Multiplying Input By Weight

(8.5 * 0.1 = 0.85)

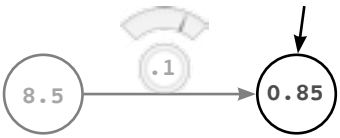


```
def neural_network(input, weight):
    prediction = input * weight
    return prediction
```

Another way to think about a neural network's weight is as a measure of *sensitivity* between the input of the network and its prediction. If the weight is very high, then even the tiniest input can create a really large prediction! If the weight is very small, then even large inputs will make small predictions. This *sensitivity* is very akin to **volume**. "Turning up the weight" amplifies our prediction relative to our input - weight is a volume knob!

4 Depositing Prediction

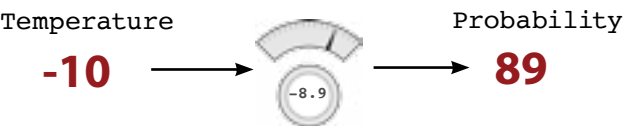
prediction



```
number_of_toes = [8.5, 9.5, 10, 9]
input = number_of_toes[0]
pred = neural_network(input,weight)
```

So in this case, what our neural network is really doing is applying a *volume knob* to our `number_of_toes` variable. In theory, this *volume knob* is able to tell us the likelihood that the team will win based on the average number of toes per player on our team. And this may or may not work. Truthfully, if the team had 0 toes, they would probably play terribly. However, baseball is much more complex than this. On the next page, we will present multiple pieces of information at the same time, so that the neural network can make more informed decisions.

Before we go, neural networks don't just predict positive numbers either, they can also *predict negative numbers*, and even take *negative numbers as input*. Perhaps you want to predict the "probability that people will wear coats today", if the temperature was -10 degrees Celsius, then a negative weight would predict a high probability that people would wear coats today.

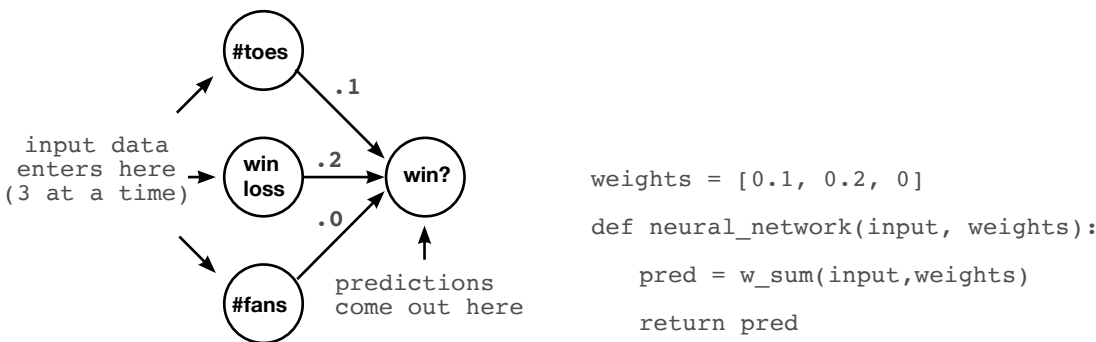


Making a Prediction with Multiple Inputs

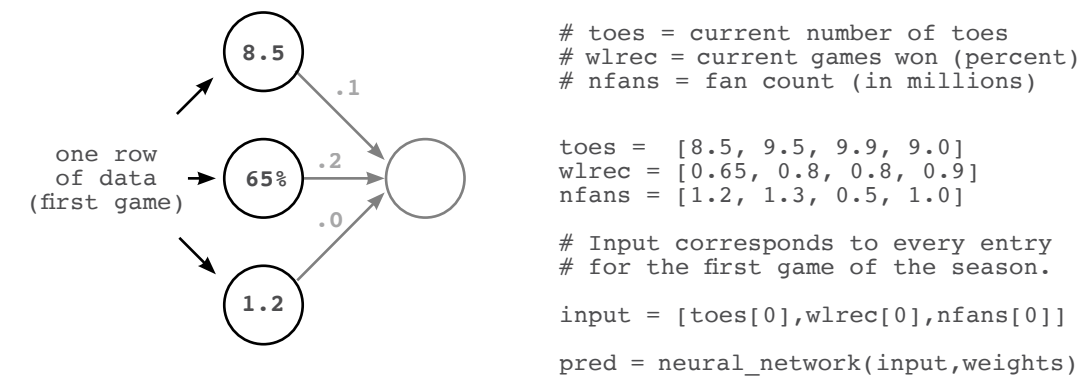
Neural Networks can combine intelligence from multiple datapoints.

Our last neural network was able to take one datapoint as input and make one prediction based on that datapoint. Perhaps you've been wondering, "is average # of toes really a very good predictor?... all by itself?" If so, you're onto something. What if we were able to give our network more information (at one time) than just the "average number of toes". It should, in theory, be able to make more accurate predictions, yes? Well, as it turns out, our network can accept multiple input datapoints at a time. See the prediction below!

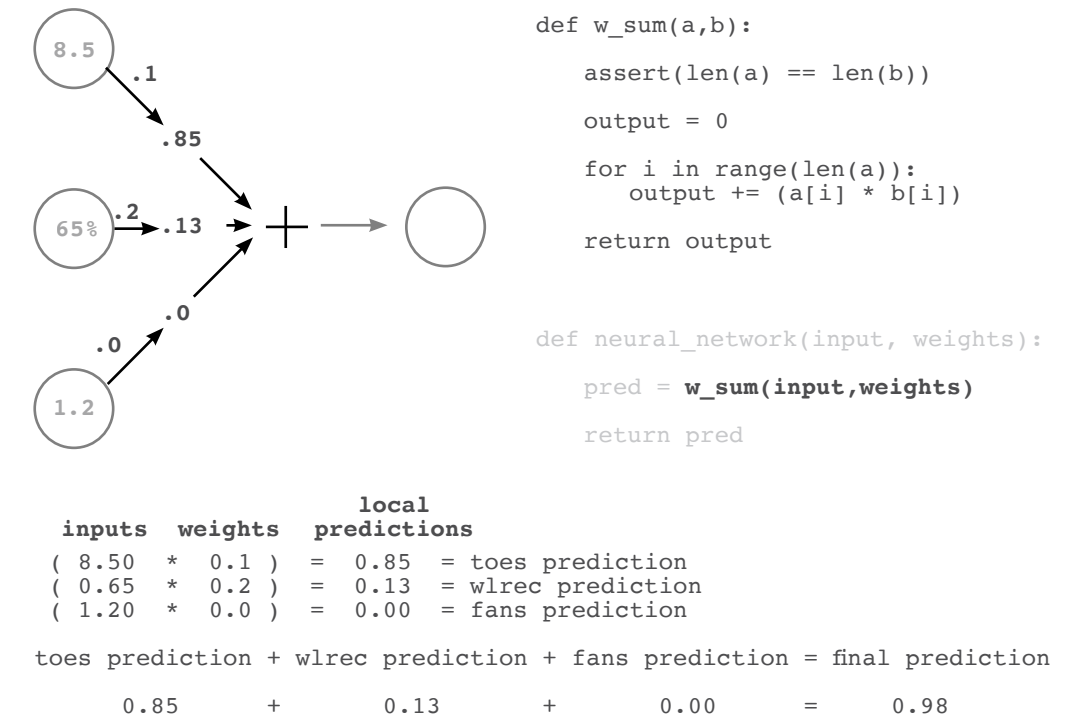
1 An Empty Network With Multiple Inputs



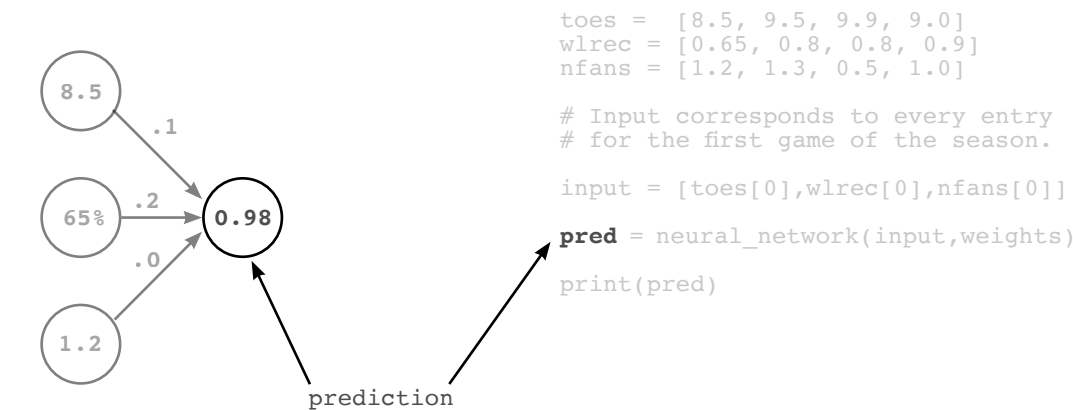
2 Inserting One Input Datapoint



3 Perform a Weighted Sum of Inputs



4 Deposit Prediction

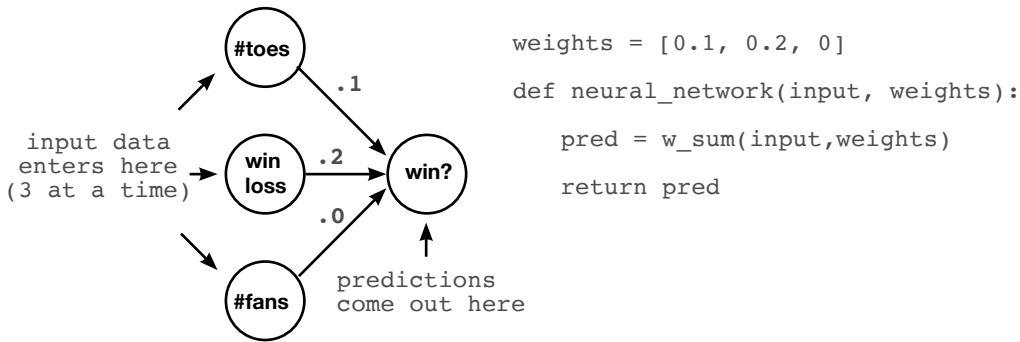


Multiple Inputs - What does this Neural Network do?

It multiplies 3 inputs by 3 knob_weights and sums them. This is a "weighted sum".

At the end of the previous section, we came to realize the limiting factor of our simple neural network, it is only a volume knob on one datapoint. In our example, that datapoint was the average number of toes on a baseball team. We realized that in order to make accurate predictions, we need to build neural networks that can *combine multiple inputs at the same time*. Fortunately, neural networks are perfectly capable of doing so.

1 An Empty Network With Multiple Inputs

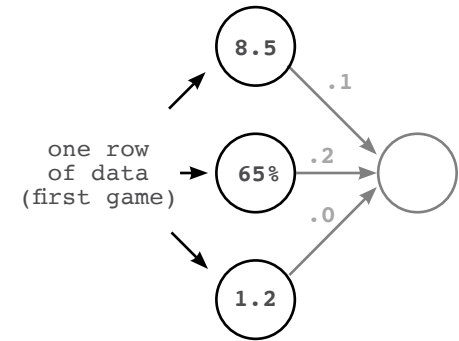


In this new neural network, we can accept *multiple inputs at a time* per prediction. This allows our network to combine various forms of information to make more well informed decisions. However, the fundamental mechanism for using our weights has not changed. We still take each input and run it through its own volume knob. In other words, we take each input and multiply it by its own weight. The new property here is that, since we have multiple inputs, we have to sum their respective predictions. Thus, we take each input, multiply it by its respective weight, and then sum all the local predictions together. This is called a "weighted sum of the input" or a "weighted sum" for short. Some also refer to this "weighted sum" as a "dot product" as we'll see.

A Relevant Reminder

The interface for our neural network is quite simple. It accepts an input variable as *information*, and a weight variable as *knowledge* and outputs a prediction.

2 Inserting One Input Datapoint



This new need to process multiple inputs at a time justifies the use of a new tool. This tool is called a **vector** and if you've been following along in your iPython notebook, you've already been using it. A vector is nothing other than a *list of numbers*. `input` is a vector and `weights` is a vector. Can you spot any more vectors in the code above (there are 3 more)?

As it turns out, vectors are incredibly useful whenever you want to perform operations involving groups of numbers. In this case, we're performing a weighted sum between two vectors (dot product). We're taking two vectors of equal length (`input` and `weights`), multiplying each number based on its position (the first position in `input` is multiplied by the first position in `weights`, etc.), and then summing the resulting output.

It turns out that whenever we perform a mathematical operation between two vectors of equal length where we "pair up" values according to their position in the vector (again... position 0 with 0, 1 with 1, and so on), we call this an **elementwise** operation. Thus "elementwise addition" sums two vectors. "Elementwise multiplication" multiplies two vectors.

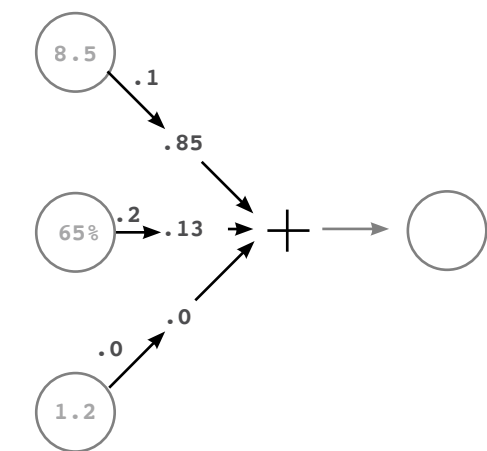
Challenge: Vector Math

Being able to manipulate vectors is a cornerstone technique for Deep Learning. See if you can write functions that perform the following operations:

```
def elementwise_multiplication(vec_a, vec_b)
def elementwise_addition(vec_a, vec_b)
def vector_sum(vec_a)
def vector_average(vec_a)
```

Then, see if you can use two of these methods to perform a dot product!

3 Perform a Weighted Sum of Inputs



```
def w_sum(a,b):
    assert(len(a) == len(b))
    output = 0
    for i in range(len(a)):
        output += (a[i] * b[i])
    return output

def neural_network(input, weights):
    pred = w_sum(input,weights)
    return pred
```

inputs	weights	local predictions	
(8.50	* 0.1)	= 0.85	= toes prediction
(0.65	* 0.2)	= 0.13	= wlrec prediction
(1.20	* 0.0)	= 0.00	= fans prediction

toes prediction + wlrec prediction + fans prediction = final prediction

0.85	+	0.13	+	0.00	=	0.98
------	---	------	---	------	---	------

The intuition behind how and why a dot product (weighted sum) works is easily one of the most important parts of truly understanding how neural networks make predictions. Loosely stated, a dot product gives us a *notion of similarity* between two vectors. Consider the examples:

```
a = [ 0, 1, 0, 1]
b = [ 1, 0, 1, 0]
c = [ 0, 1, 1, 0]
d = [.5, 0,.5, 0]
e = [ 0, 1,-1, 0]

w_sum(a,b) = 0
w_sum(b,c) = 1
w_sum(b,d) = 1
w_sum(c,c) = 2
w_sum(d,d) = .5
w_sum(c,e) = 0
```

The highest weighted sum (w_sum(c,c)) is between vectors that are exactly identical. In contrast, since a and b have no overlapping weight, their dot product is zero. Perhaps the most interesting weighted sum is between c and e, since e has a negative weight. This negative weight cancelled out the positive similarity between them. However, a dot product between e and itself would yield the number 2, despite the negative weight (double negative turns positive). Let's become familiar with these properties.

Some have equated the properties of the "dot product" to a "logical AND". Consider a and b.

```
a = [ 0, 1, 0, 1]
b = [ 1, 0, 1, 0]
```

If you asked whether both a[0] AND b[0] had value, the answer would be no. If you asked whether both a[1] AND b[1] had value, the answer would again be no. Since this is ALWAYS true for all 4 values, the final score equals 0. Each value failed the logical AND.

```
b = [ 1, 0, 1, 0]
c = [ 0, 1, 1, 0]
```

b and c, however, have one column that shares value. It passes the logical AND since b[2] AND c[2] have weight. This column (and only this column) causes the score to rise to 1.

```
c = [ 0, 1, 1, 0]
d = [.5, 0,.5, 0]
```

Fortunately, neural networks are also able to model partial ANDing. In this case, c and d share the same column as b and c, but since d only has 0.5 weight there, the final score is only 0.5. We exploit this property when modeling probabilities in neural networks.

```
d = [.5, 0,.5, 0]
e = [-1, 1, 0, 0]
```

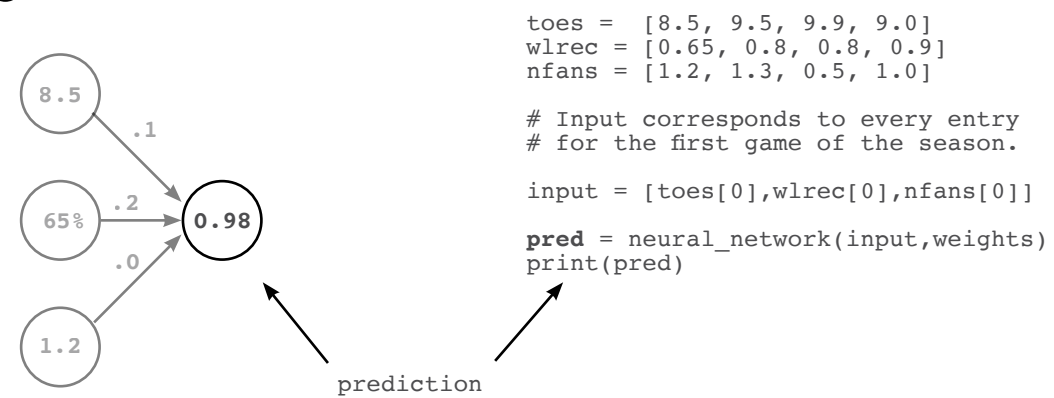
In this analogy, negative weights tend to imply a logical NOT operator, given that any positive weight paired with a negative weight will cause the score to go down. Furthermore, if both vectors have negative weights (such as w_sum(e,e)), then it will perform a *double negative* and add weight instead. Additionally, some will say that it's an OR after the AND, since if any of the rows show weight, the score is affected. Thus, for w_sum(a,b), if (a[0] AND b[0]) OR (a[1] AND b[1])...etc.. then w_sum(a,b) returns a positive score. Furthermore, if one value is negative, then that column gets a NOT. Amusingly, this actually gives us a kind of crude language to "read our weights". Let's "read" a few examples, shall we? These assume you're performing w_sum(input,weights) and the "then" to these "if statements" is just an abstract "then give high score".

```
weights = [ 1, 0, 1] => if input[0] OR input[2]
weights = [ 0, 0, 1] => if input[2]
weights = [ 1, 0, -1] => if input[0] OR NOT input[2]
weights = [ -1, 0, -1] => if NOT input[0] OR NOT input[2]
weights = [ 0.5, 0, 1] => if BIG input[0] or input[2]
```

Notice in the last row that a weight[0] = 0.5 means that the corresponding input [0] would have to be larger to compensate for the smaller weighting. And as I mentioned, this is a very very crude approximate language. However, I find it to be immensely useful when trying to picture in my head what's going on under the hood. This will help us significantly in the future, especially when putting networks together in increasingly complex ways.

So, given these intuitions, what does this mean when our neural network makes a prediction? Very roughly speaking, it means that our network gives a high score of our inputs based on *how similar they are to our weights*. Notice below that "nfans" is completely ignored in the prediction because the weight associated with it is a 0. The most sensitive predictor, in fact, is "wlrec" because its weight is a 0.2. However, the dominant force in the high score is the number of toes ("ntoes") not because the weight is the highest, but because the input combined with the weight is by far the highest.

4 Deposit Prediction



Here are a few more points that we will note for further reference. We cannot shuffle our weights. They have specific positions they need to be in. Furthermore, both the value of the weight AND the value of the input determine the overall impact on the final score. Finally, a negative weight would cause some inputs to reduce the final prediction (and vise versa).

Multiple Inputs - Complete Runnable Code

The code snippets from this example come together as follows.

Previous Code

We can create and execute our neural network using the following code. For the purposes of clarity, I have written everything out using only basic properties of Python (lists and numbers). However, there is a better way that we will start using in the future. There is a python library called "NumPy" which stands for "numerical python". It has very efficient code for creating vectors and performing common functions (such as a dot product). So, without further ado, here's the same code in numpy.

```
def w_sum(a,b):
    assert(len(a) == len(b))
    output = 0
    for i in range(len(a)):
        output += (a[i] * b[i])
    return output

weights = [0.1, 0.2, 0]

def neural_network(input, weights):
    pred = w_sum(input,weights)
    return pred
```

NumPy Code

```
import numpy as np

weights = np.array([0.1, 0.2, 0])

def neural_network(input, weights):
    # Input corresponds to every entry
    # for the first game of the season.
    pred = input.dot(weights)
    return pred

toes = np.array([8.5, 9.5, 9.9, 9.0])
wlrec = np.array([0.65, 0.8, 0.8, 0.9])
nfans = np.array([1.2, 1.3, 0.5, 1.0])

# Input corresponds to every entry
# for the first game of the season.

input = np.array([toes[0],wlrec[0],nfans[0]])
pred = neural_network(input,weights)
print(pred)
```

Both networks should simply print out:

0.98

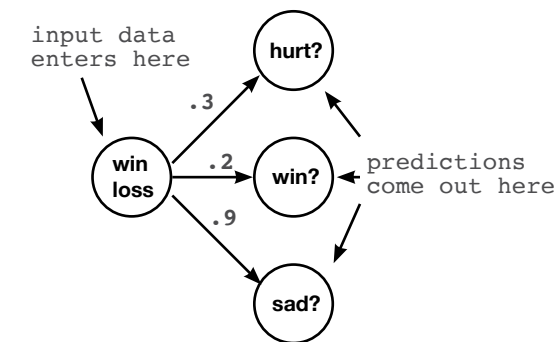
Notice that we didn't have to create a special "w_sum" function. Instead, NumPy has a special function called "dot" (short for "dot product") which we can call. Many of the functions we want to use in the future will have NumPy parallels, as we will see later.

Making a Prediction with Multiple Outputs

Neural Networks can also make multiple predictions using only a single input.

Perhaps a simpler augmentation than multiple inputs is multiple outputs. Prediction occurs in the same way as if there were 3 disconnected single-weight neural networks.

1 An Empty Network With Multiple Outputs



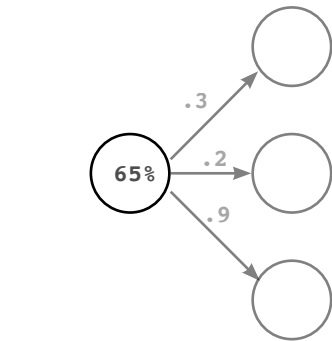
```
# Instead of predicting just
# whether the team won or lost,
# now we're also predicting
# whether they are happy/sad AND
# the percentage of the team that
# is hurt. We are making this
# prediction using only the
# current win/loss record.

weights = [0.3, 0.2, 0.9]

def neural_network(input, weights):
    pred = ele_mul(input, weights)
    return pred
```

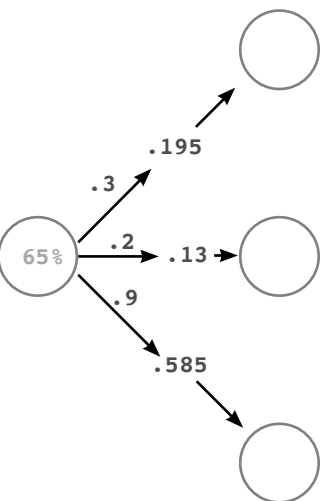
The most important commentary in this setting is to notice that the 3 predictions really are completely separate. Unlike neural networks with multiple inputs and a single output where the prediction is undeniably connected, this network truly behaves as 3 independent components, each receiving the same input data. This makes the network quite simple to implement.

2 Inserting One Input Datapoint



```
wlrec = [0.65, 0.8, 0.8, 0.9]
input = wlrec[0]
pred = neural_network(input, weights)
```

3 Perform an Elementwise Multiplication

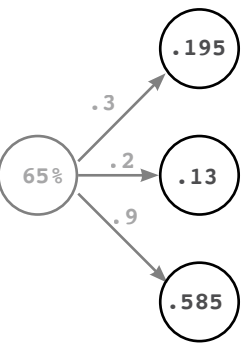


```
def ele_mul(number, vector):
    output = [0,0,0]
    assert(len(output) == len(vector))
    for i in range(len(vector)):
        output[i] = number * vector[i]
    return output

def neural_network(input, weights):
    pred = ele_mul(input, weights)
    return pred
```

inputs	weights	final predictions
(0.65 * 0.3)		= 0.195 = hurt prediction
(0.65 * 0.2)		= 0.13 = win prediction
(0.65 * 0.9)		= 0.585 = sad prediction

4 Deposit Predictions



predictions
(a vector of numbers)

```
wlrec = [0.65, 0.8, 0.8, 0.9]
input = wlrec[0]
pred = neural_network(input, weight)
print(pred)
```

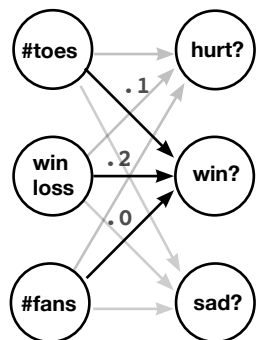
Predicting with Multiple Inputs & Outputs

Neural networks can predict multiple outputs given multiple inputs.

Finally, the way in which we built a network with multiple inputs or outputs can be combined together to build a network that has both multiple inputs AND multiple outputs. Just like before, we simply have a weight connecting each input node to each output node and prediction occurs in the usual way.

1 An Empty Network With Multiple Inputs & Outputs

inputs predictions

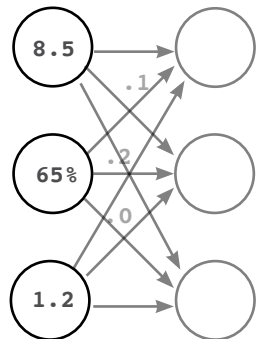


```
#toes %win #fans
weights = [ [0.1, 0.1, -0.3], #hurt?
            [0.1, 0.2, 0.0], #win?
            [0.0, 1.3, 0.1] ] #sad?

def neural_network(input, weights):
    pred = vect_mat_mul(input, weights)
    return pred
```

2 Inserting One Input Datapoint

inputs predictions



```
# This dataset is the current
# status at the beginning of
# each game for the first 4 games
# in a season.

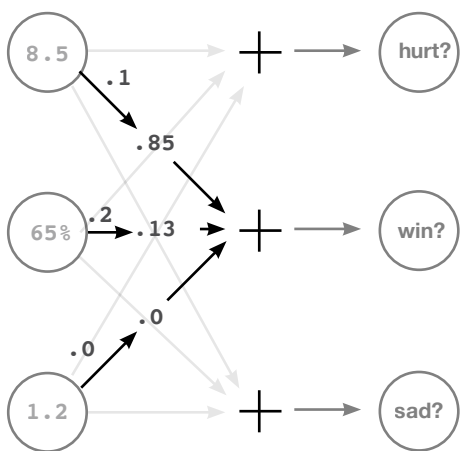
# toes = current number of toes
# wlrec = current games won (percent)
# nfans = fan count (in millions)

toes = [8.5, 9.5, 9.9, 9.0]
wlrec = [0.65, 0.8, 0.8, 0.9]
nfans = [1.2, 1.3, 0.5, 1.0]

# Input corresponds to every entry
# for the first game of the season.

input = [toes[0], wlrec[0], nfans[0]]
pred = neural_network(input, weights)
```

3 For Each Output, Perform a Weighted Sum of Inputs



```
def w_sum(a,b):
    assert(len(a) == len(b))
    output = 0
    for i in range(len(a)):
        output += (a[i] * b[i])
    return output

def vect_mat_mul(vect, matrix):
    assert(len(vect) == len(matrix))
    output = [0,0,0]

    for i in range(len(vect)):
        output[i] = w_sum(vect, matrix[i])

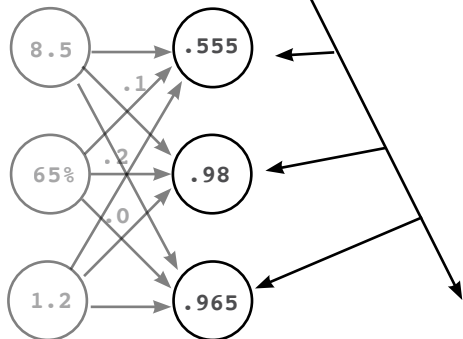
    return output

def neural_network(input, weights):
    pred = vect_mat_mul(input, weights)
    return pred
```

#toes	%win	#fans	
8.5	0.65	1.2	$(8.5 * 0.1) + (0.65 * 0.1) + (1.2 * -0.3) = 0.555 = \text{hurt prediction}$
8.5	0.65	1.2	$(8.5 * 0.1) + (0.65 * 0.2) + (1.2 * 0.0) = 0.98 = \text{win prediction}$
8.5	0.65	1.2	$(8.5 * 0.0) + (0.65 * 1.3) + (1.2 * 0.1) = 0.965 = \text{sad prediction}$

4 Deposit Predictions

inputs predictions



```
toes = [8.5, 9.5, 9.9, 9.0]
wlrec = [0.65, 0.8, 0.8, 0.9]
nfans = [1.2, 1.3, 0.5, 1.0]

# Input corresponds to every entry
# for the first game of the season.

input = [toes[0], wlrec[0], nfans[0]]
pred = neural_network(input, weight)
```

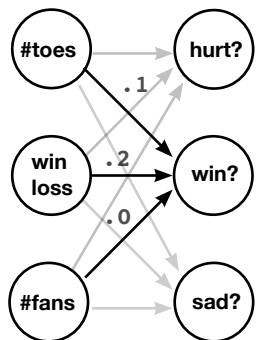

Multiple Inputs & Outputs - How does it work?

It performs 3 independent weighted sums of the input to make 3 predictions.

I find that there are 2 perspectives one can take on this architecture. You can either think of it as 3 weights coming out of each input node, or 3 weights going into each output node. For now, I find the latter to be much more beneficial. Think about this neural network as 3 independent dot products, 3 independent weighted sums of the input. Each output node takes its own weighted sum of the input and makes a prediction.

1 An Empty Network With Multiple Inputs & Outputs

inputs predictions

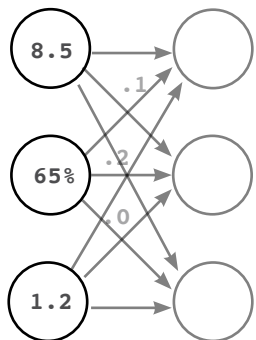


```
#toes %win #fans
weights = [ [0.1, 0.1, -0.3],#hurt?
            [0.1, 0.2, 0.0], #win?
            [0.0, 1.3, 0.1] ]#sad?

def neural_network(input, weights):
    pred=vect_mat_mul(input,weights)
    return pred
```

2 Inserting One Input Datapoint

inputs predictions



```
# This dataset is the current
# status at the beginning of
# each game for the first 4 games
# in a season.

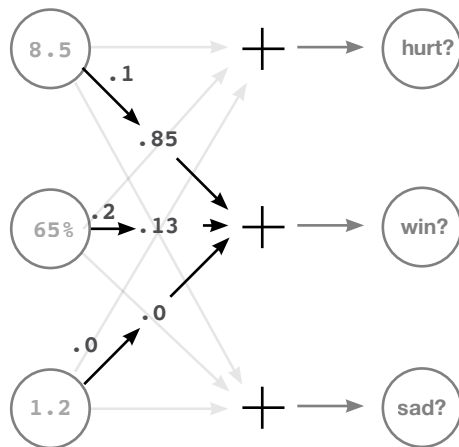
# toes = current number of toes
# wlrec = current games won (percent)
# nfans = fan count (in millions)

toes = [8.5, 9.5, 9.9, 9.0]
wlrec = [0.65,0.8, 0.8, 0.9]
nfans = [1.2, 1.3, 0.5, 1.0]

# Input corresponds to every entry
# for the first game of the season.

input = [toes[0],wlrec[0],nfans[0]]
pred = neural_network(input,weights)
```

3 For Each Output, Perform a Weighted Sum of Inputs



```
def w_sum(a,b):
    assert(len(a) == len(b))
    output = 0
    for i in range(len(a)):
        output += (a[i] * b[i])
    return output

def vect_mat_mul(vect,matrix):
    assert(len(vect) == len(matrix))
    output = [0,0,0]

    for i in range(len(vect)):
        output[i]=w_sum(vect,matrix[i])

    return output

def neural_network(input, weights):
    pred=vect_mat_mul(input,weights)
    return pred
```

#toes	%win	#fans	
8.5	0.65	1.2	$(8.5 * 0.1) + (0.65 * 0.1) + (1.2 * -0.3) = 0.555 = \text{hurt prediction}$
8.5	0.65	1.2	$(8.5 * 0.1) + (0.65 * 0.2) + (1.2 * 0.0) = 0.98 = \text{win prediction}$
8.5	0.65	1.2	$(8.5 * 0.0) + (0.65 * 1.3) + (1.2 * 0.1) = 0.965 = \text{sad prediction}$

As mentioned on the previous page, we are choosing to think about this network as a series of weighted sums. Thus, in the code above, we created a new function called "vect_mat_mul". This function iterates through each row of our weights (each row is a vector), and makes a prediction using our w_sum function. It is literally performing 3 consecutive weighted sums and then storing their predictions in a vector called "output". There's a lot more weights flying around in this one, but it isn't that much more advanced than networks we have previously seen.

I want to use this "list of vectors" and "series of weighted sums" logic to introduce you to two new concepts. See the weights variable in step (1)? It's a list of vectors. A list of vectors is simply called a **matrix**. It is as simple as it sounds. Furthermore, there are functions that we will find ourselves commonly using that leverage matrices. One of these is called **vector-matrix multiplication**. Our "series of weighted sums" is exactly that. We take a vector, and perform a dot product with every row in a matrix**. As we will find out on the next page, we even have special NumPy functions to help us out.

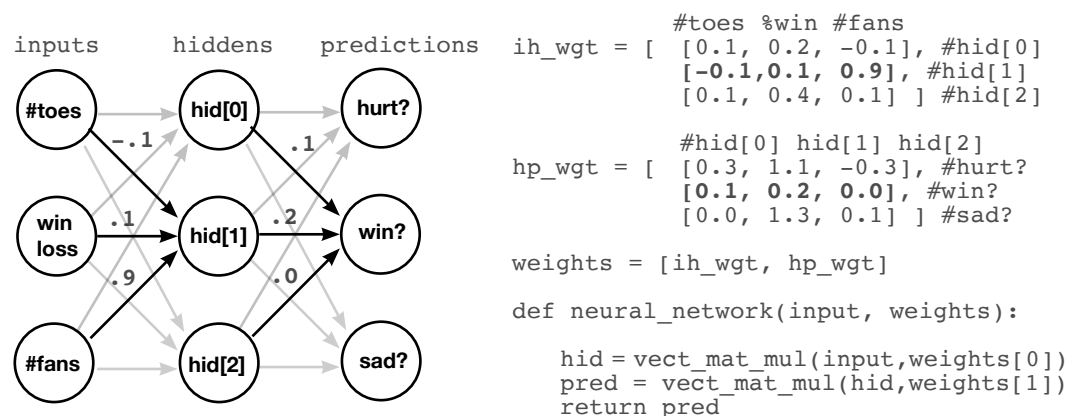
** Note: For those of you experienced with Linear Algebra, the more formal definition would store/process weights as column vectors instead of row vectors. This will be rectified shortly.

Predicting on Predictions

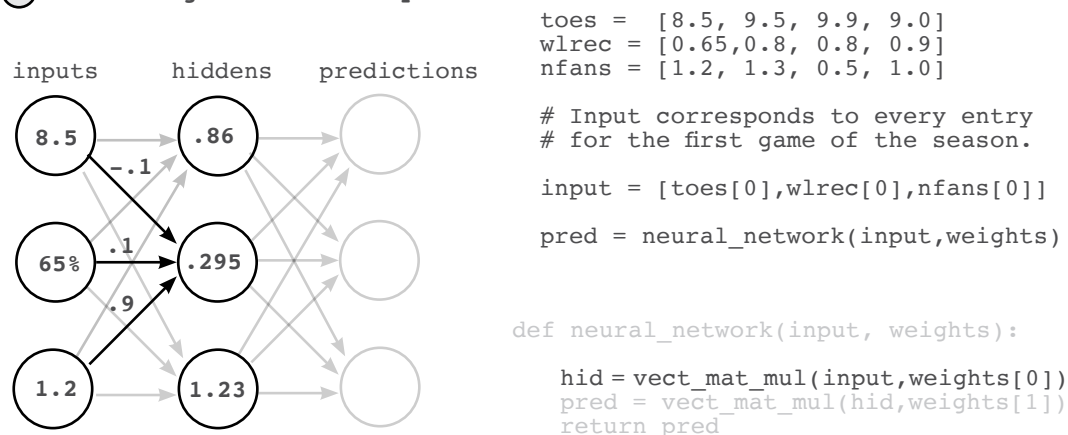
Neural networks can be stacked!

As the pictures below make clear, one can also take the output of one network and feed it as input to another network. This results in two consecutive vector-matrix multiplications. It may not yet be clear why you would predict in this way. However, some datasets (such as image classification) contain patterns that are simply too complex for a single weight matrix. Later, we will discuss the nature of these patterns. For now, it is sufficient that you know this is possible.

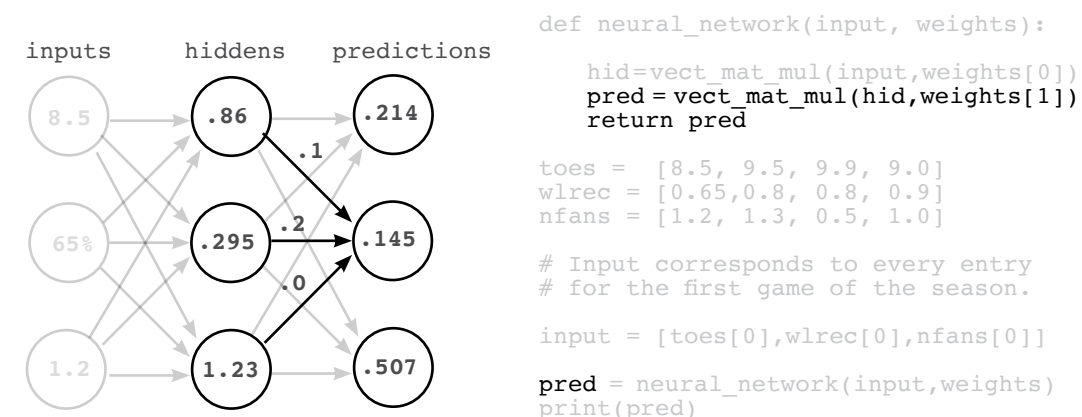
① An Empty Network With Multiple Inputs & Outputs



② Predicting the Hidden Layer



③ Predicting the Output Layer (and depositing the prediction)



NumPy Version

```

import numpy as np

#toes %win %fans
ih_wgt = np.array([
    [0.1, 0.2, -0.1], #hid[0]
    [-0.1, 0.1, 0.9], #hid[1]
    [0.1, 0.4, 0.1]])T #hid[2]

# hid[0] hid[1] hid[2]
hp_wgt = np.array([
    [0.3, 1.1, -0.3], #hurt?
    [0.1, 0.2, 0.0], #win?
    [0.0, 1.3, 0.1] ]).T #sad?

weights = [ih_wgt, hp_wgt]

def neural_network(input, weights):
    hid = input.dot(weights[0])
    pred = hid.dot(weights[1])
    return pred

toes = np.array([8.5, 9.5, 9.9, 9.0])
wlrec = np.array([0.65, 0.8, 0.8, 0.9])
nfans = np.array([1.2, 1.3, 0.5, 1.0])

input = np.array([toes[0], wlrec[0], nfans[0]])

pred = neural_network(input, weights)
print(pred)
  
```

A Quick Primer on NumPy

NumPy is so easy to use that it does a few things for you. Let's reveal the magic.

So far in this chapter, we've discussed two new types of mathematical tools: vectors and matrices. Furthermore, we have learned about different operations that occur on vectors and matrices including dot products, elementwise multiplication and addition, as well as vector-matrix multiplication. For these operations, we've written our own Python functions that can operate on simple Python "list" objects. In the short term, we will keep writing/using these functions so that we make sure we fully understand what's going on inside them. However, now that we've mentioned both "numpy" and several of the big operations, I'd like to give you a quick run-down of basic "numpy" use so that you will be ready for our transition to "only numpy" a few chapters from now. So, let's just start with the basics again, vectors and matrices.

```
import numpy as np

a = np.array([0,1,2,3]) # a vector
b = np.array([4,5,6,7]) # another vector
c = np.array([[0,1,2,3], # a matrix
              [4,5,6,7]])

d = np.zeros((2,4)) # (2x4 matrix of zeros)
e = np.random.rand(2,5) # random 2x5
# matrix with all numbers between 0 and 1

print(a)
print(b)
print(c)
print(d)
print(e)
```

Output

```
[0 1 2 3]
[4 5 6 7]
[[0 1 2 3]
 [4 5 6 7]]
[[ 0.  0.  0.  0.]
 [ 0.  0.  0.  0.]]
[[ 0.22717119  0.39712632
  0.0627734    0.08431724]
 [ 0.09675954  0.99012254
  0.45922775  0.3273326
  0.28617742]]
```

We can create vectors and matrices in multiple ways in NumPy. Most of the common ones for neural networks are listed above. Note that the processes for creating a vector and a matrix are identical. If you create a matrix with only one row, you're creating a vector. Furthermore, as in mathematics in general, you create a matrix by listing (rows,columns). I say that only so that you can remember the order. Rows comes first. Columns comes second. Let's see some operations we can do on these vectors and matrices.

```
print(a * 0.1) # multiplies every number in vector "a" by 0.1
print(c * 0.2) # multiplies every number in matrix "c" by 0.2
print(a * b) # multiplies elementwise between a and b (columns paired up)
print(a * b * 0.2) # elementwise multiplication then multiplied by 0.2
print(a * c) # since c has the same number of columns as a, this performs
# elementwise multiplication on every row of the matrix "c"

print(a * e) # since a and e don't have the same number of columns, this
# throws a "Value Error: operands could not be broadcast together with.."
```

Go ahead and run all of the code on the previous page. The first bit of "at first confusing but eventually heavenly" magic should be visible on that page. When you multiply two variables with the "*" function, NumPy automatically detects what kinds of variables you're working with and "tries" to figure out the operation you're talking about. This can be mega-convenient but sometimes makes NumPy a bit hard to read. You have to make sure you keep up with what each variable type is in your head as you go along.

The general rule of thumb for anything elementwise (+,-,*,/) is that the two variables must either have the SAME number of columns, or one of the variables must only have 1 column.

For example, "print(a * 0.1)" takes a vector and multiplies it by a single number (a scalar). NumPy goes, "Oh, I bet I'm supposed to do vector-scalar multiplication here," and then it takes the scalar (0.1) and multiplies it by every value in the vector. This looks exactly the same as "print(c * 0.2)", except that NumPy knows that "c" is a matrix. Thus, it performs scalar-matrix multiplication, multiplying every element in "c" by 0.2. Because the scalar has only one column, you can multiply it by anything (or divide, add, or subtract for that matter).

Next up: "print(a * b)". NumPy first identifies that they're both vectors. Since neither vector has only 1 column, it checks to see if they have an identical number of columns. Since they do, it knows to simply multiply each element by each element, based on their positions in the vectors. The same is true with addition, subtraction, and division.

"print(a * c)" is perhaps the most elusive. "a" is a vector with 4 columns. "c" is a (2x4) matrix. Neither have only one column, so numpy checks to see if they have the same number of columns. Since they do, numpy multiplies the vector "a" by each row of "c" (as if it was doing elementwise vector multiplication on each row).

Again, the most confusing part about this is that all of these operations look the same if you don't know which variables are scalars, vectors, or matrices. When I'm "reading NumPy", I'm really doing 2 things: reading the operations and keeping track of the "shape" (number of rows and columns) of each operation. It will take some practice, but eventually it becomes second nature.

```
a = np.zeros((1,4)) # vector of length 4
b = np.zeros((4,3)) # matrix with 4 rows & 3 columns

c = a.dot(b)
print(c.shape)
```

Output

```
(1,3)
```

There is one golden rule when using the 'dot' function. If you put the (rows,cols) description of the two variables you're "dotting" next to each other, neighboring numbers should always be the same. In this case, we're dot producting a (1,4) with a (4,3). Thus, it works fine, and outputs a (1,3). In terms of variable shape, you can think of it this way, regardless of whether you're "dotting" vectors or matrices. Their "shape" (number of rows and columns) must line up. The columns on the "left" matrix must equal rows on the "right".

$$(a,b).dot(b,c) = (a,c)$$

```

a = np.zeros((2,4)) # matrix with 2 rows and 4 columns
b = np.zeros((4,3)) # matrix with 4 rows & 3 columns

c = a.dot(b)
print(c.shape) # outputs (2,3)

e = np.zeros((2,1)) # matrix with 2 rows and 1 columns
f = np.zeros((1,3)) # matrix with 1 row & 3 columns

g = e.dot(f)
print(g.shape) # outputs (2,3)


h = np.zeros((5,4)).T # matrix with 4 rows and 5 columns
i = np.zeros((5,6)) # matrix with 6 rows & 5 columns

j = h.dot(i)
print(j.shape) # outputs (4,6)

h = np.zeros((5,4)) # matrix with 5 rows and 4 columns
i = np.zeros((5,6)) # matrix with 5 rows & 6 columns
j = h.dot(i)
print(j.shape) # throws an error

```

this ".T" "flips" the rows and columns of a matrix



Conclusion

To predict, neural networks perform repeated weighted sums of the input.

We have seen an increasingly complex variety of neural networks in this chapter. I hope that it is clear that a relatively small number of simple rules are used repeatedly to create larger, more advanced neural networks. Furthermore, the intelligence of the network really depends on what weight values we give to our network.

Everything we have done in this chapter is a form of what's called **Forward Propagation**, wherein a neural network takes input data and makes a prediction. It is called this because we are *propagating* activations *forward* through the network. In these examples, **activations** are all of the numbers that are *not* weights, and are unique for every prediction.

In the next chapter, we will be learning how to set our weights so that our neural networks make accurate predictions. We will find that in the same way that prediction is actually based on several simple techniques that are repeated/stacked on top of each other, "weight learning" is also a series of simple techniques that are combined many times across an architecture. See you there!

Introduction to Neural Learning Gradient Descent

4

IN THIS CHAPTER

Do neural networks make accurate predictions?

Why measure error?

Hot and Cold Learning

Calculating both direction and amount from error

Gradient Descent

Learning is Just Reducing Error

Derivatives and how to use them to learn

Divergence and Alpha

“ The only relevant test of the validity of a hypothesis is comparison of prediction with experience. ”

— MILTON FRIEDMAN

Predict, Compare, and Learn

This chapter is about "Compare", and "Learn"

In Chapter 3, we learned about the paradigm "Predict, Compare, Learn" and we dove deep into the first step: "Predict". In this process we learned a myriad of things, including the major parts of neural networks (nodes and weights), how datasets fit into networks (matching the number of datapoints coming in at one time), and finally how to use a neural network to make a prediction. Perhaps this process begged the question, "How do we set our weight values so that our network predicts accurately?". Answering this question will be the main focus of this chapter, covering the second two steps of our paradigm, "Compare", and "Learn".

Compare

A measurement of how much our prediction "missed".

Once we've made a prediction, the next step to learn is to evaluate how well we did. Perhaps this might seem like a rather simple concept, but we will eventually find that coming up with a good way to measure error is one of the most important and complicated subjects of Deep Learning.

In fact, there are many properties of "measuring error" that you have likely been doing your whole life without realizing it. Perhaps you (or someone you know) amplifies bigger errors while ignoring very small ones. In this chapter we will learn how to mathematically teach our network to do this. Furthermore, we will learn that error is always positive! We will consider the analogy of an "archer" hitting a target. Whether he is too low by an inch or too high by an inch, the error is still just 1 inch! In our neural network "Compare" step, we want to consider these kinds of properties when measuring error.

As a heads up, in this chapter we will only evaluate one, very simple way of measuring error called "Mean Squared Error". However, it is but one of many ways to evaluate the accuracy of your neural network.

This step will give us a sense for "how much we missed", but this isn't enough to be able to learn. The output of our "compare" logic will simply be a "hot or cold" type signal. Given some prediction, we'll calculate an error measure that will either say "a lot" or "a little". It won't tell us why we missed, what direction we missed, or what we should do to fix it. It more or less just says "big miss", "little miss", or "perfect prediction". What we do about our error is captured in the next step, "Learn".

Learn

"Learning" takes our error and tells each weight how it can change to reduce it.

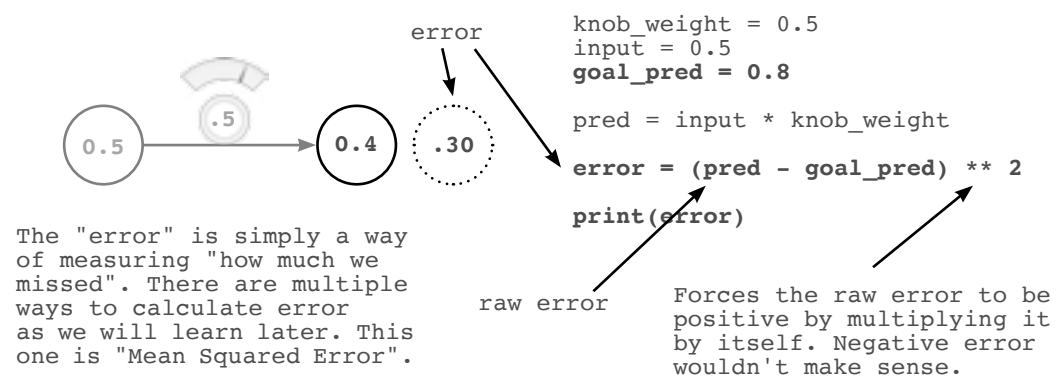
Learning is all about "error attribution", or the art of figuring out how each weight played its part in creating error. It's the "blame game" of Deep Learning. In this chapter, we will spend a great number of pages learning the most popular version of the Deep Learning "blame game" called **Gradient Descent**.

At the end of the day, it's going to result in computing a number for each of our weights. That number will represent how that weight should be higher or lower in order to reduce the error. Then we will move the weight according to that number, and we'll be done.

Compare: Does our network make good predictions?

Let's measure the error and find out!

Execute this code in your Jupyter notebook. It should print "0.3025".



What is the `goal_pred` variable?

Much like *input*, it's a number we recorded in the real world somewhere, but it's usually something that's hard to observe, like "the percentage of people who DID wear sweatshirts" given the temperature or "whether the batter DID in fact hit a home run" given his batting average.

Why is the error *squared*?

Think about an archer hitting a target. When he hits 2 inches too high, how much did he miss by? When he hits 2 inches too low, how much did he miss by? Both times he only missed by 2 inches. The primary reason why we *square* "how much we missed" is that it forces the output to be *positive*. $(pred - goal_pred)$ could be negative in some situations... *unlike actual error*.

Doesn't squaring make big errors (>1) bigger and small errors (<1) smaller?

Yeah... It is kind of a weird way of measuring error, but it turns out that **amplifying** big errors and **reducing** small errors is actually ok. Later, we'll use this error to help the network learn, and we'd rather it *pay attention* to the big errors and not worry so much about the small ones. Good parents are like this too. They practically ignore errors if they're small enough (i.e. breaking the lead on your pencil) but might go nuclear for big errors (i.e. crashing the car). See why squaring is valuable?

Why measure error?

Measuring error simplifies the problem.

The goal of training our neural network is to make correct predictions. That's what we want. And in the most pragmatic world (as mentioned in the last chapter), we want the network to take input that we can easily calculate (today's stock price), and predict things that are hard to calculate (tomorrow's stock price). That's what makes a neural network useful.

It turns out that "changing knob_weight to make the network correctly predict the goal_prediction" is *slightly* more complicated than "changing the knob_weight to make $error == 0$ ". There's something more concise about looking at the problem this way. Ultimately, both of those statements say the same thing, but trying to *get the error to 0* just seems a bit more straightforward.

Different ways of measuring error *prioritize error differently*.

If this is a bit of a stretch right now, that's ok, but think back to what I said on the last page. By *squaring* the error, numbers that are less than 1 get *smaller* whereas numbers that are greater than 1 get *bigger*. This means that we're going to change what I call "**pure error**" $(pred - goal_pred)$ so that bigger errors become VERY big and smaller errors quickly become irrelevant. By measuring error this way, we can *prioritize* big errors over smaller ones. When we have somewhat large "pure errors" (say... 10), we're going to tell ourselves we have *very* large error $(10**2 == 100)$, and in contrast, when we have small "pure errors" (say... 0.01), we're going to tell ourselves that we have *very* small error $(0.01**2 == 0.0001)$. See what I mean about *prioritizing*? It's just modifying what we *consider to be error* so that we amplify big ones and largely ignore small ones. In contrast, if we took the *absolute value* instead of *squaring* the error, we wouldn't have this type of prioritization. The error would just be the positive version of the "pure error" ... which would be fine, just different. More on this later.

Why do we only want *positive* error?

Eventually, we're going to be working with *millions* of input -> goal_prediction pairs, and we're still going to want to make accurate predictions. This means that we're going to try to take the *average error* down to 0.

This presents a problem if our error can be positive and negative. Imagine if we had two datapoints, two input -> goal_prediction pairs, that we were trying to get the neural network to correctly predict. If the first had an error of 1,000, and the second had an error of -1,000, then our *average error* would be ZERO! We would fool ourselves into thinking we predicted perfectly when we missed by 1000 each time!!! This would be really bad. Thus, we want the error of *each prediction* to always be *positive* so that they don't accidentally cancel each other out when we average them.

What's the Simplest Form of Neural Learning?

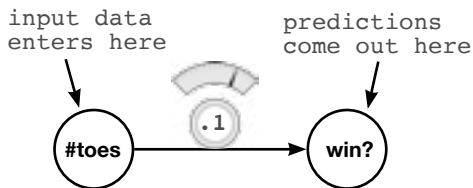
Learning using the Hot and Cold Method

At the end of the day, learning is really about one thing, adjusting our knob_weight either up or down so that our **error** reduces. If we keep doing this and our **error** goes to 0, we are done learning! So, how do we know whether to turn the knob up or down? Well, we try *both up and down* and see which one reduces the error! Whichever one reduces the error is used to actually update the **knob_weight**. It's simple, but effective. After we do this over and over again, eventually our error==0, which means our neural network is predicting with perfect accuracy.

Hot and Cold Learning

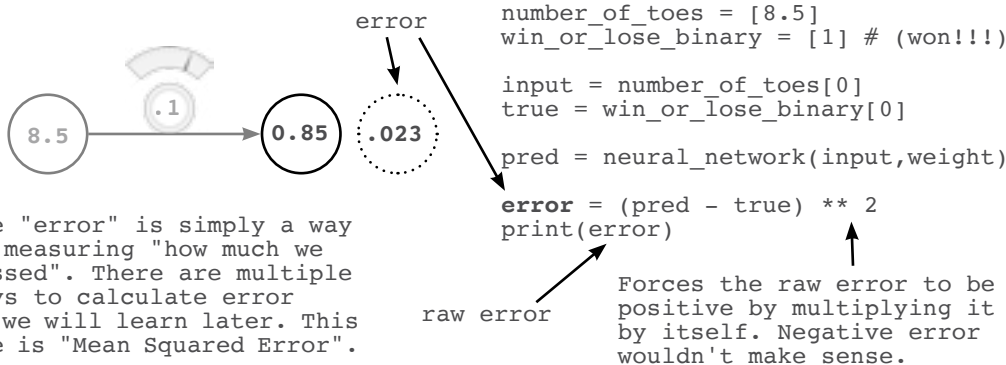
Wiggling our weights to see which direction reduces the error the most, moving our weights in that direction, and repeating until the error gets to 0.

1 An Empty Network



```
weight = 0.1
lr = 0.01
def neural_network(input, weight):
    prediction = input * weight
    return prediction
```

2 PREDICT: Making A Prediction And Evaluating Error



```
number_of_toes = [8.5]
win_or_lose_binary = [1] # (won!!!)

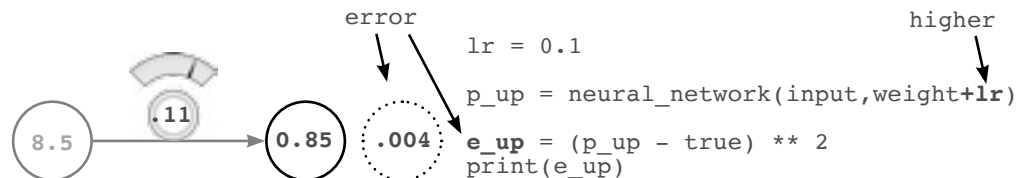
input = number_of_toes[0]
true = win_or_lose_binary[0]

pred = neural_network(input,weight)
error = (pred - true) ** 2
print(error)
```

↑ Forces the raw error to be positive by multiplying it by itself. Negative error wouldn't make sense.

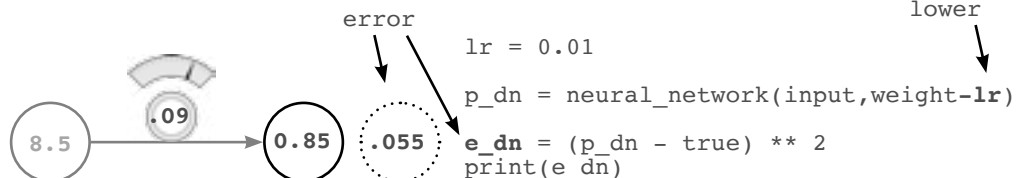
The "error" is simply a way of measuring "how much we missed". There are multiple ways to calculate error as we will learn later. This one is "Mean Squared Error".

3 COMPARE: Making A Prediction With a Higher Weight And Evaluating Error

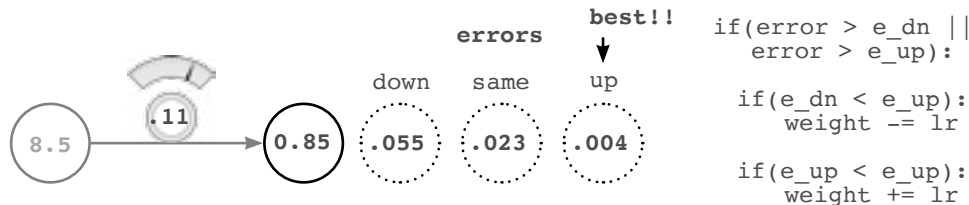


We want to move the weight so that the error goes downward, so we're going to try moving the weight up and down to see which one has the lowest error. First, we're trying moving the weight up (weight+lr).

4 COMPARE: Making A Prediction With a Lower Weight And Evaluating Error



5 COMPARE + LEARN: Comparing our Errors and Setting our New Weight



These last 5 steps comprise 1 iteration of Hot and Cold Learning. Fortunately, this iteration got us pretty close to the correct answer all by itself (the new error is only 0.004). However, under normal circumstances, we would have to repeat this process many times in order to find the correct weights. Some people even have to train their networks for weeks or months before they find a good enough weight configuration.

This reveals what learning in neural networks really is. It's a **search problem**. We are *searching* for the best possible configuration of weights so that our network's error falls to zero (and predicts perfectly). As with all other forms of search, we might not find exactly what we're looking for, and even if we do, it may take some time. On the next page, we'll use Hot and Cold Learning for a *slightly* more difficult prediction so that you can see this searching in action!

Hot and Cold Learning

Perhaps the simplest form of learning.

Execute this code in your Jupyter Notebook. (New neural network modifications are in **bold**.)

This code attempts to correctly predict 0.8.

```
weight = 0.5
input = 0.5
goal_prediction = 0.8
step_amount = 0.001

for iteration in range(1101):
    prediction = input * weight
    error = (prediction - goal_prediction) ** 2

    print("Error:" + str(error) + " Prediction:" + str(prediction))

    up_prediction = input * (weight + step_amount)
    up_error = (goal_prediction - up_prediction) ** 2

    down_prediction = input * (weight - step_amount)
    down_error = (goal_prediction - down_prediction) ** 2

    if(down_error < up_error):
        weight = weight - step_amount

    if(down_error > up_error):
        weight = weight + step_amount
```

how much to move our weights each iteration

repeat learning many times so that our error can keep getting smaller

TRY UP!

TRY DOWN!

If down is better, go down!

If up is better, go up!

When I run this code, I see the following output:

```
Error:0.3025 Prediction:0.25
Error:0.30195025 Prediction:0.2505
....
Error:2.50000000033e-07 Prediction:0.7995
Error:1.07995057925e-27 Prediction:0.8
```

Our last step correctly predicts 0.8!

Characteristics of Hot and Cold Learning

It's simple

Hot and Cold learning is simple. After making our prediction, we predict two more times, once with a slightly higher weight and again with a slightly lower weight. We then move the **weight** depending on which **direction** gave us a smaller **error**. Repeating this enough times eventually reduces our **error** down to 0.

Why did I iterate exactly 1101 times?

The neural network reaches 0.8 after exactly that many iterations. If you go past that, it wiggles back and forth between 0.8 and just above/below 0.8... making for a less pretty error log printed at the bottom of the left page. Feel free to try it out though.

PROBLEM #1: It's inefficient

We have to predict *multiple times* in order to make a single *knob_weight* update. This seems very inefficient.

PROBLEM #2: Sometimes it's impossible to predict the *exact* goal prediction.

With a set **step_amount**, unless the perfect **weight** is exactly $n \times \text{step_amount}$ away, the network will eventually overshoot by some number less than **step_amount**. When it does so, it will then start alternating back and forth between each side of the **goal_prediction**. Set the **step_amount** to 0.2 to see this in action. If you set **step_amount** to 10 you'll really break it! When I try this I see the following output. It never *remotely* comes *close* to 0.8!!!

The real problem here is that even though we know the correct **direction** to move our **weight**, we *don't know the correct amount*. Since we don't *know* the correct amount, we just pick a fixed one at random (**step_amount**). Furthermore, this *amount* has NOTHING to do with our **error**. Whether our **error** is BIG or our **error** is TINY, our **step_amount** is the same. So, Hot and Cold Learning is kind of a bummer. It's inefficient because we *predict 3 times for each weight update* and our **step_amount** is completely arbitrary, which can prevent us from learning the correct **weight** value.

```
Error:0.3025 Prediction:0.25
Error:19.8025 Prediction:5.25
Error:0.3025 Prediction:0.25
Error:19.8025 Prediction:5.25
Error:0.3025 Prediction:0.25
....
.... repeating infinitely...
```

What if we had a way of computing both **direction** and **amount** for each **weight** without having to repeatedly make predictions?

Calculating Both *direction* and *amount* from error

Let's measure the error and find out!

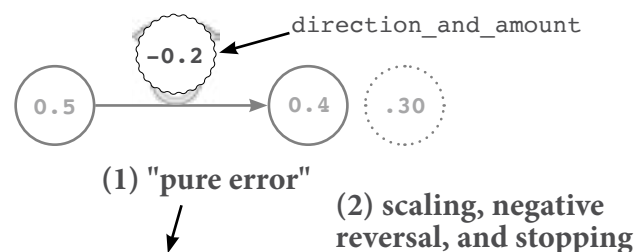
Execute this code in your Jupyter notebook.

```
weight = 0.5
goal_pred = 0.8
input = 0.5
```

```
for iteration in range(20):
    pred = input * weight
    error = (pred - goal_pred) ** 2
    direction_and_amount = (pred - goal_pred) * input
    weight = weight - direction_and_amount

    print("Error:" + str(error) + " Prediction:" + str(pred))
```

What you see above is a *superior* form of learning known as **Gradient Descent**. This method allows us to (in a single line of code, seen above in **bold**) calculate both the *direction* and the *amount* that we should change our *weight* so that we reduce our *error*.



What is the `direction_and_amount`?

It represents how we want to change our *weight*. The first part (1) is what we call "pure error" which equals `(pred - goal_pred)`. (More about this below.) The second part (2) is the multiplication by the *input* which performs scaling, negative reversal, and stopping, modifying the "pure error" so that it's ready to update our *weight*.

What is the "pure error"?

It's the `(pred - goal_pred)` which indicates "the raw direction and amount that we missed". If this is a *positive* number then we predicted too *high* and vice versa. If this is a *big* number then we missed by a *big* amount, etc.

What is "scaling, negative reversal, and stopping"?

These three attributes have the combined affect of translating our "pure error" into "the absolute amount that we want to change our *weight*". They do so by addressing three *major edge cases* at which points the "pure error" is not sufficient to make a good modification to our *weight*.

What is "stopping"?

This is the first (and simplest) effect on our "pure error" caused by multiplying it by our *input*. Imagine plugging a CD player into your stereo. If you turned the volume all the way up but the CD player was *off*... it simply wouldn't matter. "Stopping" addresses this in our neural network. If our *input* is 0, then it will force our *direction_and_amount* to also be 0. We don't learn (i.e. "change the volume") when our *input* is 0 because there's nothing to learn. Every *weight* value has the same *error*, and moving it makes no difference because the *pred* is always 0.

What is "negative reversal"?

This is probably our most difficult and important effect. Normally (when *input* is positive), moving our *weight* *upward* makes our prediction move *upward*. However, if our *input* is *negative*, then all of a sudden our *weight* changes directions!!! When our *input* is *negative*, then moving our *weight* *up* makes the prediction go *down*. It's reversed!!! How do we address this? Well, multiplying our "pure error" by our *input* will *reverse the sign* of our *direction_and_amount* in the event that our *input* is negative. This is "negative reversal", ensuring that our *weight* moves in the correct direction, even if the *input* is negative.

What is "scaling"?

Scaling is the third effect on our "pure error" caused by multiplying it by our *input*. Logically, it means that if our *input* was big, our *weight* update should also be big. This is more of a "side affect" as it can often go out of control. Later, we will use *alpha* to address when this scaling goes out of control.

When you run the code in the top left, you should see the following output.

```
Error:0.3025 Prediction:0.25
Error:0.17015625 Prediction:0.3875
Error:0.095712890625 Prediction:0.490625
...
```

```
Error:1.7092608064e-05 Prediction:0.79586567925
Error:9.61459203602e-06 Prediction:0.796899259437
Error:5.40820802026e-06 Prediction:0.797674444578
```

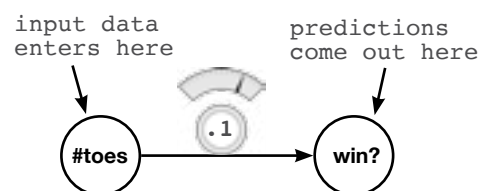
Our last steps correctly approach 0.8!

In this example, we saw **Gradient Descent** in action in a bit of an oversimplified environment. On the next page, we're going to see it in it's more native environment. Some terminology will be different, but we will code it in a way that makes it more obviously applicable to other kinds of networks (such as those with multiple inputs and outputs).

One Iteration of Gradient Descent

This performs a weight update on a single "training example" (input->>true) pair

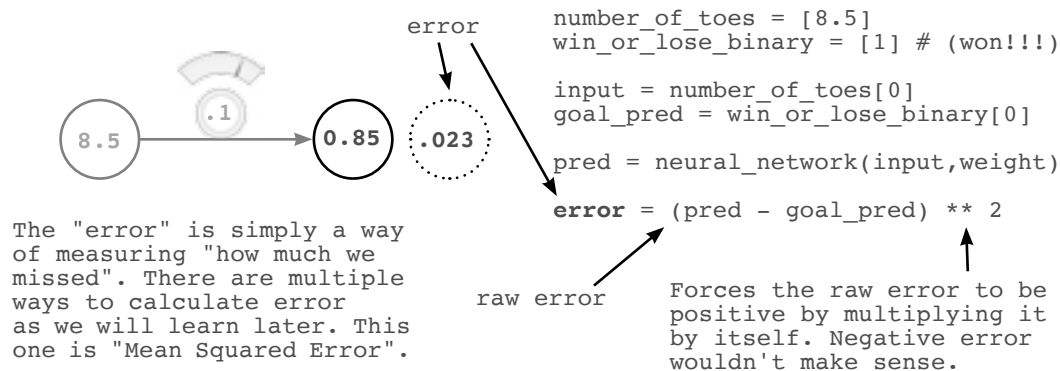
1 An Empty Network



```
weight = 0.1
alpha = 0.01

def neural_network(input, weight):
    prediction = input * weight
    return prediction
```

2 PREDICT: Making A Prediction And Evaluating Error



```
number_of_toes = [8.5]
win_or_lose_binary = [1] # (won!!!)

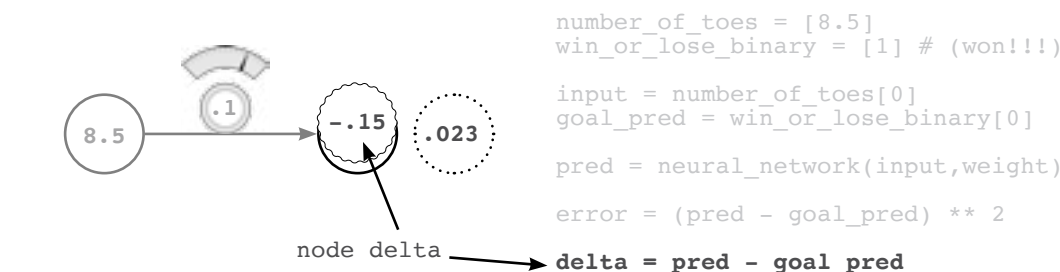
input = number_of_toes[0]
goal_pred = win_or_lose_binary[0]

pred = neural_network(input, weight)
error = (pred - goal_pred) ** 2
```

raw error

Forces the raw error to be positive by multiplying it by itself. Negative error wouldn't make sense.

3 COMPARE: Calculating "Node Delta" and Putting it on the Output Node



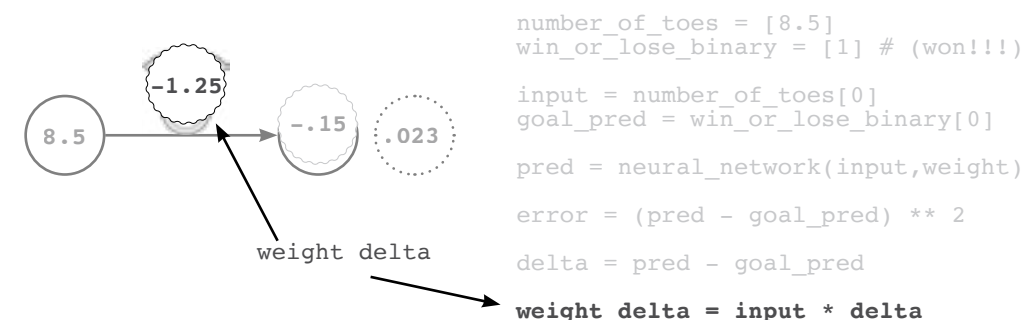
```
number_of_toes = [8.5]
win_or_lose_binary = [1] # (won!!!)

input = number_of_toes[0]
goal_pred = win_or_lose_binary[0]

pred = neural_network(input, weight)
error = (pred - goal_pred) ** 2
delta = pred - goal_pred
```

The primary difference between the gradient descent on the previous page and the implementation on this page just happened. `delta` is a new variable. It's the "raw amount that the node was too high or too low". Instead of computing `direction_and_amount` directly, we first calculate how much we wanted our output node to be different. Only then do we compute our `direction_and_amount` to change the weight (in step 4, now renamed "weight_delta").

4 LEARN: Calculating "Weight Delta" and Putting it on the Weight



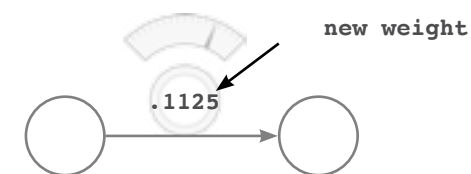
```
number_of_toes = [8.5]
win_or_lose_binary = [1] # (won!!!)

input = number_of_toes[0]
goal_pred = win_or_lose_binary[0]

pred = neural_network(input, weight)
error = (pred - goal_pred) ** 2
delta = pred - goal_pred

weight_delta = input * delta
```

5 LEARN: Updating the Weight



We multiply our `weight_delta` by a small number "alpha" before using it to update our weight. This allows us to control how fast the network learns. If it learns too fast, it can update weights too aggressively and overshoot. More on this later. Note that the weight update made the same change (small increase) as Hot and Cold Learning

```
number_of_toes = [8.5]
win_or_lose_binary = [1] # (won!!!)

input = number_of_toes[0]
goal_pred = win_or_lose_binary[0]

pred = neural_network(input, weight)

error = (pred - goal_pred) ** 2
delta = pred - goal_pred
weight_delta = input * delta

alpha = 0.01 # fixed before training
weight -= weight_delta * alpha
```


Learning Is Just Reducing Error

Modifying weight to reduce our error.

Putting together our code from the previous pages. We now have the following:

```
weight, goal_pred, input = (0.0, 0.8, 0.5)

for iteration in range(4):
    pred = input * weight
    error = (pred - goal_pred) ** 2
    delta = pred - goal_pred
    weight_delta = delta * input
    weight = weight - weight_delta
    print("Error:" + str(error) + " Prediction:" + str(pred))
```

these lines have a *secret*

The Golden Method for Learning

Adjusting each `weight` in the correct *direction* and by the correct *amount* so that our `error` reduces to 0.

All we're trying to do is figure out the right **direction** and **amount** to modify `weight` so that our `error` goes down. The secret to this lies in our `pred` and `error` calculations. Notice that we actually use our `pred` *inside* the `error` calculation. Let's replace our `pred` variable with the code we used to generate it.

```
error = ((input * weight) - goal_pred) ** 2
```

This doesn't change the value of error at all! It just combines our two lines of code so that we compute our error directly. Now, remember that our input and our goal_prediction are actually fixed at 0.5 and 0.8 respectively (we set them before the network even starts training). So, if we replace their variables names with the values... the *secret* becomes clear:

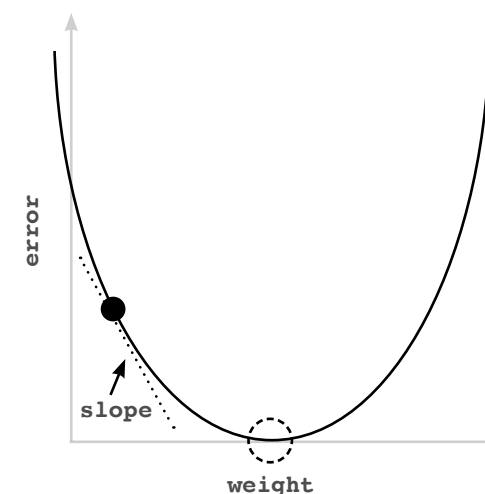
```
error = ((0.5 * weight) - 0.8) ** 2
```

The Secret

For any `input` and `goal_pred`, there is an *exact relationship* defined between our `error` and `weight`, found by combining our `prediction` and `error` formulas. In this case:

$$\text{error} = ((0.5 * \text{weight}) - 0.8) ** 2$$

Let's say that you moved `weight` up by 0.5... if there is an *exact relationship* between `error` and `weight`... we should be able to calculate how much this also *moves* the `error`! What if we wanted to *move* the `error` in a specific direction? Could it be done?



This graph represents *every value of error for every weight* according to the relationship in the formula above. Notice it makes a nice *bowl shape*. The black "dot" is at the point of BOTH our current `weight` and `error`. The dotted "circle" is where we want to be (`error == 0`).

Key Takeaway: The *slope* points to the bottom of the bowl (lowest `error`) *no matter where you are in the bowl*. We can use this *slope* to help our neural network *reduce the error*.

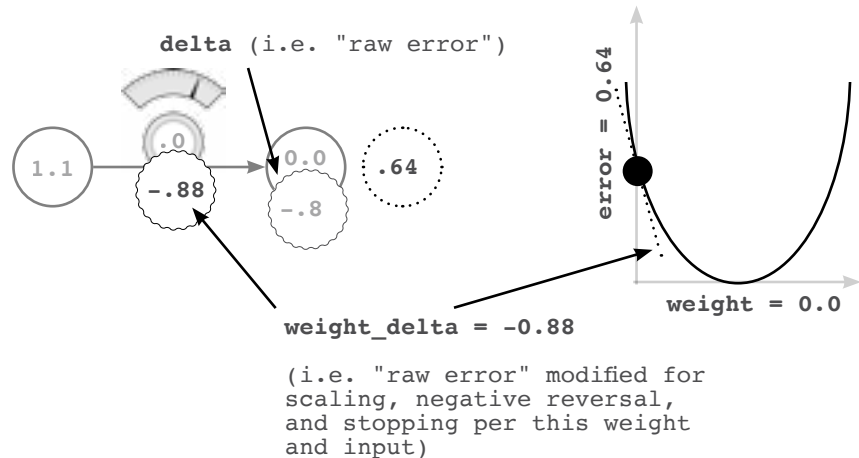
Let's Watch Several Steps of Learning

Will we eventually find the bottom of the bowl?

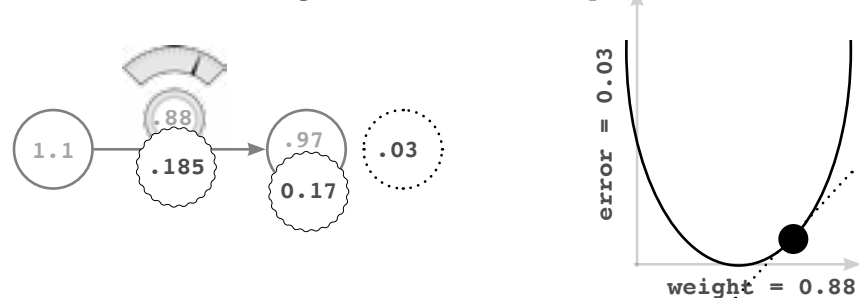
```
weight, goal_pred, input = (0.0, 0.8, 1.1)

for iteration in range(4):
    print("----\nWeight:" + str(weight))
    pred = input * weight
    error = (pred - goal_pred) ** 2
    delta = pred - goal_pred
    weight_delta = delta * input
    weight = weight - weight_delta
    print("Error:" + str(error) + " Prediction:" + str(pred))
    print("Delta:" + str(delta) + " Weight Delta:" + str(weight_delta))
```

1 A Big Weight Increase



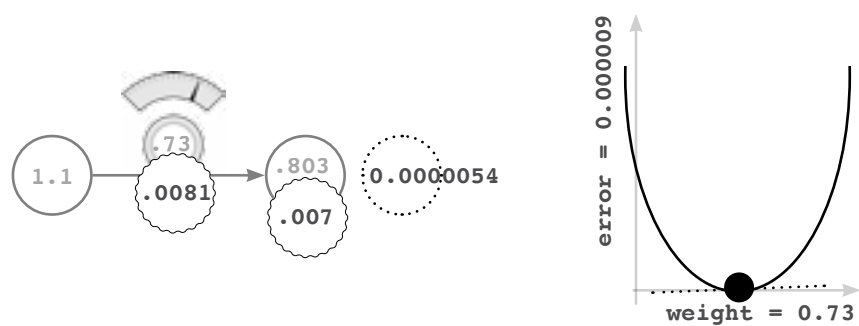
2 Overshot a bit... Let's go back the other way



3 Overshot Again! Let's go back again... but only just a little



4 Ok, we're pretty much there...



Code Output

```
-----
Weight:0.0
Error:0.64 Prediction:0.0
Delta:-0.8 Weight Delta:-0.88
-----
Weight:0.88
Error:0.028224 Prediction:0.968
Delta:0.168 Weight Delta:0.1848
-----
Weight:0.6952
Error:0.0012446784 Prediction:0.76472
Delta:-0.03528 Weight Delta:-0.038808
-----
Weight:0.734008
Error:5.489031744e-05 Prediction:0.8074088
Delta:0.0074088 Weight Delta:0.00814968
```

Why does this work? What really is `weight_delta`?

Let's back up and talk about functions. What is a function? How do we understand it?

Consider this function:

```
def my_function(x):
    return x * 2
```

A function takes some numbers as input and gives you another number as output. As you can imagine, this means that the function actually defines some sort of *relationship* between the input number(s) and the output number(s). Perhaps you can also see why the ability to *learn a function* is so powerful? It allows us to take some numbers (say, image pixels) and convert them into other numbers (say, the *probability* that the image contains a cat).

Now, every function has what you might call *moving parts*. It has pieces that we can tweak or change to make the output that the function generates *different*. Consider our "my_function" above. Ask yourself, "What is controlling the relationship between the input and the output of this function?". Well, it's the 2! Ask the same question about the function below:

```
error = ((input * weight) - goal_pred) ** 2
```

What is controlling the relationship between the `input` and the output (`error`)? Well, plenty of things are! This function is a bit more complicated! `goal_pred`, `input`, `**2`, `weight`, and all the parenthesis and algebraic operations (addition, subtraction, etc.) play a part in calculating the error. Tweaking any one of them would *change* the error. This is important to consider.

Just as a thought exercise, consider changing your `goal_pred` to reduce your error. Well, this is silly... but totally doable! In life, we might call this (setting your goals to be whatever your capability is) "giving up". It's just denying that we missed! This simply wouldn't do.

What if we changed the `input` until our error went to zero? Well, this is akin to seeing the world as you *want* to see it instead of as it actually is. This is changing your *input data* until you're predicting what you want to predict (sidenote: this is loosely how "inceptionism works").

Now consider changing the 2, or the additions, subtractions, or multiplications. This is just changing how you calculate error in the first place! Our error calculation is meaningless if it doesn't actually give us a good measure of *how much we missed* (with the right properties mentioned a few pages ago). This simply won't do either.

So, what do we have left? The only variable we have left is our `weight`. Adjusting this doesn't change our perception of the world, doesn't change our goal, and doesn't destroy our error measure. In fact, changing weight means that the function *conforms to the patterns in the data*. By forcing the rest of our function to be *unchanging*, we force our function to correctly model some pattern in our data. It is only allowed to modify how the network *predicts*.

So, at the end of the day, we're modifying specific parts of an error function until the `error` value goes to zero. This error function is calculated using a combination of variables, some of which we can change (weights) and some of which we cannot (input data, output data, and the error logic itself).

```
weight = 0.5
goal_pred = 0.8
input = 0.5

for iteration in range(20):
    pred = input * weight
    error = (pred - goal_pred) ** 2
    direction_and_amount = (pred - goal_pred) * input
    weight = weight - direction_and_amount

    print("Error:" + str(error) + " Prediction:" + str(pred))
```

We can modify *anything* in our `pred` calculation except the `input`.

In fact, we're going to spend *the rest of this book* (and many deep learning researchers will spend *the rest of their lives*) just trying everything you can imagine on that `pred` calculation so that it can make good predictions. Learning is all about automatically changing that prediction function so that it makes good predictions – aka, so that the subsequent `error` goes down to 0.

Ok, now that we know what we're *allowed* to change, how do we actually go about doing the changing? That's the good stuff! That's the *machine learning*, right? In the next, section, we're going to talk about exactly that.

Tunnel Vision on One Concept

Concept: "Learning is adjusting our weight to reduce the error to zero"

So far in this chapter, we've been hammering on the idea that learning is really just about adjusting our **weight** to reduce our error to zero. This is the secret sauce. Truth be told, knowing how to do this is *all about* understanding the **relationship** between our **weight** and our **error**. If we understand this relationship, we can know how to adjust our **weight** to reduce our **error**.

What do I mean by "understand the relationship"? Well, to understand the relationship between two variables is really just to understand *how changing one variable changes the other*. In our case, what we're really after is the **sensitivity** between these two variables. Sensitivity is really just another name for *direction* and *amount*. We want to know how sensitive the **error** is to the **weight**. We want to know the *direction* and the *amount* that the **error** changes when we change the **weight**. This is the goal. So far, we've used two different methods to attempt to understand this relationship.

You see, when we were "wiggling" our **weight** (hot and cold learning) and studying its affect on our **error**, we were really just *experimentally* studying the relationship between these two variables. It's like when you walk into a room with 15 different unlabeled light switches. You just start flipping them on and off to learn about their relationship to various lights in the room. We did the same thing to study the relationship between our **weight** and our **error**. We just wiggled the **weight** up and down and watched for how it changed the **error**. Once we knew the relationship, we could move the **weight** in the right direction using two simple if statements:

```
if(down_error < up_error):
    weight = weight - step_amount

if(down_error > up_error):
    weight = weight + step_amount
```

Now, let's go back to the formula from the previous pages, where we combined our pred and error logic. As mentioned, they quietly define an *exact relationship* between our **error** and our **weight**.

```
error = ((input * weight) - goal_pred) ** 2
```

This line of code, ladies and gentlemen, is the secret. This is a formula. This is the relationship between **error** and **weight**. This relationship is exact. It's computable. It's universal. It is and it will always be. Now, how can we use this formula to know how to change our **weight** so that our **error** moves in a *particular direction*. Now THAT is the right question! Stop. I beg you. Stop and appreciate this moment. This formula is the exact relationship between these two variables, and now we're going to figure out how to change one variable so that we move the other variable in a particular direction. As it turns out, there's a method for doing this for *any* formula. We're going to use it for reducing our error.

A Box With Rods Poking Out of It

An analogy:

Picture yourself sitting in front of a cardboard box that has two circular rods sticking through two little holes. The blue rod is sticking out of the box by 2 inches, and the red rod is sticking out of the box by 4 inches. Imagine that I told you that these rods were connected, but I wouldn't tell you in what way. You had to experiment to figure it out.

So, you take the blue rod and push it in 1 inch, and watch as... while you're pushing... the red rod also moves into the box by 2 inches!!! Then, you pull the blue rod back out an inch, and the red rod follows again, pulling out by 2 inches!! What did you learn? Well, there seems to be a *relationship* between the red and blue rods. However much you move the blue rod, the red rod will move by twice as much. You might say the following is true.

```
red_length = blue_length * 2
```

As it turns out, there's a formal definition for "when I tug on this part, how much does this other part move?". It's called a **derivative** and all it really means is "how much does rod X move when I tug on rod Y?"

In the case of the rods above, the derivative for "how much does red move when I tug on blue?" is 2. Just 2. Why is it 2? Well, that's the *multiplicative* relationship determined by the formula.

```
red_length = blue_length * 2
```

derivative

Notice that we always have the derivative *between two variables*. We're always looking to know how one variable moves when we change another one! If the derivative is *positive* then when we change one variable, the other will move in the *same* direction! If the derivative is *negative* then when we change one variable, the other will move in the *opposite* direction.

Consider a few examples. Since the derivative of red_length compared to blue_length is 2, then both numbers move in the same direction! More specifically, red will move *twice as much* as blue in the same direction. If the derivative had been -1, then red would move in the *opposite* direction by the same amount. Thus, given a function, the derivative represents the **direction** and the **amount** that one variable changes if you change the other variable. This is exactly what we were looking for!

Derivatives: Take Two

Still a little unsure about them? Let's take another perspective...

There are two ways I've heard people explain derivatives. One way is all about understanding "how one variable in a function changes when you move another variable". The other way of explaining it is "a derivative is the slope at a point on a line or curve". As it turns out, if you take a function and plot it out (draw it), the slope of the line you plot is the *same thing* as "how much one variable changes when you change the other". Let me show you by plotting our favorite function.

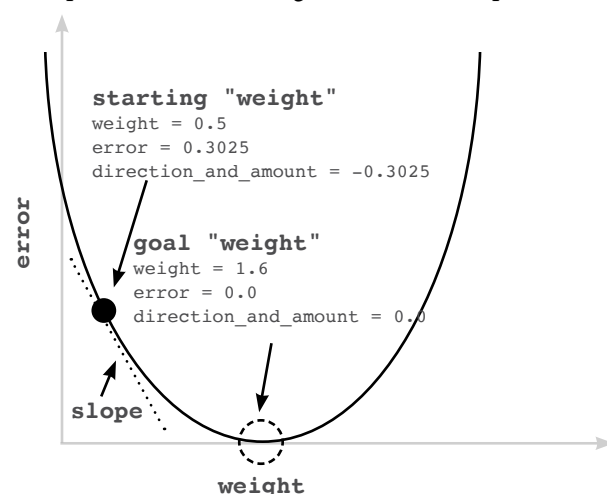
```
error = ((input * weight) - goal_pred) ** 2
```

Now remember... our `goal_pred` and `input` are fixed, so we can rewrite this function:

```
error = ((0.5 * weight) - 0.8) ** 2
```

Since there are only two variables left that actually change (all of the rest of them are fixed), we can just take every `weight` and compute the `error` that goes with it! Let's plot them:

As you can see on the right, our plot looks like a big U shaped curve! Notice that there is also a point in the middle where the `error == 0`! Also notice that to the right of that point, the slope of the line is *positive*, and to the left of that point, the slope of the line is *negative*. Perhaps even more interesting, the farther away from the **goal weight** that you move, the *steeper* the slope gets. We like all of these properties. The *slope's sign* gives us **direction** and the slope's *steepness* gives us **amount**. We can use both of these to help find the goal **weight**.



Even now, when I look at that curve, it's easy for me to lose track of what it represents. It's actually similar to our "hot and cold" method for learning. If we just tried *every possible value* for `weight` and plotted it out, we'd get this curve. And what's really remarkable about derivatives is that they can see past our big formula for computing `error` (at the top of this page) and see this curve! We can actually compute the **slope** (i.e. derivative) of the line for any value of `weight`. We can then use this slope (derivative) to figure out which **direction** reduces our `error`! Even better, based on the *steepness* we can get at least some idea for how far away we are (although not an exact one, as we'll learn more about later).

What you really need to know...

With derivatives we can pick any two variables in any formula, and know how they interact.

Take a look at this *big whopper of a function*:

$$y = (((\text{beta} * \text{gamma}) ** 2) + (\text{epsilon} + 22 - x)) ** (1/2)$$

Here's what you need to know about derivatives: For any function (even this whopper) you can pick any two variables and understand their relationship with each other. For any function, you can pick two variables and plot them on an x-y graph like we did on the last page. For any function, you can pick two variables and compute how much one changes when you change the other. Thus, for any function, we can learn how to change one variable so that we can move another variable in a direction. Sorry to harp on, but it's important you know this in your bones.

Bottom Line: In this book we're going to build neural networks. A neural network is really just one thing: a bunch of **weights** which we use to compute an **error** function. And for any error function (no matter how complicated), we can compute the relationship between any **weight** and the final **error** of the network. With this information, we can change each **weight** in our neural network to reduce our **error** down to 0... and that's exactly what we're going to do.

What you don't really need to know...

... Calculus ...

So, it turns out that learning all of the methods for taking any two variables in any function and computing their relationship takes about 3 semesters of college. Truth be told, if you went through all three semesters so that you could learn how to do Deep Learning, you would only actually find yourself *using* a very small subset of what you learned. And really, Calculus is just about memorizing and practicing every possible derivative rule for every possible function.

So, in this book I'm going to do what I typically do in real life (cuz i'm lazy?... i mean... efficient?): just look up the derivative in a reference table. All you really *need to know* is what the derivative *represents*. It's the relationship between two variables in a function so that you can know how much one changes when you change the other. It's just the sensitivity between two variables. I know that was a lot of information to just say, "It's the sensitivity between two variables"... but it is. Note that this can include both "positive" sensitivity (when variables move together), and "negative" sensitivity (when they move in opposite directions), and "zero" sensitivity (where one stays fixed regardless of what you do to the other). For example, $y = 0 * x$. Move x ... and y is always 0. Ok, enough about derivatives. Let's get back to Gradient Descent.

How to use a Derivative to Learn

"weight_delta" is our derivative.

What is the difference between the **error**, and the derivative of our **error** and **weight**? The error is just a *measure* of how much we missed. The derivative defines the *relationship* between each weight and how much we missed. In other words, it tells *how much changing a weight contributed to the error*. So, now that we know this, how do we use it to move the error in a particular direction?

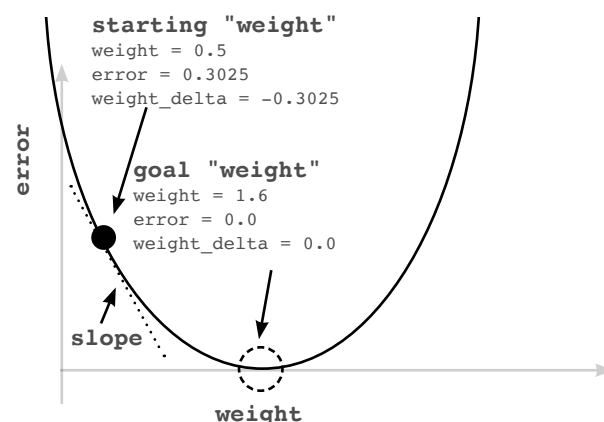
So, we've learned the relationship between two variables in a function, but how do we exploit that relationship? As it turns out, this is incredibly visual and intuitive. Check out our error curve again. The black dot is where our **weight** starts out at (0.5). The dotted circle is where we want it to go, our goal **weight**. Do you see the dotted line attached to our black dot? That's our *slope* otherwise known as our *derivative*. It tells us *at that point in the curve* how much the **error** changes when we change the **weight**. Notice that it's pointed downward! It's a negative slope!

The slope of a line or curve *always* points in the opposite direction of the *lowest point* of the line or curve. So, if you have a negative slope, you *increase* your **weight** to find the minimum of the **error**. Check it out!

So, how do we use our *derivative* to find the **error** minimum (lowest point in the **error** graph)? We just move the opposite direction of the slope! We move in the opposite direction of the derivative! So, we can take each weight, calculate the derivative of that weight with respect to the error (so we're comparing two variables: the weight and the error) and then change the weight in the *opposite* direction of that slope! That will move us to the minimum!

Let's remember back to our goal again. We are trying to figure out the **direction** and the **amount** to change our **weight** so that our **error** goes down. A derivative gives us the relationship between any two variables in a function. We use the derivative to determine the relationship between any **weight** and the **error**. We then move our **weight** in the *opposite* direction of the derivative to find the lowest **weight**. Voilà! Our neural network learns!

This method for learning (finding error minimums) is called **Gradient Descent**. This name should seem intuitive! We move the **weight** value *opposite the gradient* value, which *descends* our error to 0. By *opposite*, I simply mean that we increase our **weight** when we have a negative gradient and vice versa. It's like gravity!



Look Familiar?

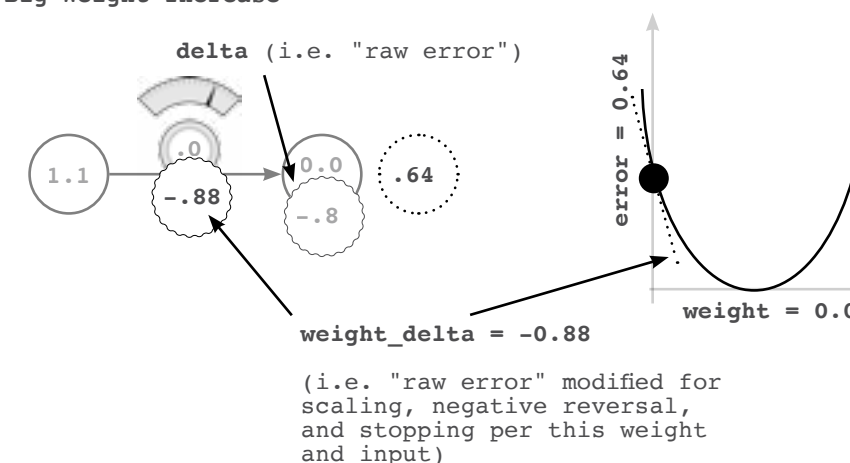
```
weight = 0.0
goal_pred = 0.8
input = 1.1
```

```
for iteration in range(4):
    pred = input * weight
    error = (pred - goal_pred) ** 2
    delta = pred - goal_pred
    weight_delta = delta * input
    weight = weight - weight_delta

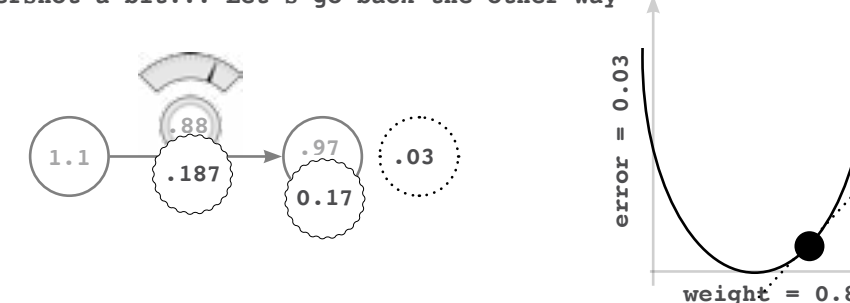
print("Error:" + str(error) + " Prediction:" + str(pred))
```

derivative
 (i.e., how fast the error changes given changes in the weight)

① A Big Weight Increase



② Overshot a bit... Let's go back the other way



Breaking Gradient Descent

Just Give Me The Code

```
weight = 0.5
goal_pred = 0.8
input = 0.5

for iteration in range(20):
    pred = input * weight
    error = (pred - goal_pred) ** 2
    delta = pred - goal_pred
    weight_delta = input * delta
    weight = weight - weight_delta
    print("Error:" + str(error) + " Prediction:" + str(pred))
```

When I run this code, I see the following output:

```
Error:0.3025 Prediction:0.25
Error:0.17015625 Prediction:0.3875
Error:0.095712890625 Prediction:0.490625
...
Error:1.7092608064e-05 Prediction:0.79586567925
Error:9.61459203602e-06 Prediction:0.796899259437
Error:5.40820802026e-06 Prediction:0.797674444578
```

Now that it works... let's break it! Play around with the starting `weight`, `goal_pred`, and `input` numbers. You can set them all to just about anything and the neural network will figure out how to predict the output given the input using the weight. See if you can find some combinations that the neural network cannot predict! I find that trying to break something is a great way to learn about it.

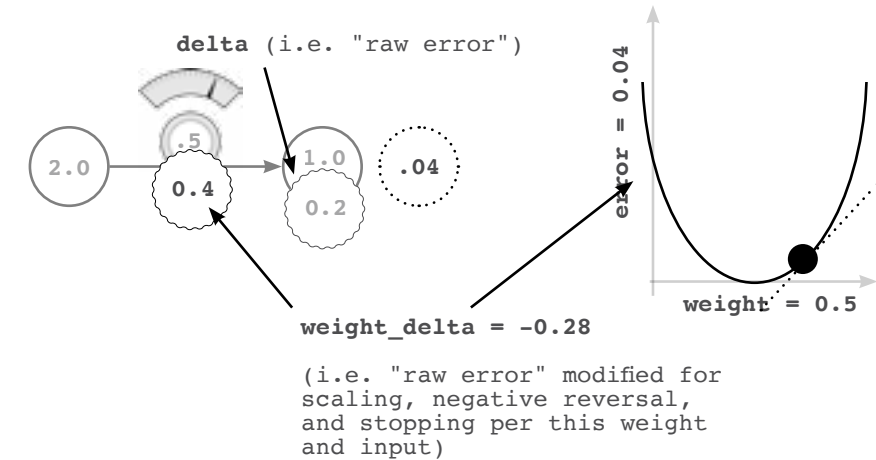
Let's try setting `input` to be equal to 2, but still try to get the algorithm to predict 0.8. What happens? Well, take a look at the output.

```
Error:0.04 Prediction:1.0
Error:0.36 Prediction:0.2
Error:3.24 Prediction:2.6
...
Error:6.67087267987e+14 Prediction:-25828031.8
Error:6.00378541188e+15 Prediction:77484098.6
Error:5.40340687069e+16 Prediction:-232452292.6
```

Woah! That's not what we wanted! Our predictions exploded! They alternate from negative to positive and negative to positive, getting farther away from the true answer at every step! In other words, every update to our weight **overcorrects**! In the next section, we'll learn more about how to combat this phenomenon.

Visualizing the Overcorrections

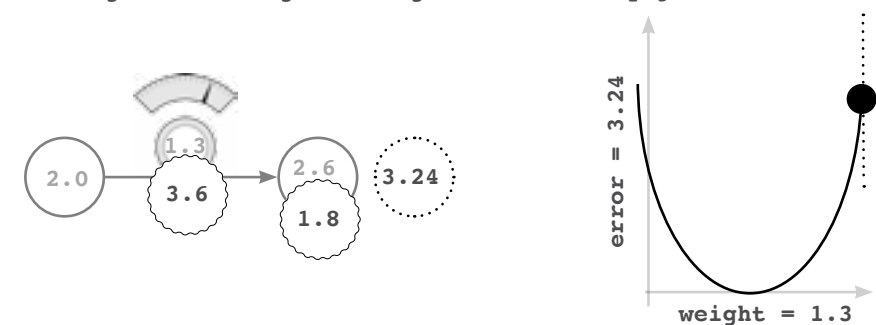
① A Big Weight Increase



② Overshot a bit... Let's go back the other way

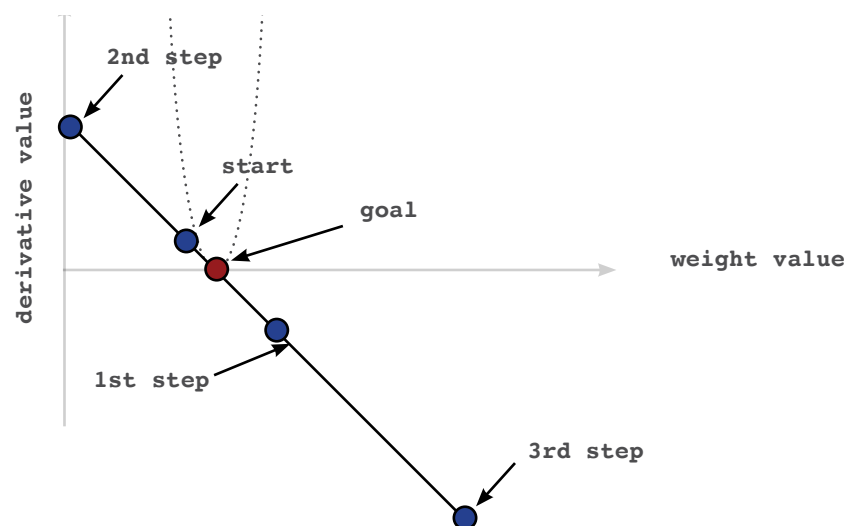


③ Overshot Again! Let's go back again... but only just a little



Divergence

Sometimes neural networks explode in value... Oops?



So what really happened? The explosion in error on the previous page is caused by the fact that we made the input larger. Consider how we're updating our weight:

$$\text{weight} = \text{weight} - (\text{input} * (\text{pred} - \text{goal_pred}))$$

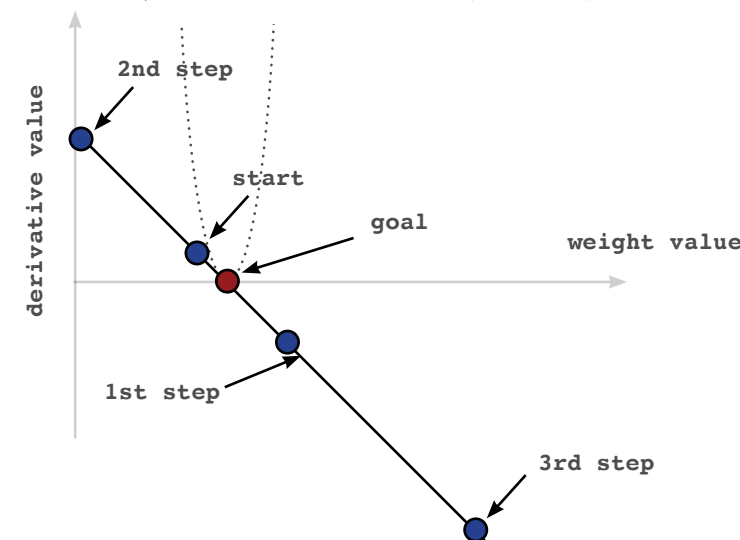
If our input is sufficiently large, this can make our weight update large even when our error is small. What happens when you have a large weight update and a small error? It overcorrects!!! If the new error is even bigger, it overcorrects even more!!! This causes the phenomenon that we saw on the previous page, called **divergence**.

You see, if we have a BIG input, then the prediction is very sensitive to changes in the weight (since $\text{pred} = \text{input} * \text{weight}$). This can cause our network to overcorrect. In other words, even though our **weight** is still only starting at 0.5, our derivative at that point is *very steep*. See how tight the u shaped error curve is in the graph above?

This is actually really intuitive. How do we predict? Well, we predict by *multiplying* our input by our **weight**. So, if our input is *huge*, then small changes in our **weight** are going to cause BIG changes in our prediction!! The error is very *sensitive* to our **weight**. Aka... the derivative is really big! So, how do we make it smaller?

Introducing: Alpha

The simplest way to prevent overcorrecting our weight updates.



So, what was the problem we're trying to solve? The problem is that if the input is too big, then our weight update can overcorrect. What is the symptom? The symptom is that when we overcorrect, our new derivative is even *larger in magnitude* than when we started (although the sign will be the opposite). Stop and consider this for a second. Look at the graph above to understand the symptom. The 2nd step is even farther away from the goal, which means the *derivative* is even greater in magnitude! This causes the 3rd step to be even farther away from the goal than the second step, and the neural network continues like this, demonstrating **divergence**.

The symptom is this overshooting. The solution is to *multiply the weight update by a fraction* to make it smaller. In most cases, this involves multiplying our weight update by a single real-valued number between 0 and 1, known as **alpha**. One might note, this has no effect on the *core issue*, which is that our input is larger. It will also reduce the **weight** updates for inputs that aren't too large. In fact, finding the appropriate alpha, even for state-of-the-art neural networks, is often done simply by guessing. You watch your **error** over time. If it starts diverging (going up), then your alpha is too high, and you decrease it. If learning is happening too slowly, then your alpha is too low, and you increase it. There are other methods than simple gradient descent that attempt to counter for this, but gradient descent is still very popular.

Alpha In Code

Where does our "alpha" parameter come in to play?

So we just learned that **alpha** reduces our weight update so that it doesn't overshoot. How does this affect our code? Well, we were updating our weights according to the following formula:

$$\text{weight} = \text{weight} - \text{derivative}$$

Accounting for alpha is a rather small change, pictured below. Notice that if alpha is small (say, 0.01), it will reduce our weight update considerably, thus preventing it from overshooting.

$$\text{weight} = \text{weight} - (\text{alpha} * \text{derivative})$$

Well, that was easy! So, let's install alpha into our tiny implementation from the beginning of this chapter and run it where input = 2 (which previously didn't work):

```
weight = 0.5
goal_pred = 0.8
input = 2
alpha = 0.1

for iteration in range(20):
    pred = input * weight
    error = (pred - goal_pred) ** 2
    derivative = input * (pred - goal_pred)
    weight = weight - (alpha * derivative)

    print("Error:" + str(error) + " Prediction:" + str(pred))
```

```
Error:0.04 Prediction:1.0
Error:0.0144 Prediction:0.92
Error:0.005184 Prediction:0.872
```

...

```
Error:1.14604719983e-09 Prediction:0.800033853319
Error:4.12576991939e-10 Prediction:0.800020311991
Error:1.48527717099e-10 Prediction:0.800012187195
```

What happens
when you make
alpha crazy
small or big?
What about
making it
negative?

Voilà! Our tiniest neural network can now make good predictions again! How did I know to set alpha to 0.1? Well, to be honest, I just tried it and it worked. And despite all the crazy advancements of deep learning in the past few years, most people just try several orders of magnitude of alpha (10,1,0.1,0.01,0.001,0.0001) and then tweak from there to see what works best. It's more art than science. There are more advanced ways which we can get to later, but for now, just try various alphas until you get one that seems to work pretty well. Play with it!

Memorizing

Ok, it's time to really learn this stuff.

This may sound like something that's a bit intense, but I can't stress enough the value I have found from this exercise. The code on the previous page, see if you can build it in an iPython notebook (or a .py file if you must) from *memory*. I know that might seem like overkill, but I (personally) didn't have my *click* moment with neural networks until I was able to perform this task.

Why does this work? Well, for starters, the only way to *know* that you have gleaned all the information necessary from this chapter is to try to produce it just from your head. Neural networks have lots of small moving parts, and it's easy to miss one.

Why is this important for the rest of the chapters? In the following chapters, I will be referring to the concepts discussed in this chapter at a *faster pace* so that I can spend plenty of time on the newer material. It is *vitaly important* that when I say something like "add your alpha parameterization to the weight update" that it is at least immediately apparent to which concepts from this chapter I'm referring.

All that is to say, memorizing small bits of neural network code has been hugely beneficial for me personally, as well as to many individuals who have taken my advice on this subject in the past.

IN THIS CHAPTER •

- Gradient Descent Learning with Multiple Inputs
- Freezing One Weight - What does it do?
- Gradient Descent Learning with Multiple Outputs
- Gradient Descent Learning with Multiple Inputs and Outputs
- Visualizing Weight Values
- Visualizing Dot Products

“ You don't learn to walk by following rules. You learn by doing, and by falling over. ”

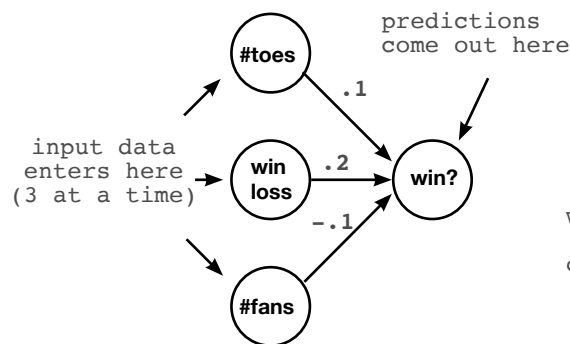
— RICHARD BRANSON

Gradient Descent Learning with Multiple Inputs

Gradient Descent Also Works with Multiple Inputs

In the last chapter, we learned how to use Gradient Descent to update a weight. In this chapter, we will more or less reveal how the same techniques can be used to update a network that contains multiple weights. Let's start by just jumping in the deep end, shall we? The following diagram lists out how a network with multiple inputs can learn!

1 An Empty Network With Multiple Inputs

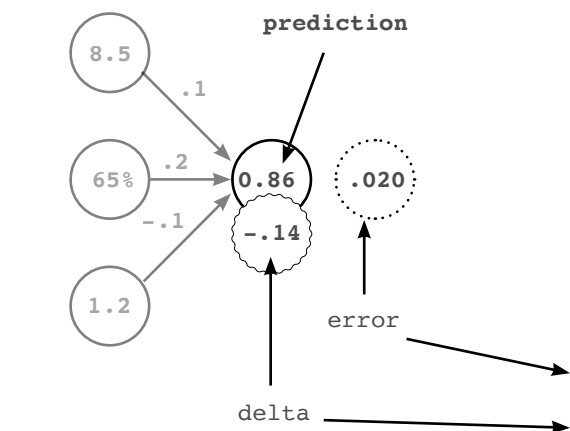


```
def w_sum(a,b):
    assert(len(a) == len(b))
    output = 0
    for i in range(len(a)):
        output += (a[i] * b[i])
    return output

weights = [0.1, 0.2, -.1]

def neural_network(input, weights):
    pred = w_sum(input,weights)
    return pred
```

2 PREDICT+COMPARE: Making a Prediction and Calculating Error And Delta



```
toes = [8.5 , 9.5, 9.9, 9.0]
wlrec = [0.65, 0.8, 0.8, 0.9]
nfans = [1.2 , 1.3, 0.5, 1.0]

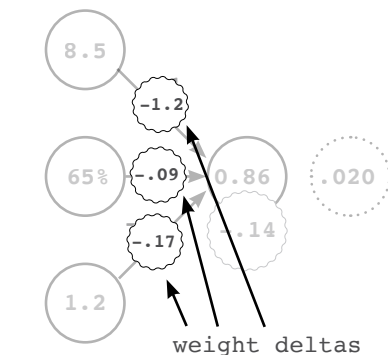
win_or_lose_binary = [1, 1, 0, 1]

true = win_or_lose_binary[0]

# Input corresponds to every entry
# for the first game of the season.

input = [toes[0],wlrec[0],nfans[0]]
pred = neural_network(input,weights)
error = (pred - true) ** 2
delta = pred - true
```

3 LEARN: Calculating Each "Weight Delta" and Putting It on Each Weight



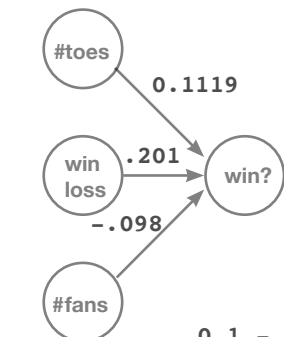
```
def ele_mul(number,vector):
    output = [0,0,0]
    assert(len(output) == len(vector))
    for i in range(len(vector)):
        output[i] = number * vector[i]
    return output
```

```
input = [toes[0],wlrec[0],nfans[0]]
pred = neural_network(input,weight)
error = (pred - true) ** 2
delta = pred - true
weight_deltas = ele_mul(delta,input)
```

There is actually nothing new in this diagram. Each weight delta is calculated by taking its output delta and multiplying it by its input. In this case, since our three weights share the same output node, they also share that node's delta. However, our weights have different weight deltas owing to their different input values. Notice further that we were able to re-use our ele_mul function from before as we are multiplying each value in weights by the same value delta.

$$\begin{aligned} 8.5 * -0.14 &= -1.19 = \text{weight_deltas}[0] \\ 0.65 * -0.14 &= -0.091 = \text{weight_deltas}[1] \\ 1.2 * -0.14 &= -0.168 = \text{weight_deltas}[2] \end{aligned}$$

4 LEARN: Updating the Weights



```
input = [toes[0],wlrec[0],nfans[0]]
pred = neural_network(input,weight)
error = (pred - true) ** 2
delta = pred - true

weight_deltas = ele_mul(delta,input)

alpha = 0.01

for i in range(len(weights)):
    weights[i] -= alpha * weight_deltas[i]
print("Weights:" + str(weights))
print("Weight Deltas:" + str(weight_deltas))
```

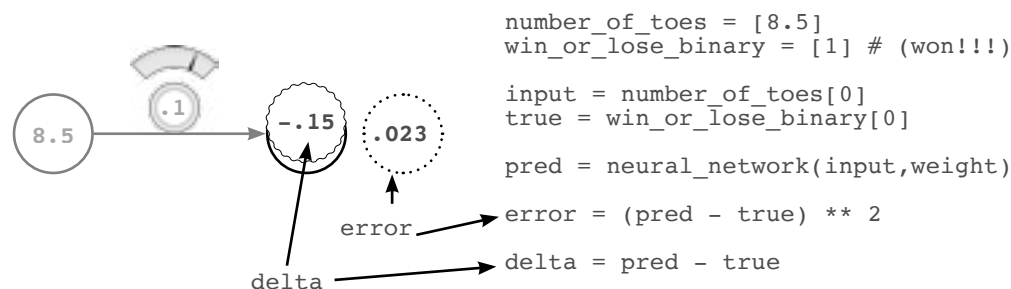
$$\begin{aligned} 0.1 - (-1.19 * 0.01) &= 0.1119 = \text{weights}[0] \\ 0.2 - (-0.091 * 0.01) &= 0.2009 = \text{weights}[1] \\ -0.1 - (-0.168 * 0.01) &= -0.098 = \text{weights}[2] \end{aligned}$$

Gradient Descent with Multiple Inputs - Explained

Simple to execute, fascinating to understand.

When put side by side with our single-weight neural network, Gradient Descent with multiple inputs seems rather obvious in practice. However, the properties involved are quite fascinating and worthy of discussion. First, let's take a look at them side by side.

② Single Input: Making a Prediction and Calculating Error and Delta



How do we turn a single delta (on the node) into 3 weight_delta values?

Let's remember what the definition and purpose of delta is vs weight_delta. Delta is a measure of "how much we want a node's value to be different". In this case, we compute it by a direct subtraction between the node's value and what we wanted the node's value to be (pred - true). Positive delta indicates the node's value was too high, and negative that it was too low.

delta

A measure of how much we want a node's value to be higher or lower to predict "perfectly" given the current training example.

weight_delta, on the other hand, is an *estimate* for the direction and amount we should move our weights to reduce our node delta, inferred by the derivative. How do we transform our delta into a weight_delta? We multiply delta by a weight's input.

weight_delta

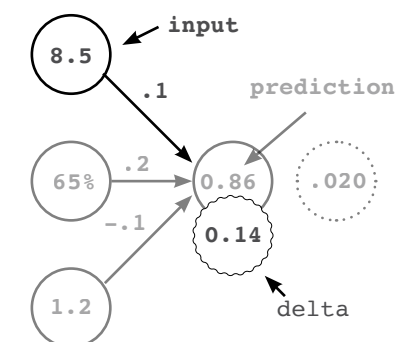
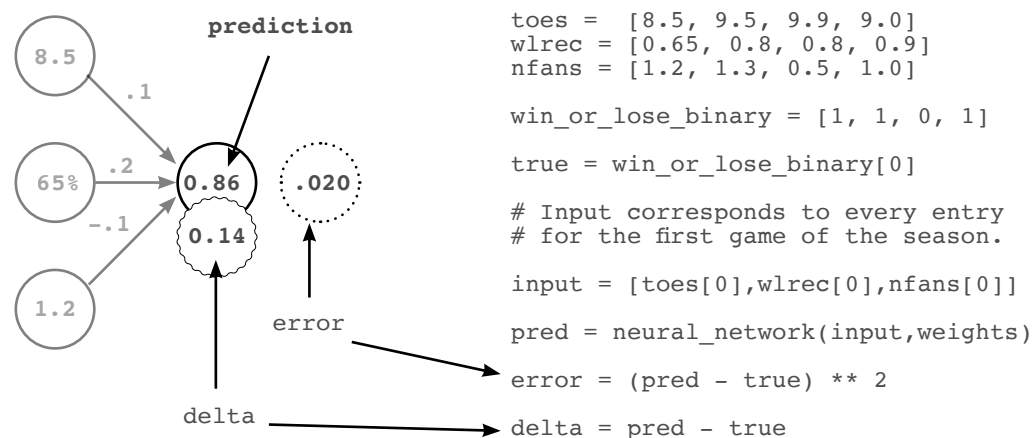
A derivative based *estimate* for the direction and amount we should move a weight to reduce our node_delta, accounting for scaling, negative reversal, and stopping.

Consider this from the perspective of a single weight, highlighted on the right. The delta says, "Hey inputs! ... Yeah you 3!!! Next time, predict a little higher!". Then, our single weight says, "hmm, if my input was 0, then my weight wouldn't have mattered and I wouldn't change a thing (stopping). If my input was negative, then I'd want to decrease my weight instead of increase (negative reversal). However, my input is positive and quite large, so I'm *guessing* that my personal prediction mattered a lot to the aggregated output, so I'm going to move my weight up a lot to compensate (scaling)!". It then increases its weight.

So, what did those three properties/statements really say? They all three (stopping, negative reversal, and scaling) made an observation of how the weight's role in the delta was affected by its input! Thus, each weight_delta is a sort of "input modified" version of the delta.

Bringing us back to our original question, how do we turn one (node) delta into three weight_delta values? Well, since each weight has a unique input and a shared delta, we simply use each respective weight's input multiplied by the delta to create each respective weight_delta. It's really quite simple. Let's see this process in action on the next page.

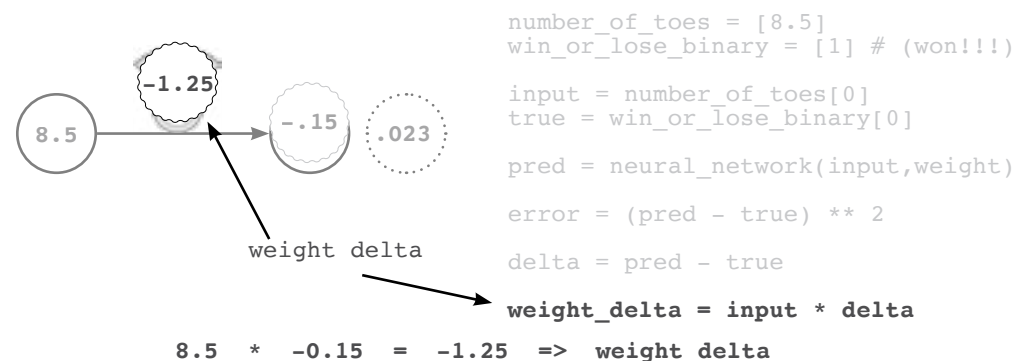
② Multi Input: Making a Prediction and Calculating Error And Delta



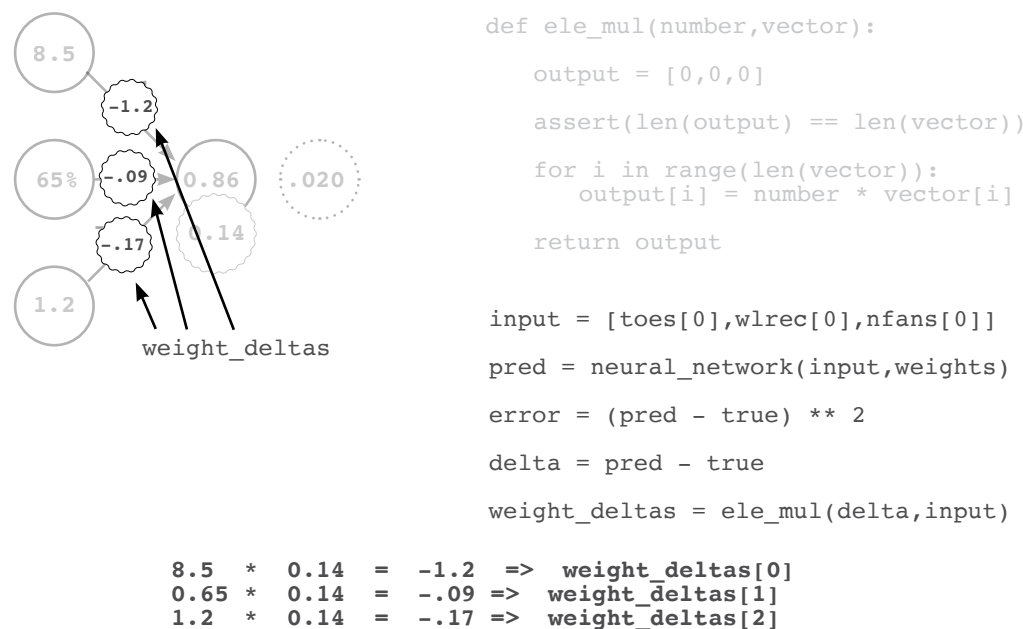
Indeed, up until the generation of "delta" on the output node, single input and multi-input Stochastic Gradient Descent are identical (other than the prediction differences we studied in Chapter 3). We make a prediction, and calculate the error and delta in identical ways. However, the following problem remains: when we only had one weight, we only had one input (one weight_delta to generate). Now we have 3! How do we generate 3 weight_deltas?

Below you can see the generation of `weight_delta` variables for the previous single-input architecture and for our new multi-input architecture. Perhaps the easiest way to see how similar they are is by reading the pseudo-code at the bottom of each section. Notice that the multi-weight version (bottom of the page), simply multiplies the delta (0.14) by every input to create the various `weight_deltas`. It's really quite a simple process.

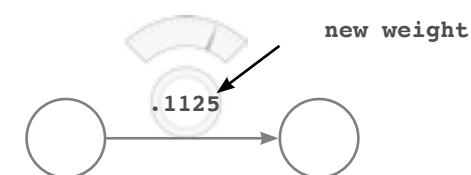
③ Single Input: Calculating "Weight Delta" and Putting it on the Weight



③ Multi: Calculating Each "Weight Delta" and Putting It on Each Weight



④ Updating the Weight



We multiply our `weight_delta` by a small number "alpha" before using it to update our weight. This allows us to control how fast the network learns. If it learns too fast, it can update weights too aggressively and overshoot. Note that the weight update made the same change (small increase) as Hot and Cold Learning.

```

number_of_toes = [8.5]
win_or_lose_binary = [1] # (won!!!)

input = number_of_toes[0]
true = win_or_lose_binary[0]

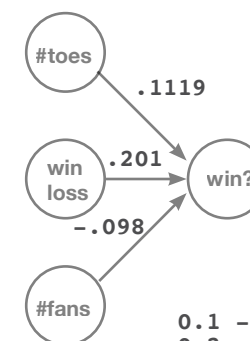
pred = neural_network(input,weight)

error = (pred - true) ** 2
delta = pred - true
weight_delta = input * delta

alpha = 0.01 # fixed before training
weight -= weight_delta * alpha

```

④ Updating the Weights



```

input = [toes[0],wlrec[0],nfans[0]]
pred = neural_network(input,weights)
error = (pred - true) ** 2
delta = pred - true
weight_deltas = ele_mul(delta,input)
alpha = 0.01

for i in range(len(weights)):
    weights[i] -= alpha * weight_deltas[i]

0.1 - (1.19 * 0.01) = 0.1119 = weights[0]
0.2 - (.091 * 0.01) = 0.2009 = weights[1]
-0.1 - (.168 * 0.01) = -0.098 = weights[2]

```

Finally, the last step of our process is also nearly identical to the single-input network. Once we have our `weight_delta` values, we simply multiply them by alpha and subtract them from our weights. It's literally the same process as before, repeated across multiple weights instead of just a single one.

Let's Watch Several Steps of Learning

```
def neural_network(input, weights):
    out = 0
    for i in range(len(input)):
        out += (input[i] * weights[i])
    return out

def ele_mul(scalar, vector):
    out = [0,0,0]
    for i in range(len(out)):
        out[i] = vector[i] * scalar
    return out

toes = [8.5, 9.5, 9.9, 9.0]
wlrec = [0.65, 0.8, 0.8, 0.9]
nfans = [1.2, 1.3, 0.5, 1.0]

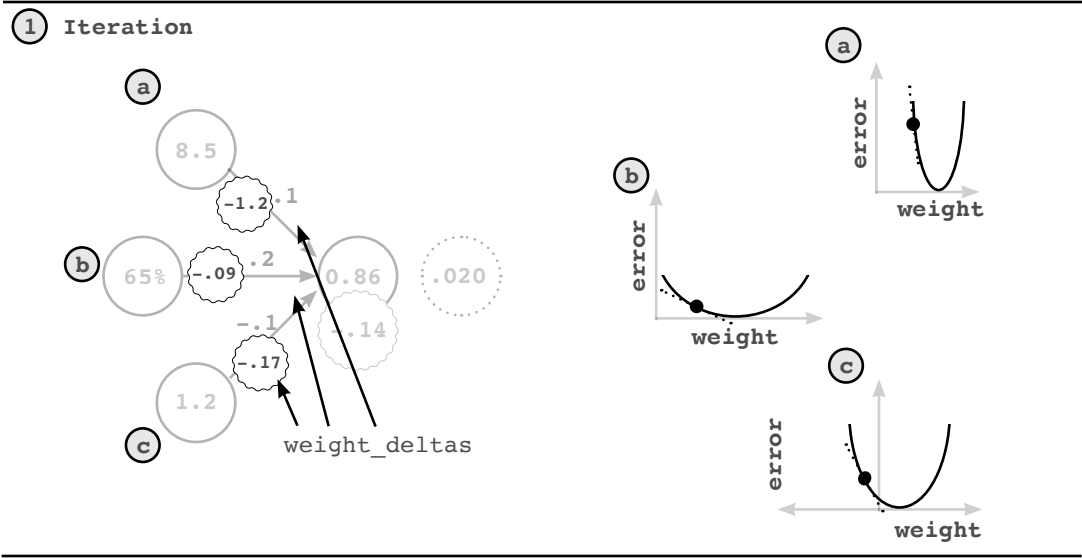
win_or_lose_binary = [1, 1, 0, 1]
true = win_or_lose_binary[0]

alpha = 0.01
weights = [0.1, 0.2, -.1]
input = [toes[0],wlrec[0],nfans[0]]

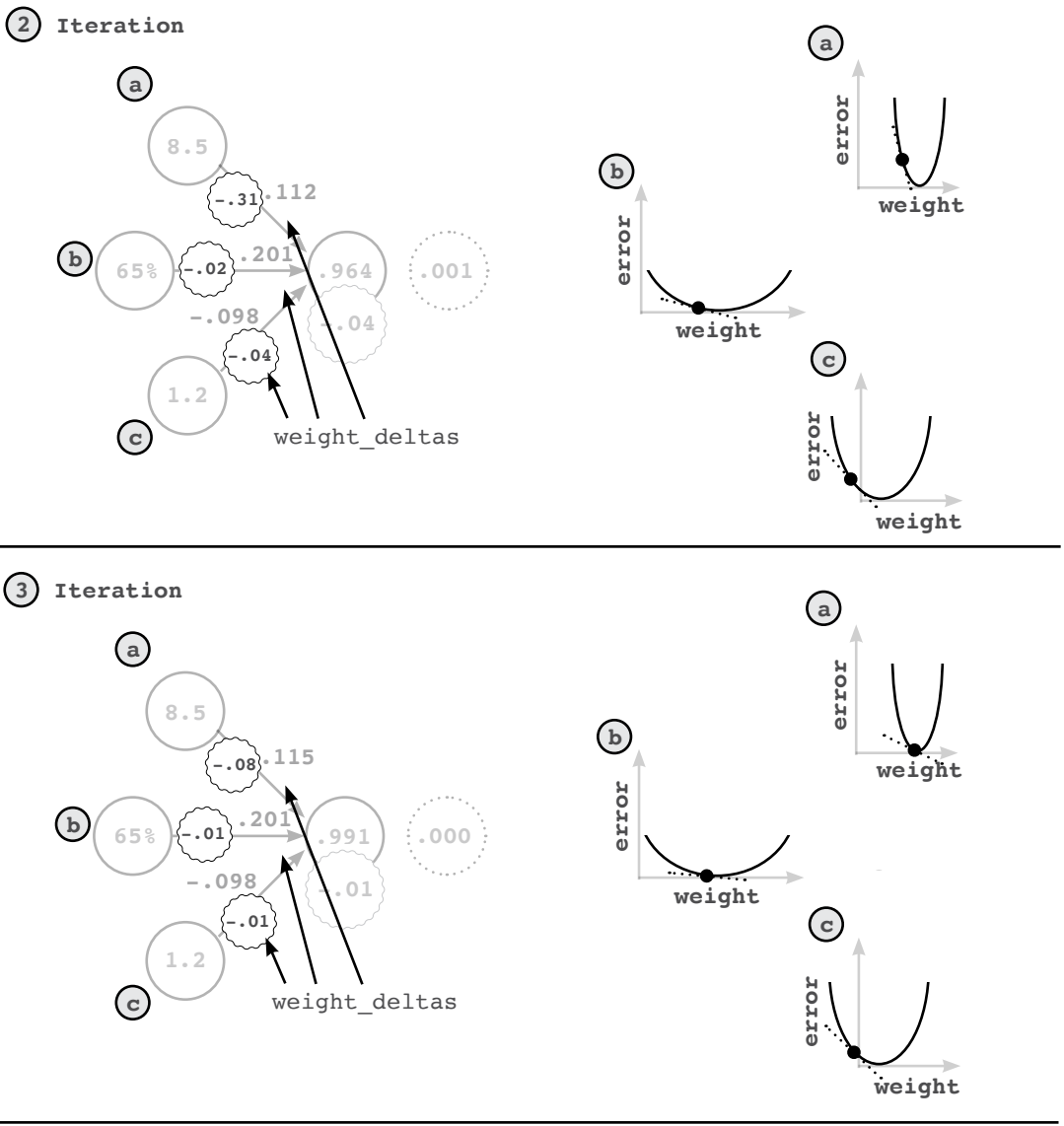
for iter in range(3):
    pred = neural_network(input,weights)
    error = (pred - true) ** 2
    delta = pred - true
    weight_deltas=ele_mul(delta,input)

    print("Iteration:" + str(iter+1))
    print("Pred:" + str(pred))
    print("Error:" + str(error))
    print("Delta:" + str(delta))
    print("Weights:" + str(weights))
    print("Weight_Deltas:")
    print(str(weight_deltas))

    for i in range(len(weights)):
        weights[i]-=alpha*weight_deltas[i]
```



Notice that on the right, we can picture three individual error/weight curves, one for each weight. As before, the slopes of these curves (the dotted lines) are reflected by the "weight_delta" values. Furthermore, notice that (a) is steeper than the others. Why is the weight_delta steeper for (a) than the others if they share the same output delta and error measure? Well, (a) has an input value that is significantly higher than the others. Thus, a higher derivative.



A few additional takeaways: most of the learning (weight changing) was performed on the weight with the largest input (a), because the input changes the slope significantly. This isn't necessarily advantageous in all settings. There is a sub-field called "normalization" that helps encourage learning across all weights despite dataset characteristics such as this. In fact, this significant difference in slope forced me to set the alpha to be lower than I wanted (0.01 instead of 0.1). Try setting alpha to 0.1. Do you see how (a) causes it to diverge?

Freezing One Weight - What Does It Do?

This experiment is perhaps a bit advanced in terms of theory, but I think that it's a great exercise to understand how the weights affect each other. We're going to train again, except weight (a) won't ever be adjusted. We'll try to learn the training example using only weights (b) and (c) (weights[1] and weights[2]).

```
def neural_network(input, weights):
    out = 0
    for i in range(len(input)):
        out += (input[i] * weights[i])
    return out

def ele_mul(scalar, vector):
    out = [0,0,0]
    for i in range(len(out)):
        out[i] = vector[i] * scalar
    return out

toes = [8.5, 9.5, 9.9, 9.0]
wlrec = [0.65, 0.8, 0.8, 0.9]
nfans = [1.2, 1.3, 0.5, 1.0]

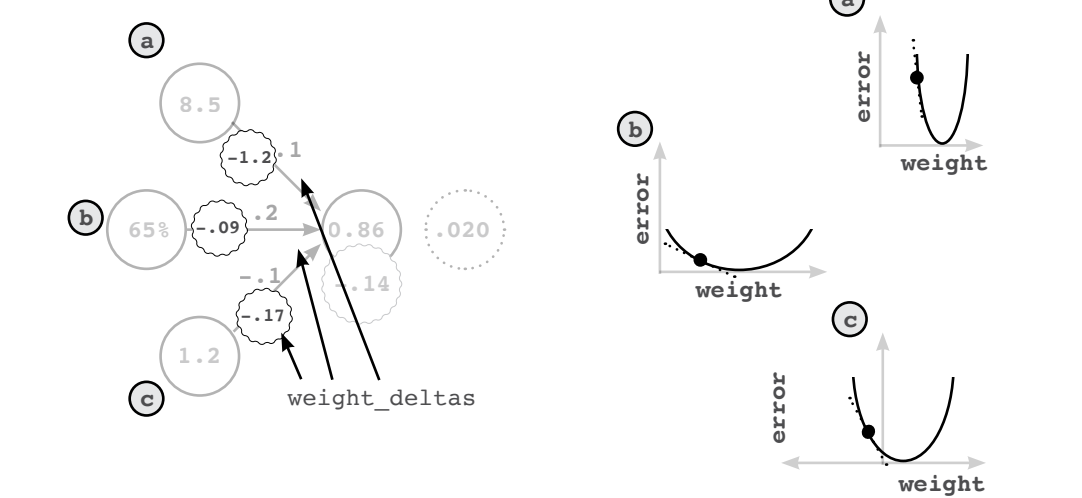
win_or_lose_binary = [1, 1, 0, 1]
true = win_or_lose_binary[0]

alpha = 0.3
weights = [0.1, 0.2, -.1]
input = [toes[0],wlrec[0],nfans[0]]

for iter in range(3):
    pred = neural_network(input,weights)
    error = (pred - true) ** 2
    delta = pred - true
    weight_deltas=ele_mul(delta,input)
    weight_deltas[0] = 0

    print("Iteration:" + str(iter+1))
    print("Pred:" + str(pred))
    print("Error:" + str(error))
    print("Delta:" + str(delta))
    print("Weights:" + str(weights))
    print("Weight Deltas:")
    print(str(weight_deltas))
    print(
    )
    for i in range(len(weights)):
        weights[i]-=alpha*weight_deltas[i]
```

1 Iteration



Perhaps you will be surprised to see that (a) still finds the bottom of the bowl? Why is this? Well, the curves are a measure of each individual weight relative to the global error. Thus, since the error is shared, when one weight finds the bottom of the bowl, all the weights find the bottom of the bowl.

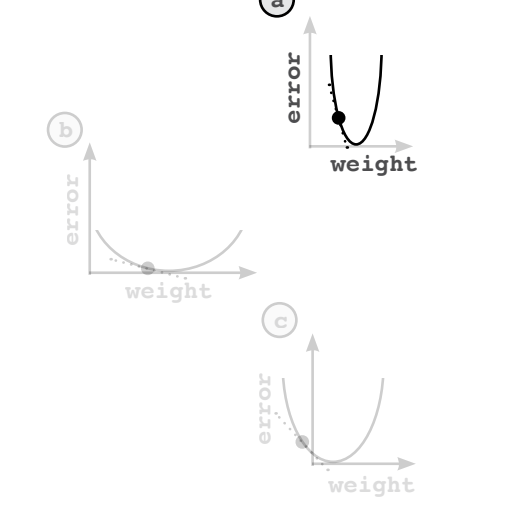
This is actually an extremely important lesson. First of all, if we converged (reached error = 0) with (b) and (c) weights and then tried to train (a), (a) wouldn't move! Why? error = 0 which means weight_delt is 0! This reveals a potentially damaging property of neural networks. (a) might be a really powerful input with lots of predictive power, but if the network accidentally figures out how to predict accurately on the training data without it, then it will never learn to incorporate (a) into its prediction.

Furthermore, notice "how" (a) finds the bottom of the bowl. Instead of the black dot moving, the curve seems to move to the left instead! What does this mean? Well, the black dot can only move horizontally if the weight is updated. Since the weight for (a) is frozen for this experiment, the dot must stay fixed. However, the error clearly goes to 0.

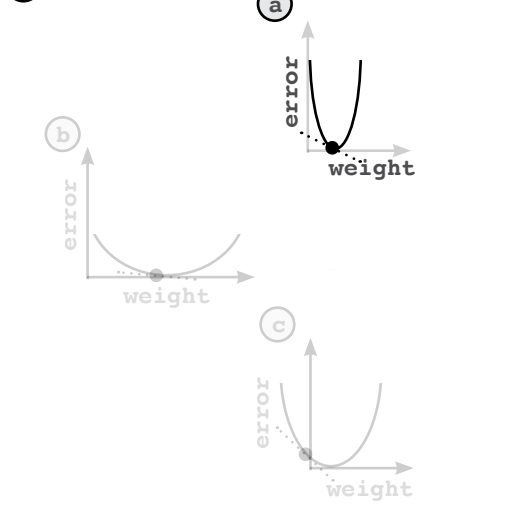
This tells us what the graphs really are. In truth, these are 2-d slices of a 4-dimensional shape. 3 of the dimensions are the weight values, and the 4th dimension is the error. This shape is called the "error plane" and, believe it or not, its curvature is determined by our training data! Why is it determined by our training data?

Well, our error is determined by our training data. Any network can have any weight value, but the value of the "error" given any particular weight configuration is 100% determined by data. We have already seen how the steepness of the "U" shape is affected by our input data (on several occasions). Truth be told, what we're really trying to do with our neural network is find the lowest point on this big "error plane", where the lowest point refers to the "lowest error". Interesting, eh? We're going to come back to this idea later, so just file it away for now.

2 Iteration



3 Iteration

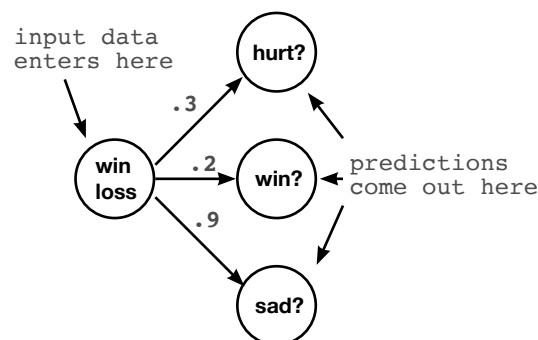


Gradient Descent Learning with Multiple Outputs

Neural Networks can also make multiple predictions using only a single input.

Perhaps this one will seem a bit obvious. We calculate each delta in the same way, and then multiply them all by the same, single input. This becomes each weight's weight_delta. At this point, I hope it is clear that a rather simple mechanism (Stochastic Gradient Descent) is consistently used to perform learning across a wide variety of architectures.

1 An Empty Network With Multiple Outputs

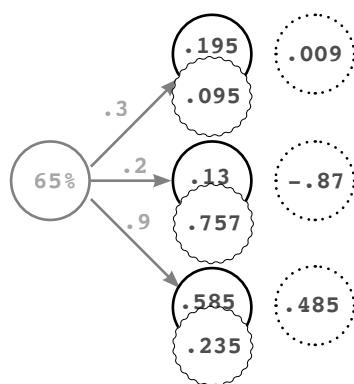


```
# Instead of predicting just
# whether the team won or lost,
# now we're also predicting whether
# they are happy/sad AND the
# percentage of the team that is
# hurt. We are making this
# prediction using only
# the current win/loss record.
```

```
weights = [0.3, 0.2, 0.9]
```

```
def neural_network(input, weights):
    pred = ele_mul(input, weights)
    return pred
```

2 PREDICT: Make a Prediction and Calculate Error and Delta



```
wlrec = [0.65, 1.0, 1.0, 0.9]
```

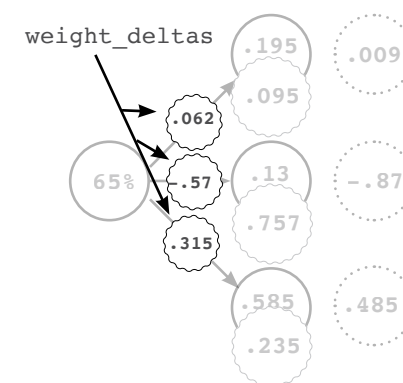
```
hurt = [0.1, 0.0, 0.0, 0.1]
win = [1, 1, 0, 1]
sad = [0.1, 0.0, 0.1, 0.2]
```

```
input = wlrec[0]
true = [hurt[0], win[0], sad[0]]
pred = neural_network(input, weights)
```

```
error = [0, 0, 0]
delta = [0, 0, 0]
```

```
for i in range(len(true)):
    error[i] = (pred[i] - true[i]) ** 2
    delta[i] = pred[i] - true[i]
```

3 COMPARE: Calculating Each "Weight Delta" and Putting It on Each Weight



```
def scalar_ele_mul(number, vector):
    output = [0, 0, 0]
    assert(len(output) == len(vector))
    for i in range(len(vector)):
        output[i] = number * vector[i]
    return output
```

```
wlrec = [0.65, 1.0, 1.0, 0.9]
```

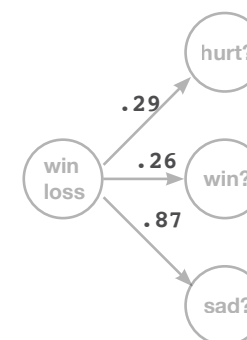
```
hurt = [0.1, 0.0, 0.0, 0.1]
win = [1, 1, 0, 1]
sad = [0.1, 0.0, 0.1, 0.2]
```

As before, weight_deltas are computed by multiplying the input node value with the output node delta for each weight. In this case, our weight_deltas share the same input node and have unique output node (deltas). Note also that we are able to re-use our ele_mul function.

```
input = wlrec[0]
true = [hurt[0], win[0], sad[0]]
pred = neural_network(input, weights)
error = [0, 0, 0]
delta = [0, 0, 0]
for i in range(len(true)):
    error[i] = (pred[i] - true[i]) ** 2
    delta[i] = pred[i] - true[i]
```

```
weight_deltas = scalar_ele_mul(input, weights)
```

4 LEARN: Updating the Weights



```
input = wlrec[0]
true = [hurt[0], win[0], sad[0]]
pred = neural_network(input, weights)
```

```
error = [0, 0, 0]
delta = [0, 0, 0]
```

```
for i in range(len(true)):
    error[i] = (pred[i] - true[i]) ** 2
    delta[i] = pred[i] - true[i]
```

```
weight_deltas = scalar_ele_mul(input, weights)
alpha = 0.1
```

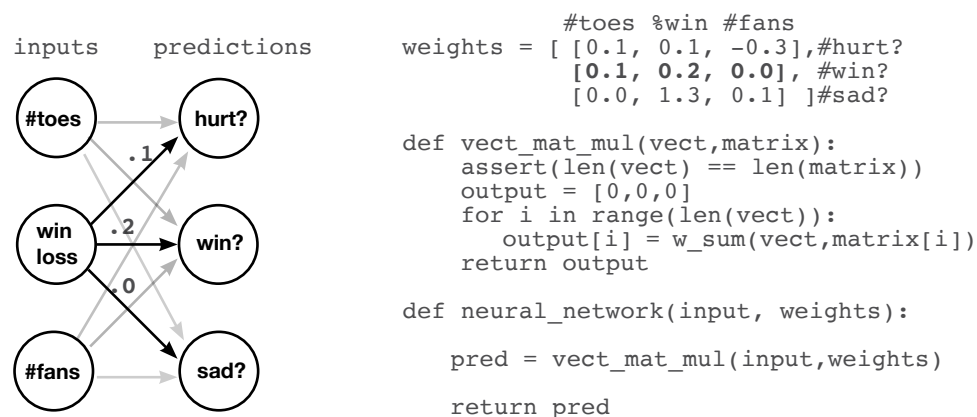
```
for i in range(len(weights)):
    weights[i] -= (weight_deltas[i] * alpha)
```

```
print("Weights:" + str(weights))
print("Weight Deltas:" + str(weight_deltas))
```

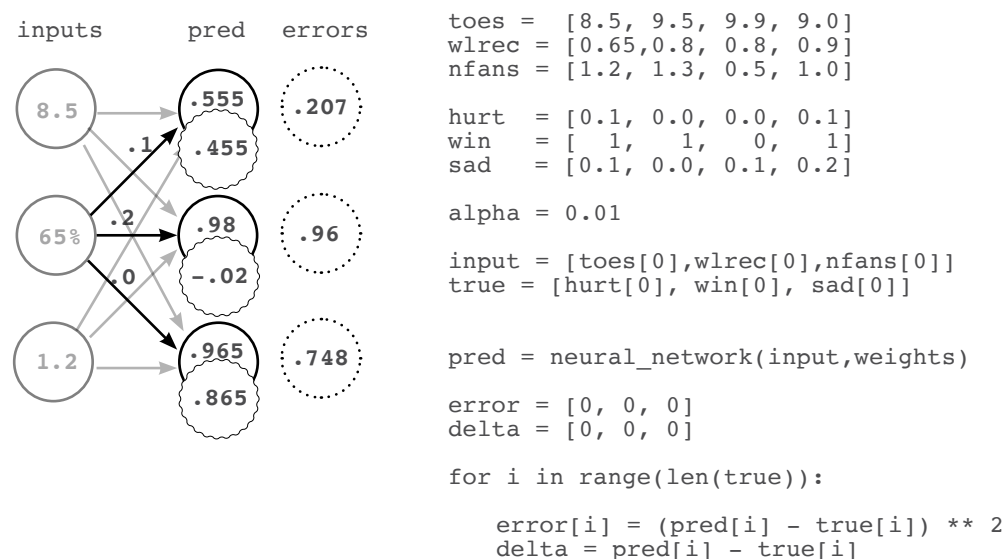
Gradient Descent with Multiple Inputs & Outputs

Gradient Descent generalizes to arbitrarily large networks.

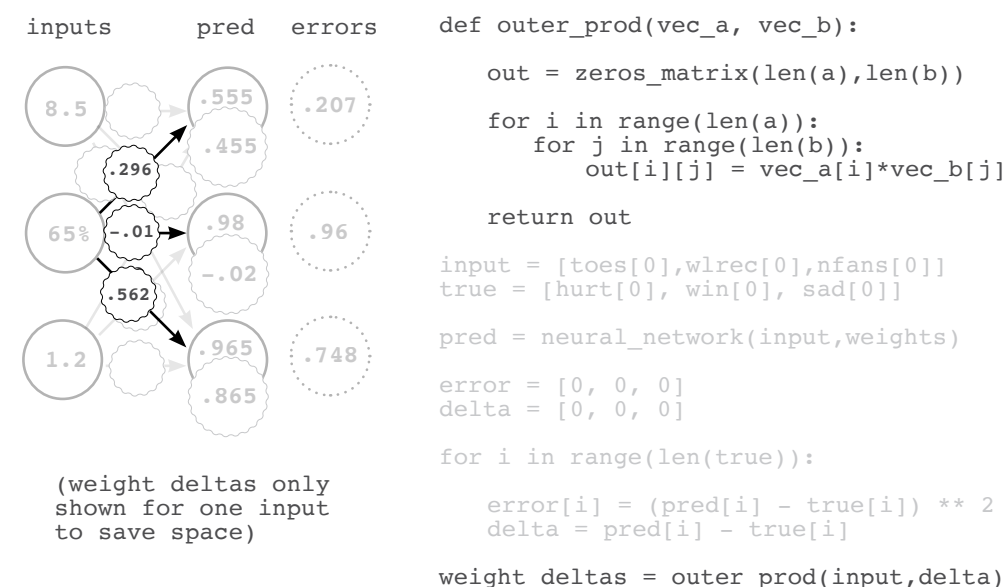
1 An Empty Network With Multiple Inputs & Outputs



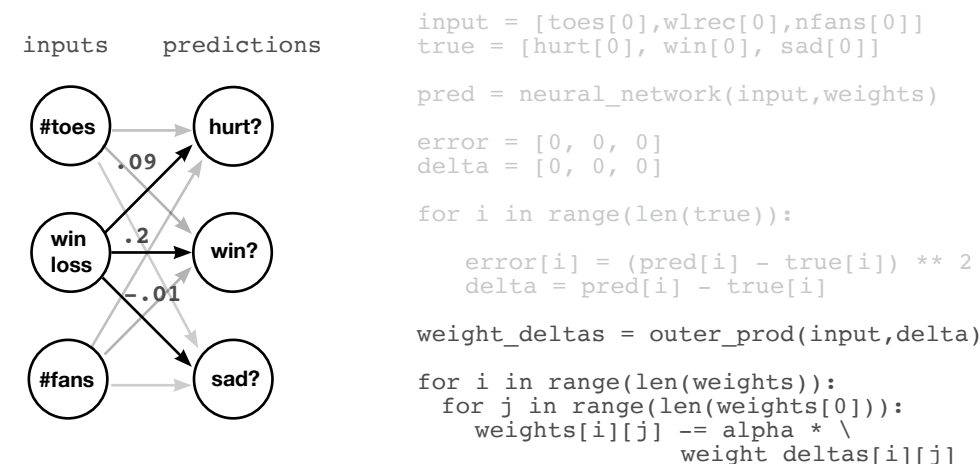
2 PREDICT: Make a Prediction and Calculate Error and Delta



3 COMPARE: Calculating Each "Weight Delta" and Putting It on Each Weight



4 LEARN: Updating the Weights



What do these weights learn?

Each weight tries to reduce the error, but what do they learn in aggregate?

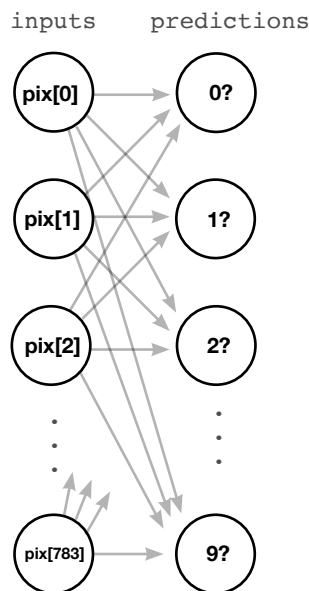
Congratulations! This is the part of the book where we move onto our first *real world dataset*. As luck would have it, it's one with historical significance!

Our new dataset is called the MNIST dataset, which is a dataset comprised of digits that high school students and employees of the US Census bureau hand wrote some years ago. The interesting bit is that these handwritten digits are simply black and white images of people's handwriting. Accompanying each digit image is the actual number that they were writing (0-9). For the last few decades, people have been using this dataset to train neural networks to read human handwriting, and today, you're going to do the same!

Each image is only 784 pixels (28 x 28). So, given that we have 784 pixels as input and 10 possible labels as output, you can imagine the shape of our neural network. So, now that each training example contains 784 values (one for each pixel), our neural network must have 784 input values. Pretty simple, eh? We just adjust the number of input nodes to reflect how many data points are in each training example. Furthermore, we want to predict 10 probabilities, one for each digit. In this way, given an input drawing, our neural network will produce these 10 probabilities, telling us which digit is most likely to be what was drawn.

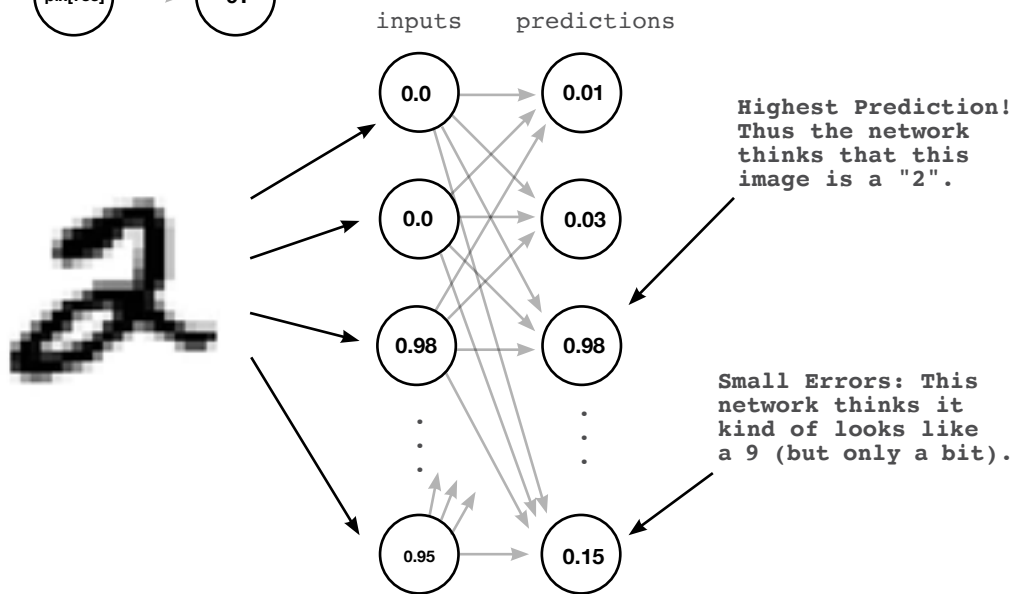
So, how do we configure our neural network to produce ten probabilities? Well, on the last page, we saw a diagram for a neural network that could take multiple inputs at a time and make multiple predictions based on that input. Thus, we should be able to simply modify this network to have the correct number of inputs and outputs for our new MNIST task. We'll just tweak it to have 784 inputs and 10 outputs.

In the notebook entitled "MNISTPreprocessor", you'll see a script to pre-process the MNIST dataset and load the first 1000 *images* and *labels* into two NumPy matrices called `images` and `labels`. You may be wondering, "images are 2-dimensional... How do we load the (28 x 28) pixels into a flat neural network?" For now, the answer is quite simple. We "flatten" the images into a vector of 1 x 784. So, we take the first row of pixels and concatenate them with the second row, and third row, and so on until we have one long list of pixels per image (784 pixels long in fact).



This picture on the left represents our new "MNIST Classification" neural network. It most closely resembles the network we trained with "Multiple Inputs and Outputs" a few pages ago. The only difference is the number of inputs and outputs, which has increased substantially. This network has 784 inputs (one for each pixel in a 28x28 image) and 10 outputs (one for each possible digit in the image).

If this network was able to predict perfectly, it would take in an image's pixels (say a 2 like the one on the previous page), and predict a 1.0 in the correct output position (the third one) and a 0 everywhere else. If it was able to do this correctly for all of the images in our dataset, it would have no error.



Over the course of training, the network will adjust the weights between the "input" and "prediction" nodes so that the error falls toward 0 in training. However, what does this actually do? What does it mean to modify a bunch of weights to learn a pattern in aggregate?

Visualizing Weight Values

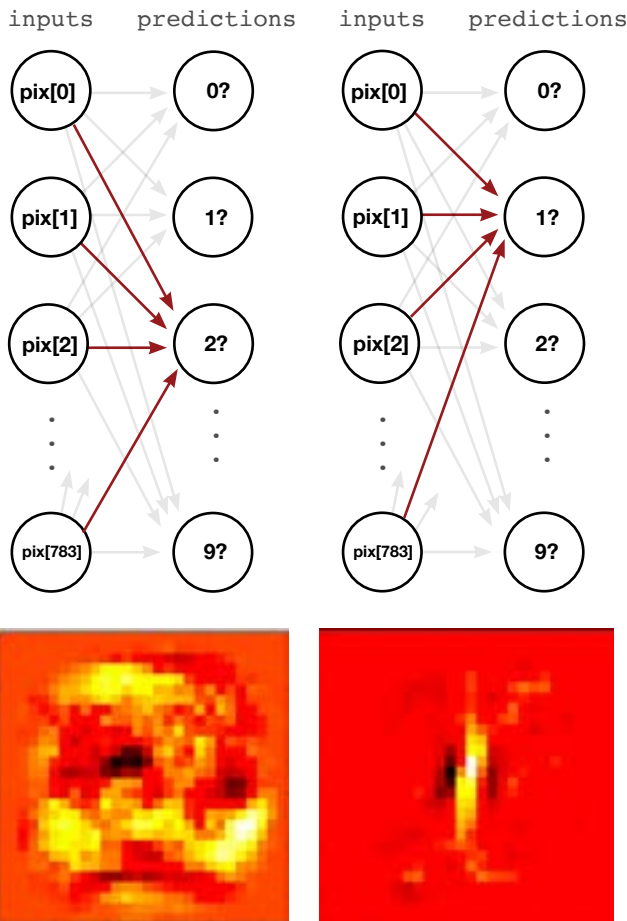
Each weight tries to reduce the error, but what do they learn in aggregate?

Perhaps an interesting and intuitive practice in neural network research (particularly for image classifiers) is to visualize the weights as if they were an image. If you look at the diagram on the right, you will see why.

Each output node has a weight coming from every pixel. For example, our "2?" node has 784 input weights, each mapping the relationship between a pixel and the number "2". What is this relationship? Well, if the weight is high, it means that the model believes there's a high degree of *correlation* between that pixel and the number 2. If the number is very low (negative), then the network believes there is a very low correlation (perhaps even negative correlation) between that pixel and the number two.

Thus, if we take our weights and print them out into an image that's the same shape as our input dataset images, we can "see" which pixels have the highest correlation with a particular output node. As you can see above, there is a very vague "2" and "1" in our two images, which were created using the weights for "2" and "1" respectively. The "bright" areas are high weights, and the dark areas are negative weights. The neutral color (red if you're reading this in color) represents 0 in the weight matrix. This illustrates that our network generally knows the shape of a 2 and of a 1.

Why does it turn out this way? Well, this takes us back to our lesson on "dot products". Let's have a quick review, shall we?



Visualizing Dot Products (weighted sums)

Each weight tries to reduce the error, but what do they learn in aggregate?

Recall how dot products work. They take two vectors, multiply them together (elementwise), and then sum over the output. So, in the example below:

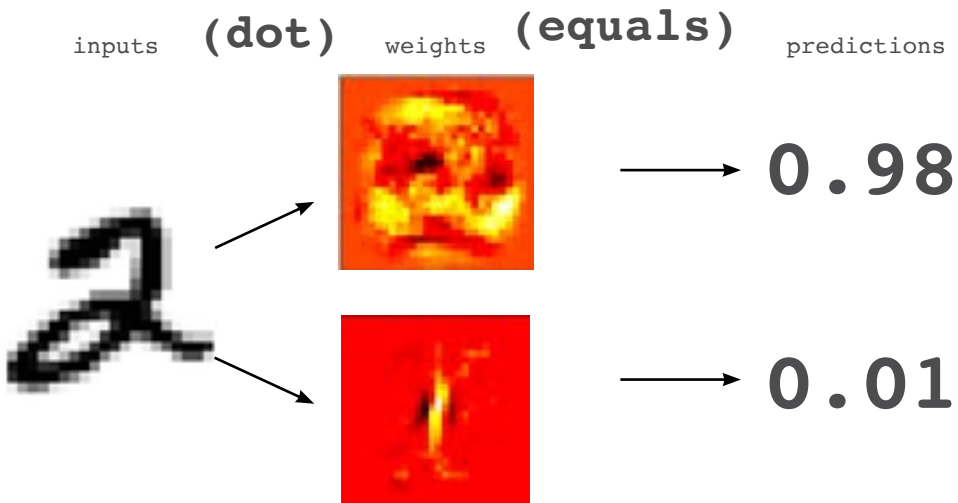
$$\begin{aligned} a &= [0, 1, 0, 1] \\ b &= [1, 0, 1, 0] \\ &\quad [0, 0, 0, 0] \rightarrow 0 \end{aligned} \quad \text{score}$$

First, you would multiply each element in a and b by each other, in this case creating a vector of 0s. The sum of this vector is also zero. Why? Because the vectors had nothing in common.

$$\begin{aligned} c &= [0, 1, 1, 0] & b &= [1, 0, 1, 0] \\ d &= [.5, 0, .5, 0] & c &= [0, 1, 1, 0] \end{aligned}$$

However, dot products between c and d return higher scores, because there is overlap in the columns that have positive values. Furthermore, performing dot products between two identical vectors tend to result in higher scores as well. The takeaway? **A dot product is a loose measurement of similarity between two vectors.**

What does this mean for our weights and inputs? Well, if our weight vector is similar to our input vector for "2", then it's going to output a high score because the two vectors are similar! Inversely, if our weight vector is NOT similar to our input vector for 2, then it's going to output a low score. You can see this in action below! Why is the top score (0.98) higher than the lower one (0.01)?



Conclusion

Gradient Descent is a General Learning Algorithm

Perhaps the most important subtext of this chapter is that Gradient Descent is a very flexible learning algorithm. If you combine weights together in a way that allows you to calculate an error function and a delta, gradient descent can show you how to move your weights to reduce your error. We will spend the rest of this book exploring different types of weight combinations and error functions for which Gradient Descent is useful. The next chapter is no exception.

Building Your First "Deep" Neural Network

Introduction to Backpropagation

6

IN THIS CHAPTER •

- The Streetlight Problem
- Matrices and the Matrix Relationship
- Full / Batch / Stochastic Gradient Descent
- Neural Networks Learn Correlation
- Overfitting
- Creating our Own Correlation
- Backpropagation: Long Distance Error Attribution
- Linear vs Non-Linear
- The Secret to Sometimes Correlation
- Our First "Deep" Network
- Backpropagation in Code / Bringing it all Together

“

Who invented backpropagation?

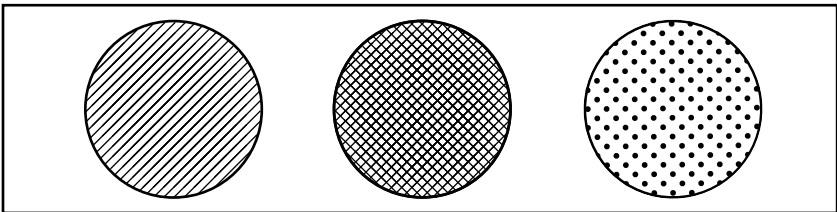
”

— JURGEN SCHMIDHUBER

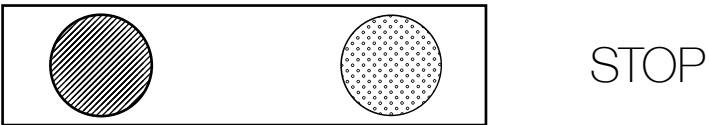
The Street Light Problem

A toy problem for us to consider how a network learns entire datasets.

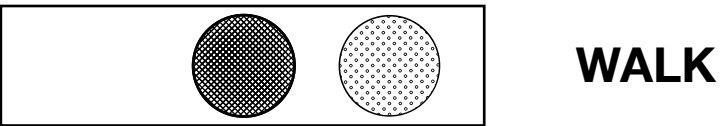
Consider yourself approaching a street corner in a foreign country. As you approach, you look up and realize that the street light is quite unfamiliar. How can you know when it is safe to cross the street?



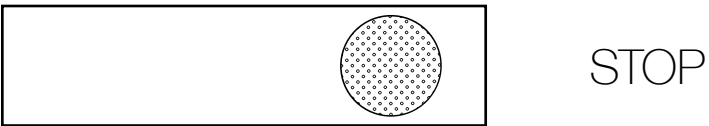
You can know when it is safe to cross the street by interpreting the streetlight. However, in this case, we don't know how to interpret it! Which "light" combinations indicate when it is time to **walk**? Which indicate when it is time to **stop**? To solve this problem, you might sit at the street corner for a few minutes observing correlation between each light combination and whether or not people around you choose to walk or stop. You take a seat and record the following pattern:



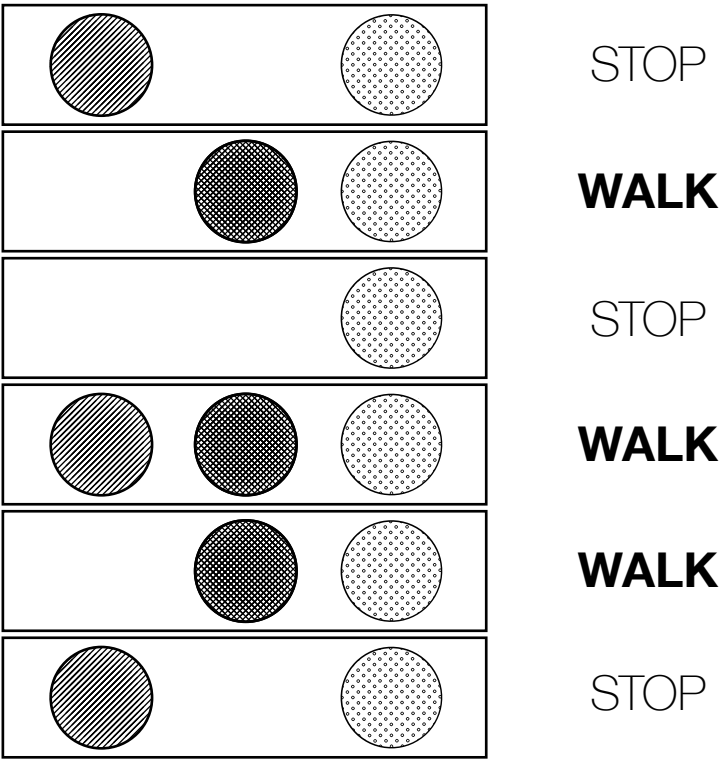
Ok, nobody walked at the first light. At this point you're thinking, "man, this pattern could be anything. The left light or the right light could be correlated with stopping, or the central light could be correlated with walking." There's no way to know. Let's take another datapoint.



People walked! Ok, so something changed with this light that changed the signal. The only thing we know for sure is that the far right light doesn't seem to indicate one way or another. Perhaps it is irrelevant. Let's collect another datapoint.



Now we're getting somewhere! Only the middle light changed this time, and we got the opposite pattern. Our working hypothesis is that the **middle** light indicates when people feel safe to walk. Over the next few minutes, we recorded the following six light patterns, noting when people seemed to walk or stop. Do you notice a pattern overall?



As hypothesized on the previous page, there is a **perfect correlation** between the middle (criss-cross) light and whether or not it is safe to walk. You were able to learn this pattern by observing all of the individual datapoints and *searching for correlation*. This is what we're going to train our neural network to do.

Preparing our Data







Neural Networks Don't Read Streetlight

In the previous chapters, we learned about supervised algorithms. We learned that they can take one dataset and turn it into another. More importantly, they can take a dataset of **what we know** and turn it into a dataset of **what we want to know**.

So, how do we train a supervised neural network? Well, we present it with two datasets and ask it to learn how to transform one into the other. Think back to our streetlight problem. Can you identify two datasets? Which one do we always know? Which one do we want to know?

We do indeed have two datasets. On the one hand, we have six streetlight states. On the other hand, we have 6 observations of whether people walked or not. These are our two datasets.

So, we can train our neural network to convert from the dataset we **know** to the dataset that we **want to know**. In this particular "real world example", we know the state of the streetlight at any given time, and we want to know whether it is safe to cross the street.







What We Know	What We Want to Know
	STOP
	WALK
	STOP
	WALK
	WALK
	STOP

So, in order to prepare this data for our neural network, we need to first split it into these two groups (what we know and what we want to know). Note that we could attempt to go backwards if we swapped which dataset was in which group. For some problems, this works.

Matrices and the Matrix Relationship

Translating your streetlight into math.

Math doesn't understand streetlights. As mentioned in the previous section, we want to teach our neural network to translate a streetlight pattern into the correct stop/walk pattern. The operative word here is **pattern**. What we really want to do is mimick the pattern of our streetlight in the form of numbers. Let me show you what I mean.

Streetlights	Streetlight Pattern
	1 0 1
	0 1 1
	0 0 1
	1 1 1
	0 1 1
	1 0 1

Notice in the matrix on the right that we have mimicked the pattern from our streetlights in the form of 1s and 0s. Notice that each of the lights gets a column (3 columns total since there are three lights). Notice also that there are 6 rows representing the 6 different streetlights that we observed.

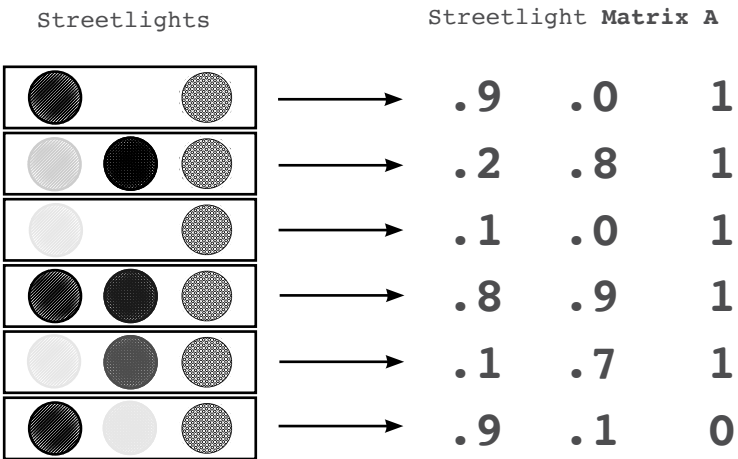
This structure of 1s and 0s is called a **matrix**. Furthermore, this relationship between the rows and columns is very common in matrices, especially matrices of data (like our streetlights).

In data matrices, it is convention to give each *recorded example* a single **row**. It is also convention to give each *thing being recorded* a single **column**. This makes it easy to read.

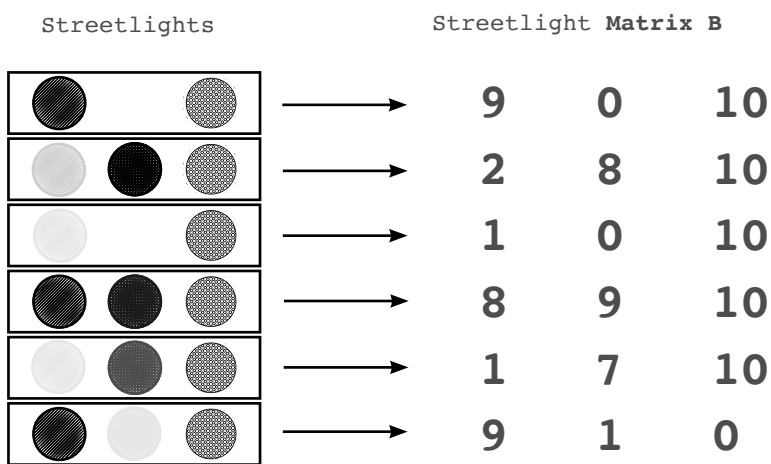
So, a column contains every state we recorded a thing in. In this case, a column contains every on/off state we recorded of a particular light. Each row contains the simultaneous state of every light at a particular moment in time. Again, this is common.

Good data matrices perfectly mimic the outside world.

Our data matrix doesn't have to be all 1s and 0s. What if the streetlights were on "dimmers" and they turned on and off at varying degrees of intensity. Perhaps our streetlight matrix would look more like this:



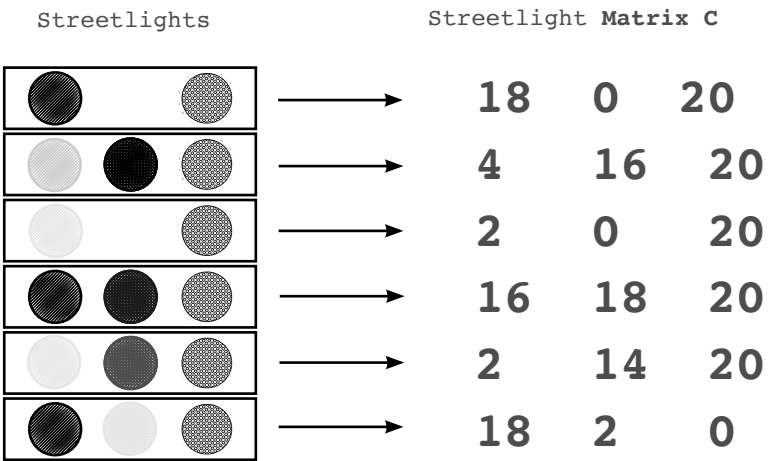
Matrix A above is a perfectly valid matrix. It is mimicking the patterns that exist in the real world (streetlight) so that we can ask our computer to interpret them. Would the following matrix still be valid?



In fact, this matrix (B) is still valid. It adequately captures the relationships between various training examples (rows) and lights (columns). Note that "Matrix A" * 10 == "Matrix B" (A * 10 == B). This actually means that these matrices are **scalar multiples** of each other.

Matrix A and B both contain the same underlying pattern.

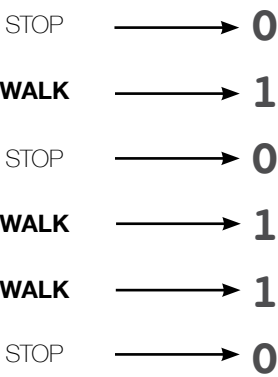
The important takeaway here is that there are an **infinite** number of matrices that perfectly reflect the streetlight patterns in our dataset. Even the one below is still perfect.



It's important to recognize that "the underlying pattern" is not the same as "the matrix". It's a property of the matrix. In fact, it's a property of all three of these matrices (A, B, and C). The pattern is what each of these matrices is *expressing*. The pattern also existed in the streetlights. This *input data pattern* is what we want our neural network to learn to transform into the *output data pattern*.

However, in order to learn the *output data pattern*, we also need to capture the pattern in the form of a matrix. Let's do that below.

Note that we could reverse the 1s and 0s here and the output matrix would still be capturing the underlying STOP/WALK pattern that's present in our data. We know this because regardless of whether we assign a 1 to WALK or to STOP, we can still decode the 1s and 0s into the underlying STOP/WALK pattern. We call this resulting matrix a **lossless representation** because we can perfectly convert back and forth between our stop/walk notes and the matrix.



Creating a Matrix or Two in Python

Importing our matrices into Python

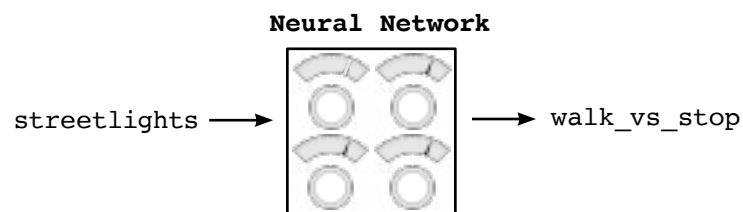
So, we've converted our streetlight pattern into a matrix (one with just 1s and 0s). Now, we want to create that matrix (and more importantly, it's underlying pattern) in Python so that our neural network can read it. Python has a special library built just for handling matrices called **NumPy**. Let's see it in action:

```
import numpy as np
streetlights = np.array( [ [ 1, 0, 1 ],
                           [ 0, 1, 1 ],
                           [ 0, 0, 1 ],
                           [ 1, 1, 1 ],
                           [ 0, 1, 1 ],
                           [ 1, 0, 1 ] ] )
```

If you're a regular Python user, something should be very striking from this code. A matrix is really just a list of lists! It's an array of arrays! What is NumPy? NumPy is really just a fancy wrapper for an array of arrays that gives us special, matrix-oriented functions. Let's create a NumPy matrix for our output data too:

```
walk_vs_stop = np.array( [ [ 0 ],
                           [ 1 ],
                           [ 0 ],
                           [ 1 ],
                           [ 1 ],
                           [ 0 ] ] )
```

So, what will we want our neural network to do? Well, we will want it to take our `streetlights` matrix and learn to transform it into our `walk_vs_stop` matrix. More importantly, we will want our neural network to take **any matrix containing the same underlying pattern as `streetlights`** and transform it into a matrix that contains the underlying pattern of `walk_vs_stop`. More on that later. Let's start by trying to transform `streetlights` into `walk_vs_stop` using a neural network.



Building Our Neural Network

Ok, so we've been learning about neural networks for several chapters now. We've got a new dataset, and we're going to create a neural network to solve it. Below, I've written out some example code to learn the first streetlight pattern. This should look very familiar.

```
import numpy as np
weights = np.array([0.5,0.48,-0.7])
alpha = 0.1

streetlights = np.array( [ [ 1, 0, 1 ],
                           [ 0, 1, 1 ],
                           [ 0, 0, 1 ],
                           [ 1, 1, 1 ],
                           [ 0, 1, 1 ],
                           [ 1, 0, 1 ] ] )

walk_vs_stop = np.array( [ 0, 1, 0, 1, 1, 0 ] )

input = streetlights[0] # [1,0,1]
goal_prediction = walk_vs_stop[0] # equals 0... i.e. "stop"

for iteration in range(20):
    prediction = input.dot(weights)
    error = (goal_prediction - prediction) ** 2
    delta = prediction - goal_prediction
    weights = weights - (alpha * (input * delta))

    print("Error:" + str(error) + " Prediction:" + str(prediction))
```

Perhaps this code example will bring back several nuances we learned in Chapter 3. First, the use of the function "dot" was a way to perform a dot product (weighted sum) between two vectors. However, not included in Chapter 3 is the way that NumPy matrices can perform elementwise addition and multiplication.

```
import numpy as np

a = np.array([0,1,2,1])
b = np.array([2,2,2,3])

print(a*b) #elementwise multiplication
print(a+b) #elementwise addition
print(a * 0.5) # vector-scalar multiplication
print(a + 0.5) # vector-scalar addition
```

One could say that NumPy makes these operations very easy. When you put a "+" sign between two vectors, it does what you would expect it to. It adds the two vectors together. Other than these nice NumPy operators and our new dataset, the neural network above is the same as ones we built before.

Learning the whole dataset!

So, in the last few pages we've only been learning one streetlight. Don't we want to learn them all?

So far in this book, we've trained neural networks that learned how to model a single training example (input -> goal_pred pair). However, now we're trying to build a neural network that tells us "whether or not it is safe to cross the street". We need it to know more than one streetlight! How do we do this? We train it on all the streetlights at once!

```
import numpy as np

weights = np.array([0.5, 0.48, -0.7])
alpha = 0.1

streetlights = np.array( [[ 1, 0, 1 ],
                          [ 0, 1, 1 ],
                          [ 0, 0, 1 ],
                          [ 1, 1, 1 ],
                          [ 0, 1, 1 ],
                          [ 1, 0, 1 ] ] )

walk_vs_stop = np.array( [ 0, 1, 0, 1, 1, 0 ] )

input = streetlights[0] # [1,0,1]
goal_prediction = walk_vs_stop[0] # equals 0... i.e. "stop"

for iteration in range(40):
    error_for_all_lights = 0
    for row_index in range(len(walk_vs_stop)):
        input = streetlights[row_index]
        goal_prediction = walk_vs_stop[row_index]

        prediction = input.dot(weights)

        error = (goal_prediction - prediction) ** 2
        error_for_all_lights += error

        delta = prediction - goal_prediction
        weights = weights - (alpha * (input * delta))
        print("Prediction:" + str(prediction))
    print("Error:" + str(error_for_all_lights) + "\n")
```

```
Error:2.6561231104
Error:0.962870177672
...
Error:0.000614343567483
Error:0.000533736773285
```

Full / Batch / Stochastic Gradient Descent

Stochastic Gradient Descent: Updating weights one example at a time.

As it turns out, this idea of learning "one example at a time" is a variant on Gradient Descent called **Stochastic** Gradient Descent, and it is one of the handful of methods that can be used for learning an entire dataset.

How does Stochastic Gradient Descent work? As exemplified on the previous page, it simply performs a prediction and weight update for each training example separately. In other words, it takes the first streetlight, tries to predict it, calculates the weight_delta, and updates the weights. Then it moves onto the second streetlight, etc. It iterates through the entire dataset many times until it can find a weight configuration that works well for all of the training examples.

(Full) Gradient Descent: Updating weights one dataset at a time.

Another method for learning an entire dataset is simply called Gradient Descent (or "Average/Full Gradient Descent" if you like). Instead of updating the weights once for each training example, the network simply calculates the average weight_delta over the entire dataset, only actually changing the weights each time it computes a full average.

Batch Gradient Descent: Updating weights after "n" examples.

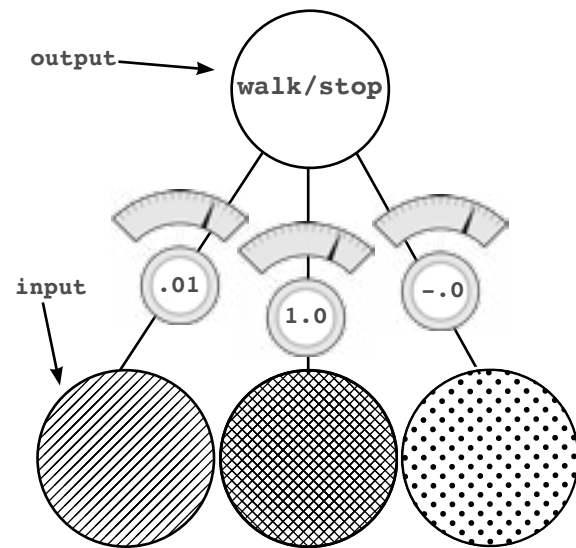
This will be covered in more detail later, but there is also a third configuration that sort of "splits the difference" between Stochastic Gradient Descent and Full Gradient Descent. Instead of updating the weights after just one or after the entire dataset of examples, you choose a "batch size" (typically between 8 and 256) after which the weights are updated.

We will discuss this more later in the book, but for now, simply recognize that on the previous page we created a neural network that can learn our entire "Streetlights" dataset by training on each example, one at a time.

Neural Networks Learn Correlation

What did our last neural network learn?

We just got done training a single-layer neural network to take a streetlight pattern and identify whether or not it was safe to cross the street. Let's take on the neural network's perspective for a moment. The neural network doesn't know that it was processing streetlight data. All it was trying to do was identify which input (out of the 3 possible) correlated with the output. You can see that it correctly identified the middle light by analyzing the final weight positions of the network.



Notice that the middle weight is very near 1 while the far left and right weights are very near 0. At a high level, all the iterative, complex processes for learning we identified actually accomplished something rather simple. The network *identified correlation* between the middle input and output. The correlation is located wherever the weights were set to high numbers. Inversely, *randomness* with respect to the output was found at the far left and far right weights (where the weight values are very near 0).

How did it identify correlation? Well, in the process of **Gradient Descent**, each training example either asserts *up pressure* or *down pressure* on our weights. On average, there was more *up pressure* for our middle weight and more *down pressure* for our other weights. Where does the pressure come from? Why is it different for different weights?

Up and Down Pressure

It comes from our data.

Each neuron is individually trying to correctly predict the output given the input. For the most part, each neuron ignores all the other neurons when attempting to do so. The only *cross communication* occurs in that all three weights must share the same error measure. Our *weight update* is nothing more than taking this *shared* error measure and multiplying it by each *respective* input. Why do we do this? Well, a key part of why neural networks learn is by **error attribution**, which means that given a shared error, the network needs to figure out which weights contributed (so they can be adjusted) and which weights did NOT contribute (so they can be left alone).

Training Data				Weight Pressure			
1	0	1	0	-	0	-	0
0	1	1	→ 1	0	+	+	→ 1
0	0	1	→ 0	0	0	-	→ 0
1	1	1	→ 1	+	+	+	→ 1
0	1	1	→ 1	0	+	+	→ 1
1	0	1	→ 0	-	0	-	→ 0

Consider the first training example. Because the middle input is a 0, then the middle weight is *completely irrelevant* for this prediction. No matter what the weight is, it's going to be multiplied by zero (the input). Thus, any error at that training example (regardless of whether its too high or too low), can only be **attributed** to the far left and right weights.

Consider the pressure of this first training example. If the network should predict 0, and two inputs are 1s, then this is going to cause error which drives the weight values *towards 0*.

The "Weight Pressure" table helps describe the effect that each training example has on each respective weight. + indicates that it has pressure *towards 1* whereas the - indicates that it has pressure *towards 0*. Zeroes (0) indicate that there is *no pressure* because the input datapoint is 0, so that weight won't be changed at all. Notice that the far left weight has 2 negatives and 1 positive, so on average the weight will move towards 0. The middle weight has 3 positives, so on average the weight will move towards 1.

Up and Down Pressure (cont.)

Training Data				Weight Pressure			
1	0	1	0	-	0	-	0
0	1	1	→ 1	0	+	+	→ 1
0	0	1	→ 0	0	0	-	→ 0
1	1	1	→ 1	+	+	+	→ 1
0	1	1	→ 1	0	+	+	→ 1
1	0	1	0	-	0	-	0

So, each individual weight is attempting to compensate for error. In our first training example, we see *discorrelation* between the far right and left inputs and our desired output. This causes those weights to experience *down pressure*. This same phenomenon occurs throughout all 6 training examples, rewarding correlation with pressure *towards 1* and penalizing de-correlation with pressure *towards 0*. On average, this causes our network to find the correlation that is present between our middle weight and the output to be the dominant predictive force (heaviest weight in the weighted average of our input) making our network quite accurate.

Bottom Line

Our prediction is a weighted sum of our inputs. Our learning algorithm rewards inputs that correlate with our output with *upward pressure* (towards 1) on their weight while penalizing inputs with no-correlation with *downward pressure*. So, the weighted sum of our inputs will find perfect correlation between our input and our output by weighting de-correlated inputs to 0.

Now, the mathematician in you might be cringing a little bit. "upward pressure" and "downward pressure" are hardly precise mathematical expressions, and they have plenty of edge cases where this logic doesn't hold (which we'll address in a second). However, we will later find that this is an *extremely* valuable approximation, allowing us to temporarily overlook all of the complexity of Gradient Descent and just remember that *learning rewards correlation* with larger weights (or more generally, *learning finds correlation between our two datasets*).

Edge Case: Overfitting

Sometimes correlation happens accidentally...

Consider again the first example in the training data. What if our far left weight was 0.5 and our far right weight was -0.5. Their prediction would equal 0! The network would predict perfectly! However, it hasn't remotely learned how to safely predict streetlights (i.e. those weights would fail in the real world). This phenomenon is known as **overfitting**.

Deep Learning's Greatest Weakness: Overfitting

Error is shared between all of our weights. If a particular configuration of weights *accidentally* creates perfect correlation between our prediction and the output dataset (such that error == 0) without actually giving the heaviest weight to the best inputs, **the neural network will stop learning**.

In fact, if it wasn't for our other training examples, this fatal flaw would cripple our neural network. What do the other training examples do? Well, let's take a look at the second training example. It would bump the far right weight *upward* while not changing the far left weight. This throws off the equilibrium that stopped the learning in our first example. So, as long as we don't train *exclusively on the first example*, the rest of the training examples will help the network avoid getting stuck in these edge case configurations that exist for any one training example.

This is **very** important. Neural networks are so flexible that they can find many, many different weight configurations that will correctly predict for a subset of your training data. In fact, if we trained our neural network on the first 2 training examples, it would likely stop learning at a point where it did NOT work well for our other training examples. In essence, it *memorized* the two training examples instead of actually finding the *correlation* that will *generalize* to any possible streetlight configuration.

If we only train on two streetlights and the network just finds these edge case configurations, it could FAIL to tell us whether it is safe to cross the street when it sees a streetlight that wasn't in our training data!

The greatest challenge you will face with deep learning is convincing your neural network to *generalize* instead of just *memorize*. We will see this again.

Edge Case: Conflicting Pressure

Sometimes correlation fights itself.

Consider the far right column in the "Weight Pressure" table below. What do you see? This column seems to have an equal number of *upward* and *downward* pressure moments. However, we have seen that the network correctly pushes this (far right) weight down to 0 which means that the *downward* pressure moments must be larger than the *upward* ones. How does this work?

Training Data				Weight Pressure			
1	0	1	0	-	0	-	0
0	1	1	→ 1	0	+	+	→ 1
0	0	1	→ 0	0	0	-	→ 0
1	1	1	→ 1	+	+	+	→ 1
0	1	1	→ 1	0	+	+	→ 1
1	0	1	0	-	0	-	0

The left and middle weights have enough signal to converge on their own. The left weight falls to 0 and the middle weight moves towards 1. As the middle weight moves higher and higher, the error for positive examples continues to decrease. However, as they approach their optimal positions, the de-correlation on the far right weight becomes more apparent. Let's consider the extreme example of this, where the left and middle weights are perfectly set to 0 and 1 respectively. What happens to our network? Well, if our right weight is above 0, then our network predicts too *high* and if our right weight is beneath 0, our network predicts too *low*.

In short, as other neurons learn, they absorb some of the *error*; they absorb some part of the *correlation*. They cause the network to predict with *moderate* relative power which reduces the error. The other weights then only try to adjust their weights to correctly predict what's left! In this case, because the middle weight has consistent signal to absorb *all* of the correlation (because of the 1:1 relationship between the middle input and the output), the error when we want to predict 1 becomes very *small*, but the error to predict 0 becomes large, pushing our middle weight downward.

Edge Case: Conflicting Pressure (cont.)

It doesn't always work out like this.

In some ways, we kind of got lucky. If our middle node hadn't been so perfectly correlated, our network might have struggled to silence our far right weight. In fact, later we will learn about **Regularization** which forces weights with conflicting pressure to move towards 0.

As a preview, *regularization* is advantageous because if a weight has equal pressure *upward* and *downward*, then it isn't really good for anything. It's not helping either direction. In essence, regularization aims to say "only weights with really strong correlation can stay on, everything else should be silenced because it's contributing noise". It's sort of like natural selection, and as a side effect it would cause our neural network to train *faster* (fewer iterations) because our far right weight has this "both positive and negative" pressure problem. In this case, because our far right node isn't *definitively* *correlative*, the network would immediately start driving it towards 0. Without regularization (like we trained it before), we won't end up learning that the far right input is useless until after the left and middle start to figure their patterns out. More on this later.

So, if networks look for correlation between an input column of data and our output column, what would our neural network do with this dataset?

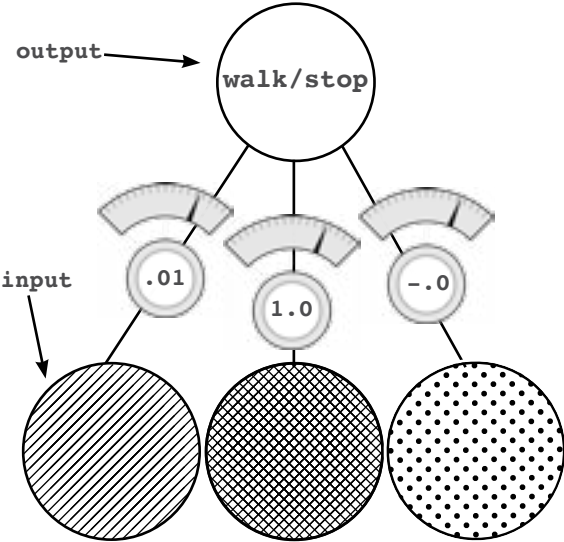
Training Data				Weight Pressure			
1	0	1	1	+	0	+	1
0	1	1	→ 1	0	+	+	→ 1
0	0	1	→ 0	0	0	-	→ 0
1	1	1	→ 0	-	-	-	→ 0

There is no correlation between any input column and the output column. Every weight has an equal amount of upward pressure as it does downward pressure. *This dataset is a real problem for our neural network.* Previously, we could solve for input datapoints that had *both* upward and downward pressure because other neurons would start solving for either the positive or negative predictions, drawing our balanced neuron to favor up or down. However, in this case, *all of the inputs* are equally balanced between *positive* and *negative* pressure. What do we do?

Learning Indirect Correlation

If your data doesn't have correlation, let's create intermediate data that does!

Previously, I have described neural networks as an instrument that searches for correlation between input and output *datasets*. I should refine this just a touch. In reality, neural networks actually search for correlation between their input and output *layers*. We set the values of our input *layer* to be individual rows of our input data, and we try to train the network so that our output *layer* equals our output dataset. Funny enough, the neural network actually doesn't "know" about data. It just searches for correlation between the input and output layers.



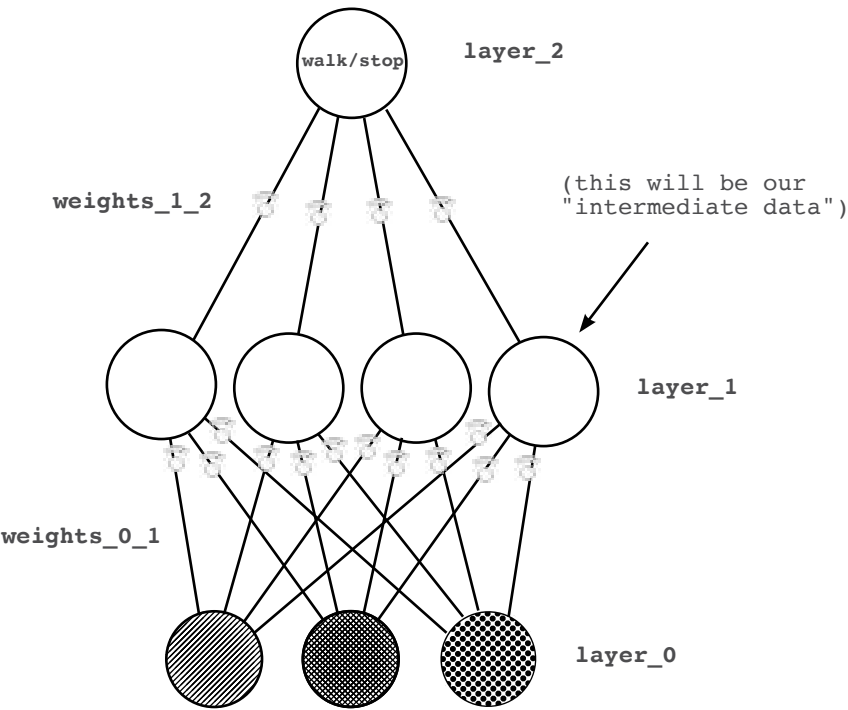
Unfortunately, we just encountered a new streetlights dataset where there isn't *any* correlation between our input and output. The solution is simple. Let's use *two* of these networks. The first one will create an *intermediate dataset* that has *limited* correlation with our output. The second will then use that *limited* correlation to correctly predict our output!

Since our input dataset doesn't correlate with our output dataset, we're going to use our input dataset to create an *intermediate dataset* that *DOES* have correlation with our output. It's kind of like cheating!

Creating Our Own Correlation

If your data doesn't have correlation, let's create intermediate data that does!

Below, you'll see a picture of our new neural network. Notice that we basically just stacked two neural networks on top of each other. The middle layer of nodes (*layer_1*) represents our *intermediate dataset*. Our goal is to train this network so that even though there's no correlation between our input dataset and output dataset (*layer_0* and *layer_2*) that our *layer_1* dataset that we create *using layer_0* will have correlation with *layer_2*.

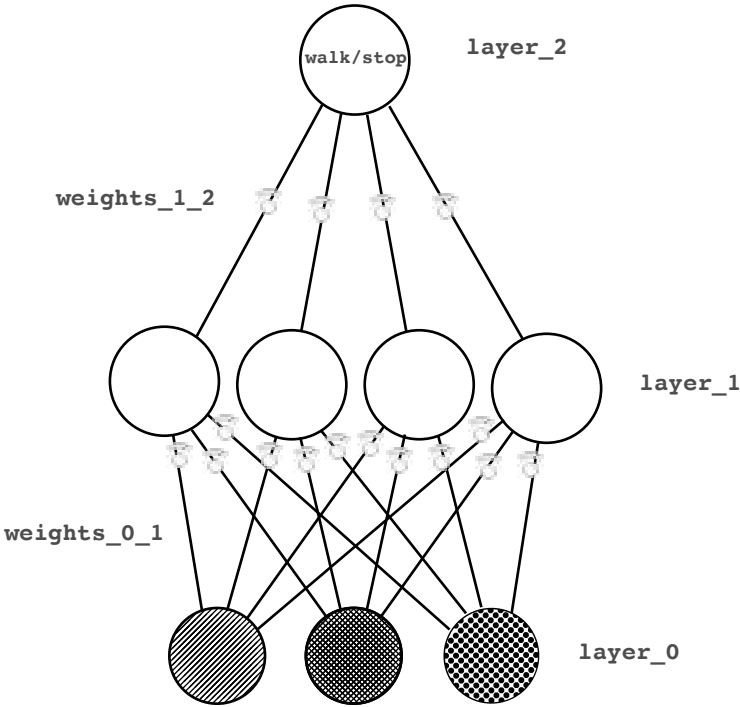


Things to notice: This network is still just a function! It still just has a bunch of weights that are collected together in a particular way. Furthermore, **Gradient Descent** still works because we can calculate how much each weight contributes to the error and adjust it to reduce the error to 0. And that's exactly what we're going to do.

Stacking Neural Networks - A Review

In Chapter 3, we briefly mentioned stacked neural networks. Let's review.

So, when you look at the architecture below, the prediction occurs exactly as you might expect when I say "stack neural networks". The *output* of the first "lower" network (layer_0 to layer_1) is the *input* to the second "upper" neural network (layer_1 to layer_2). The prediction for each of these networks is identical to what we saw before.



So, as we start to think about how this neural network learns, we actually already know a great deal. If we ignored the lower weights and just considered their output to be our training set, then the top half of the neural network (layer_1 to layer_2) is just like the networks we trained in the last chapter. We can use all the same learning logic to help them learn. In fact, this is the case.

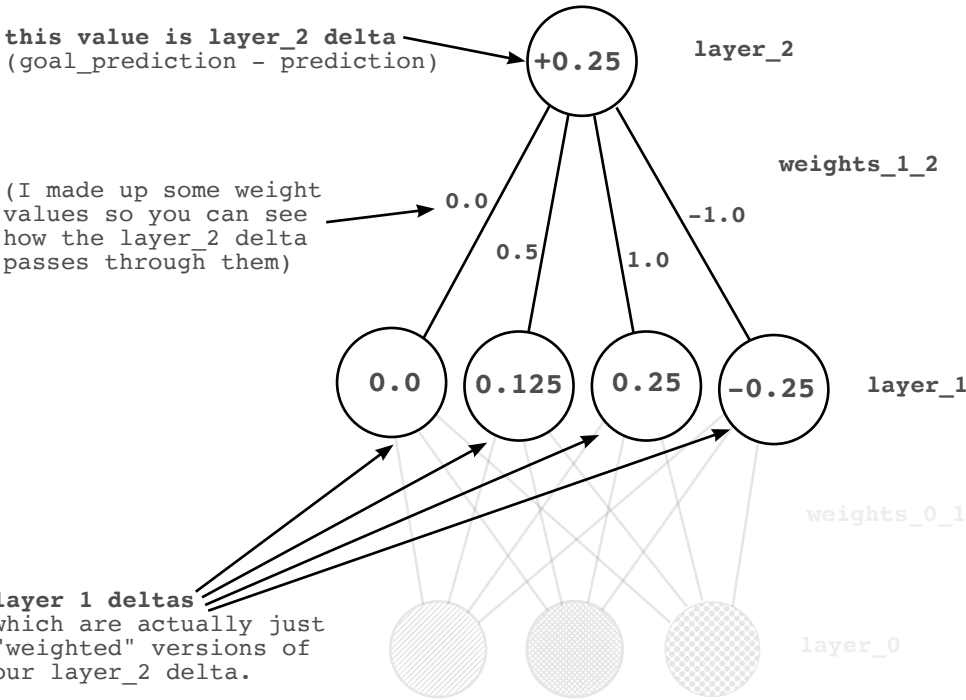
So, the part that we *don't yet understand* is how to update the weights between layer_0 and layer_1. What do they use as their error measure? If you remember from the last chapter, our *cached/normalized error measure* was called *delta*. In our case, we want to figure out how to know the *delta* values at layer_1 so that they can help layer_2 make accurate predictions.

Backpropagation: Long Distance Error Attribution

The "weighted average error"

What is the prediction from layer_1 to layer_2? It's just a weighted average of the values at layer_1. So, if layer_2 is too high by "x" amount, how do we know which values at layer_1 contributed to the error? Well, the ones with *higher weights* (weights_1_2) contributed more! The ones with *lower weights* from layer_1 to layer_2 contributed less! Consider the extreme. Let's say that the far *left* weight from layer_1 to layer_2 was zero. How much did that node at layer_1 cause the network's error? ZERO!

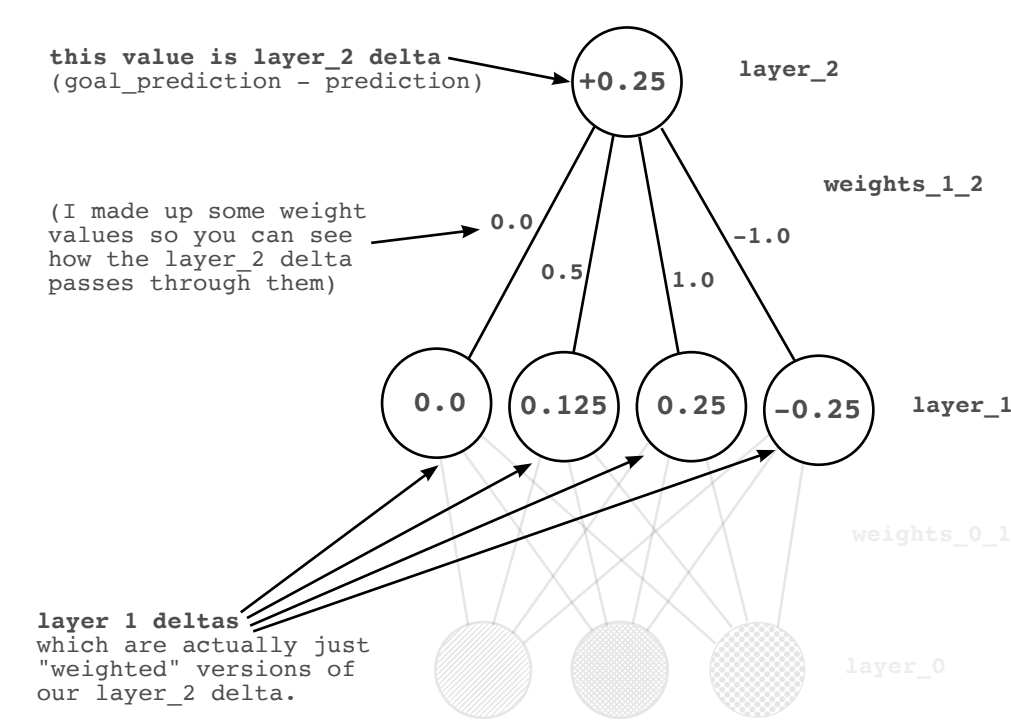
It's so simple it's almost hilarious. Our weights from layer_1 to layer_2 exactly describe how much each layer_1 neuron contributes to the layer_2 prediction. This means that those weights *ALSO exactly describe* how much each layer_1 neuron contributes to the layer_2 **error**! So, how do we use the delta at layer_2 to figure out the delta at layer_1? We just multiply it by each of the respective weights for layer_1!!! It's like our prediction logic in reverse! This process of "moving delta signal around" is called **backpropagation**.



Backpropagation: Why does this work?

The "weighted average delta"

In our neural network from the previous chapter, our *delta* variable told us "the **direction** and **amount** we want the value of this neuron to change next time". All backpropagation lets us do is say "Hey, if you want this neuron to be X amount higher, then each of these previous 4 neurons need to be X*weights_1_2 amount higher/lower, because these weights were *amplifying* the prediction by weights_1_2 times". When used in *reverse*, our weights_1_2 matrix *amplifies the error* by the appropriate amount. It *amplifies the error* so that we know how much each layer_1 node should move up or down. Once we know this, we can just update each weight matrix just like we did before. For each weight, multiply its output *delta* by its input *value*... and adjust our weight by that much (or we can scale it with *alpha*).



Linear vs Non-Linear

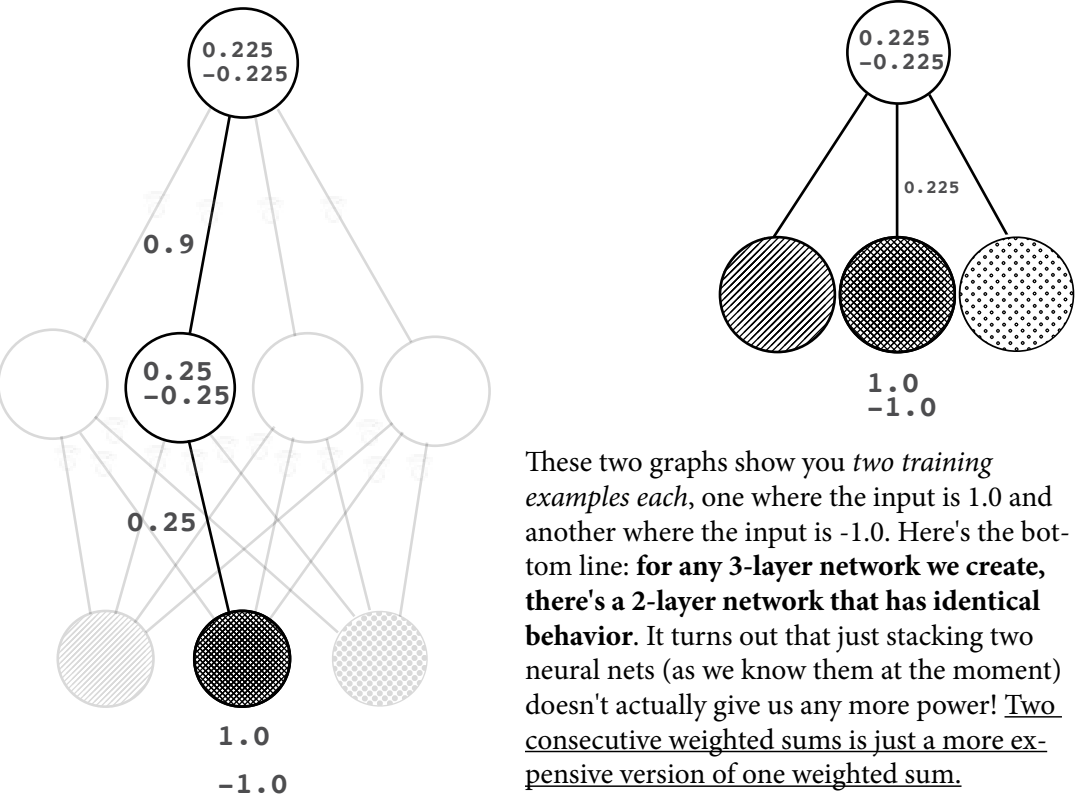
This is probably the hardest concept in the book. Let's take it slow.

I'm going to show you a phenomenon. As it turns out, we need one more "piece" to make this neural network train. We're going to take it from two perspectives. The first is going to show you why *the neural network can't train without it*. In other words, first I'm going to show you why our neural network is currently broken. Then, once we add this piece, I'm going to show you what it does to fix this problem. For now, check out this simple algebra:

$$\begin{matrix} 1 & * & 10 & * & 2 & = & 100 \\ 5 & * & 20 & = & 100 \end{matrix}$$

$$\begin{matrix} 1 & * & 0.25 & * & 0.9 & = & 0.225 \\ 1 & * & 0.225 & = & 0.225 \end{matrix}$$

Here's the takeaway, for any *two multiplications* that I do, I can actually accomplish the *same thing* using a single multiplication. As it turns out, this is bad. Check out the following:



These two graphs show you *two training examples each*, one where the input is 1.0 and another where the input is -1.0. Here's the bottom line: **for any 3-layer network we create, there's a 2-layer network that has identical behavior**. It turns out that just stacking two neural nets (as we know them at the moment) doesn't actually give us any more power! Two consecutive weighted sums is just a more expensive version of one weighted sum.

Why The Neural Network Still Doesn't Work

If we trained the 3 layer network as it is now, it would NOT converge.

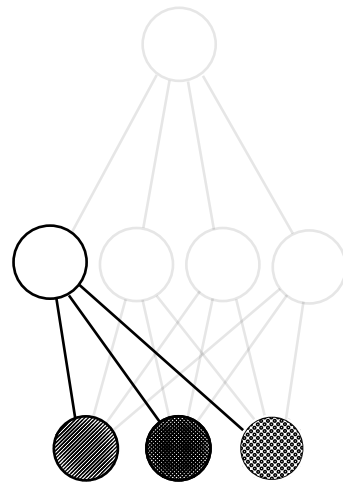
Problem: For *any* two consecutive *weighted sums* of the input, there exists a *single weighted sum* with exactly identical behavior. Aka, anything that our 3 layer network can do, our 2 layer network can also do.

Let's talk about the *middle layer* (layer_1) that we have at present (before we fix it). Right now, each node (out of the 4 we have), has a weight coming to it from each of the inputs. Let's think about this from a *correlation* standpoint. Each node in our *middle layer* subscribes to a *certain amount of correlation* with each input node. If the weight from an input to the middle layer is a 1.0, then it subscribes to *exactly* 100% of that node's movement. If that node goes up by 0.3, our middle node will follow. If the weight connecting two nodes is 0.5, it subscribes to *exactly* 50% of that node's movement.

The only way that our middle node can escape the correlation of one particular *input node* is if it *subscribes to additional correlation from another input node*. So, you can see, there's *nothing new* being contributed to this neural network. Each of our hidden nodes simply subscribe to a little bit of correlation from our input nodes. Our middle nodes don't actually get to *add anything to the conversation*. They don't get to have *correlation of their own*. They're just more or less correlated to various *input nodes*. However, since we KNOW that in our new dataset there is NO correlation between ANY of our inputs and our output, then how could our middle layer help us at all?!

It just gets to mix up a bunch of correlation

that was already *useless*! What we really need is for our *middle layer* to be able to *selectively* correlate with our input. We want it to *sometimes* correlate with an input, and *sometimes not correlate*. That gives it correlation of its own! This gives our middle layer the *opportunity* to not just "always be X% correlated to one input and Y% correlated to another input". Instead, it can be "X% correlated to one input only when it wants to be, but other times not be correlated at all". This is "conditional correlation" or "sometimes correlation".



The Secret to "Sometimes Correlation"

We're going to simply turn our node "off" when the value would be below 0.

This might seem too simple to work, but consider this: if the node's value dropped below 0, normally the node would still have *just as much correlation* to the input as it always did! It would just happen to be negative in value. However, if we *turn off the node* (setting it to 0) when it would be negative, then it has *ZERO CORRELATION* to *ANY INPUTS* whenever it's negative.

What does this mean? It means that our node can now selectively pick and choose when it wants to be correlated to something. This allows it to say something like "make me perfectly correlated to the left input but *ONLY* when the right input is turned OFF". How would it do this? Well, if the weight from the left input is a 1.0, and the weight from the right input is a HUGE NEGATIVE NUMBER, then turning on both the left and right inputs would cause the node to just be 0 all the time. However, if just the left node was on, the node would take on the value of the left node.

This wasn't possible before! Before our middle node was either ALWAYS correlated to an input or ALWAYS not correlated. Now it can be conditional. Now it can speak for itself!

Solution: By turning any middle node off whenever it would be negative, we allow the network to *sometimes* subscribe to correlation from various inputs. This is *impossible* for 2-layer neural networks, thus adding power to 3-layer nets.

The fancy term for this "if the node would be negative then set it to 0" logic is called a **nonlinearity**. This is because without this tweak, our neural network is *linear*. Without this technique, our output layer only gets to pick from the same *correlation* that it had in the 2-layer network. It's still just subscribing to pieces of the *input layer*, which means that it can't solve our new streetlights dataset.

There are **many kinds of nonlinearities**. However, the one we discussed above is, in many cases, the best one to use. It's also the simplest. (It's called "relu".)

For what it's worth, most other books/courses simply say, "consecutive matrix multiplication is still just a linear transformation". I find this to be very unintuitive. Furthermore, it makes it harder to understand what nonlinearities actually *contribute* and why you choose one over the other (which we'll get to later). It just says "without the nonlinearity two matrix multiplications might as well be 1". So, this page's explanation, while not the most concise answer, is an intuitive explanation of why we need nonlinearities.

A Quick Break

That last part probably felt a little abstract, and that's totally ok.

So, here's the deal. In previous chapters we were working with very simple algebra. This meant that everything was ultimately grounded in fundamentally simple tools. This chapter has started building on the premises we learned previously. In other words, previously we learned lessons like:

We can compute the relationship between our *error* and any one of our *weights* so that we know how changing the weight changes the error. We can then use this to reduce our error down to 0.

That was a **massive lesson!** However, now we're moving past it. Since we already worked through why that works, we can just trust it. We take the statement at face value. The next big lesson came at the beginning of this chapter:

Adjusting our weights to reduce our error over a *series of training examples* ultimately just searches for correlation between our *input* and our *output* layers. If no correlation exists, then error will never reach 0.

This is an even **bigger lesson!** Why? Well, it largely means that we can put the previous lesson out of our minds for now. We don't actually need it. Now we're focused on *correlation*. The takeaway for you is that you can't constantly think about *everything all at once*. You take each one of these lessons and you let yourself trust it. When it's a more *concise* summarization (a higher abstraction) of more granular lessons, we can set aside the granular and only focus on understanding the higher summarizations.

This is akin to a professional swimmer, biker, or really any other skill that requires a *combined fluid knowledge* of a bunch of really small lessons. A baseball player who *swings a bat* actually learned thousands of little lessons to ultimately culminate in a great bat swing. However, he doesn't think of *all of them* when he goes to the plate! He just lets it be fluid, subconscious even. It is the same way for studying these math concepts.

Neural networks look for correlation between input and output... and you no longer have to worry about *how that happens*. We just know that it does. Now we're building on that idea! Let yourself relax and trust the things you've already learned.

Our First "Deep" Neural Network

How to Make the Prediction

In the code below, we initialize our weights and make a forward propagation. **New is bold.**

```
import numpy as np
np.random.seed(1)

def relu(x):
    return (x > 0) * x

alpha = 0.2
hidden_size = 4

streetlights = np.array( [[ 1, 0, 1 ],
                          [ 0, 1, 1 ],
                          [ 0, 0, 1 ],
                          [ 1, 1, 1 ] ] )

walk_vs_stop = np.array([[ 1, 1, 0, 0 ]]).T

weights_0_1 = 2*np.random.random((3,hidden_size)) - 1
weights_1_2 = 2*np.random.random((hidden_size,1)) - 1

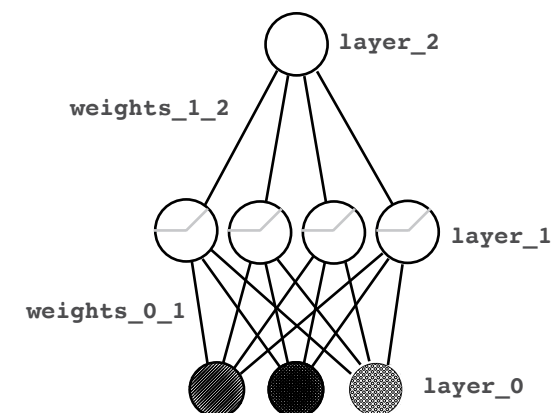
layer_0 = streetlights[0]
layer_1 = relu(np.dot(layer_0,weights_0_1))
layer_2 = np.dot(layer_1,weights_1_2)
```

this function sets all negative numbers to 0

2 sets of weights now to connect our 3 layers (randomly initialized)

the output of layer_1 is sent through "relu" where negative values become 0. This is then the input for the next layer, layer_2.

Take each piece and follow along with the picture on the bottom right. Input data comes into layer_0. Via the "dot" function, the signal travels up the weights from layer_0 to layer_1 (performing a weighted sum at each of the 4 layer_1 nodes). These weighted sums at layer_1 are then passed through the "relu" function, which converts all negative numbers to zero. We then perform a final weighted sum into the final node, layer_2.



Backpropagation in Code

How we can learn the amount that each weight contributes to the final error.

At the end of the previous chapter, I made an assertion that it would be very important to memorize the 2-layer neural network code so that you could quickly and easily recall it when I reference the more advanced concepts. This is when that memorization matters! We're about to look at the new learning code and it is absolutely essential that you recognize and understand the parts that were addressed in the previous chapters. If you get lost, go back to the last chapter and memorize the code and come back. It'll save your life someday.

```
import numpy as np
np.random.seed(1)

def relu(x):
    return (x > 0) * x # returns x if x > 0
                        # return 0 otherwise

def relu2deriv(output):
    return output>0 # returns 1 for input > 0
                    # return 0 otherwise

alpha = 0.2
hidden_size = 4

weights_0_1 = 2*np.random.random((3,hidden_size)) - 1
weights_1_2 = 2*np.random.random((hidden_size,1)) - 1

for iteration in range(60):
    layer_2_error = 0
    for i in range(len(streetlights)):
        layer_0 = streetlights[i:i+1]
        layer_1 = relu(np.dot(layer_0,weights_0_1))
        layer_2 = np.dot(layer_1,weights_1_2)

        layer_2_error += np.sum((layer_2 - walk_vs_stop[i:i+1]) ** 2)

        layer_2_delta = (walk_vs_stop[i:i+1] - layer_2)
        layer_1_delta=layer_2_delta.dot(weights_1_2.T)*relu2deriv(layer_1)

        weights_1_2 += alpha * layer_1.T.dot(layer_2_delta)
        weights_0_1 += alpha * layer_0.T.dot(layer_1_delta)

    if(iteration % 10 == 9):
        print("Error:" + str(layer_2_error))
```

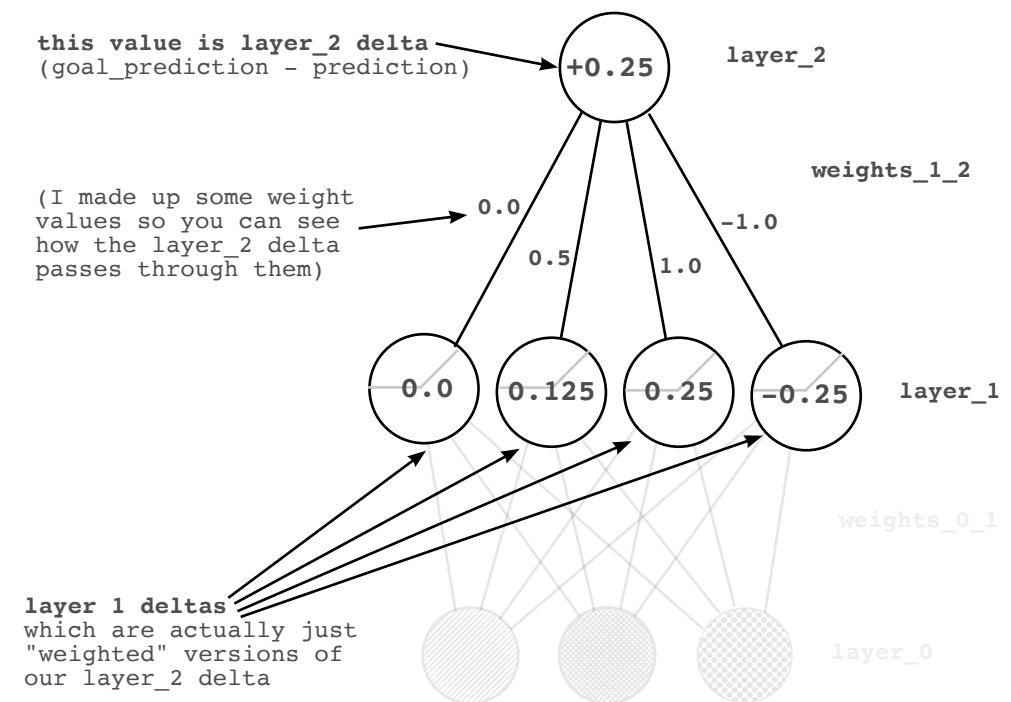
this line computes the delta at layer_1 given the delta at layer_2 by taking the layer_2_delta and multiplying it by its connecting weights_1_2

Believe it or not, the only truly new code is in bold. Everything else is fundamentally the same as in previous pages. The "relu2deriv" function returns 1 when "output" is > 0 and it returns 0 otherwise. This is actually the *slope* of our relu function. It's the *derivative* of our relu function. It serves a very important purpose as we'll see in a moment.

Remember, the goal here is **error attribution**. It's all about figuring out how much each weight contributed to the final error. In our first (2-layer) neural network, we calculated a *delta* variable, which told us how much higher or lower we wanted the output prediction to be. Look at the code here. We compute our `layer_2_delta` in the same way. Nothing new here! (Again, go back to the previous chapter if you've forgotten how that part works.)

So, now that we have how much we want the final prediction to move up or down (delta), we need to figure out how much we want each middle (layer_1) node to move up or down. These are effectively *intermediate predictions*. Once we have the delta at layer_1, we can use all the same processes we used before for calculating a weight update (for each weight, multiply its input value by its output delta and increase the weight value by that much).

So, how do we calculate the deltas for layer_1? Well, first we do the obvious as mentioned on the previous pages, we multiply the output delta by each weight attached to it. This gives us a weighting of how much each weight contributed to that error. There's one more thing we need to factor in. If the relu set the output to a layer_1 node to be 0, then it didn't contribute to the error at all. So, when this was true, we should also set the delta of that node to be zero. Multiplying each layer_1 node by the `relu2deriv` function accomplishes this. `relu2deriv` is either a 1 or a 0 depending on whether the layer_1 value was > 0 or not.

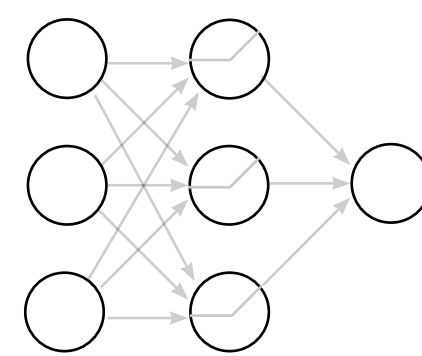


One Iteration of Backpropagation

1 Initialize the Network's Weights and Data

inputs

layer_0



```
import numpy as np
np.random.seed(1)

def relu(x):
    return (x > 0) * x

def relu2deriv(output):
    return output>0

lights = np.array( [[ 1, 0, 1 ],
                    [ 0, 1, 1 ],
                    [ 0, 0, 1 ],
                    [ 1, 1, 1 ] ] )

walk_stop = np.array([[ 1, 1, 0, 0]]).T

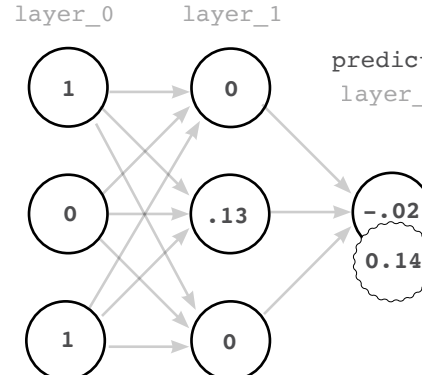
alpha = 0.2
hidden_size = 3

weights_0_1 = 2*np.random.random(\
    (3,hidden_size)) - 1
weights_1_2 = 2*np.random.random(\
    (hidden_size,1)) - 1
```

2 PREDICT & COMPARE: Make a Prediction, Calculate Output Error and Delta

inputs

layer_0



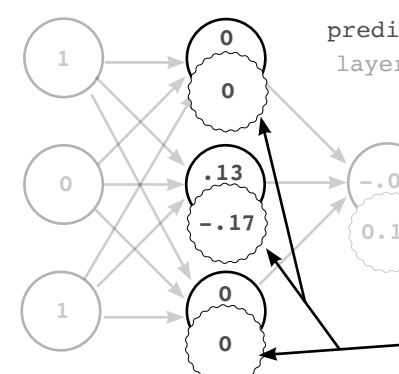
```
layer_0 = lights[0:1]
layer_1 = np.dot(layer_0,weights_0_1)
layer_1 = relu(layer_1)
layer_2 = np.dot(layer_1,weights_1_2)

error = (layer_2-walk_stop[0:1])**2
layer_2_delta=(layer_2-walk_stop[0:1])
```

3 LEARN: Backpropagate From layer_2 to layer_1

inputs

layer_0



```
layer_0 = lights[0:1]
layer_1 = np.dot(layer_0,weights_0_1)
layer_1 = relu(layer_1)
layer_2 = np.dot(layer_1,weights_1_2)

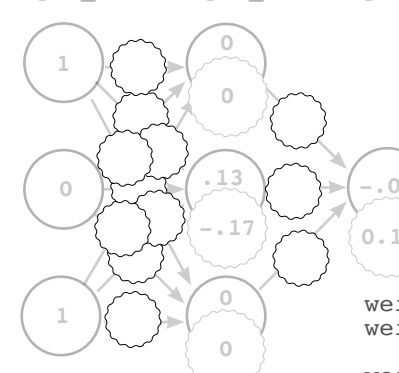
error = (layer_2-walk_stop[0:1])**2
layer_2_delta=(layer_2-walk_stop[0:1])

layer_1_delta=layer_2_delta.dot(weights_1_2.T)
layer_1_delta *= relu2deriv(layer_1)
```

4 LEARN: Generate Weight Deltas and Update Weights

inputs

layer_0



```
layer_0 = lights[0:1]
layer_1 = np.dot(layer_0,weights_0_1)
layer_1 = relu(layer_1)
layer_2 = np.dot(layer_1,weights_1_2)
error = (layer_2-walk_stop[0:1])**2
layer_2_delta=(layer_2-walk_stop[0:1])

layer_1_delta=layer_2_delta.dot(weights_1_2.T)
layer_1_delta *= relu2deriv(layer_1)

weight_delta_1_2 = layer_1.T.dot(layer_2_delta)
weight_delta_0_1 = layer_0.T.dot(layer_1_delta)

weights_1_2 -= alpha * weight_delta_1_2
weights_0_1 -= alpha * weight_delta_0_1
```

As we can see, backpropagation in its entirety is about calculating deltas for intermediate layers so that we can perform Gradient Descent. In order to do so, we simply take the weighted average delta on layer_2 for layer_1 (weighted by the weights in between them). We then turn off (set to 0) nodes that weren't participating in the forward prediction, since they could not have contributed to the error.

Putting it all Together

Here's the self sufficient program you should be able to run (runtime output below)

```
import numpy as np

np.random.seed(1)

def relu(x):
    return (x > 0) * x # returns x if x > 0
                        # return 0 otherwise

def relu2deriv(output):
    return output>0 # returns 1 for input > 0
                    # return 0 otherwise

streetlights = np.array( [[ 1, 0, 1 ],
                           [ 0, 1, 1 ],
                           [ 0, 0, 1 ],
                           [ 1, 1, 1 ] ] )

walk_vs_stop = np.array([[ 1, 1, 0, 0]]).T

alpha = 0.2
hidden_size = 4

weights_0_1 = 2*np.random.random((3,hidden_size)) - 1
weights_1_2 = 2*np.random.random((hidden_size,1)) - 1

for iteration in range(60):
    layer_2_error = 0
    for i in range(len(streetlights)):
        layer_0 = streetlights[i:i+1]
        layer_1 = relu(np.dot(layer_0,weights_0_1))
        layer_2 = np.dot(layer_1,weights_1_2)

        layer_2_error += np.sum((layer_2 - walk_vs_stop[i:i+1]) ** 2)

        layer_2_delta = (layer_2 - walk_vs_stop[i:i+1])
        layer_1_delta=layer_2_delta.dot(weights_1_2.T)*relu2deriv(layer_1)

        weights_1_2 -= alpha * layer_1.T.dot(layer_2_delta)
        weights_0_1 -= alpha * layer_0.T.dot(layer_1_delta)

    if(iteration % 10 == 9):
        print("Error:" + str(layer_2_error))
```

```
Error:0.634231159844
Error:0.358384076763
Error:0.0830183113303
Error:0.0064670549571
Error:0.000329266900075
Error:1.50556226651e-05
```

Why do deep networks matter?

What's the point of creating "intermediate datasets" that have correlation?

Consider the cat picture below. Consider further that we had a dataset of images "with cats" and "without cats" (and we labeled them as such). If we wanted to train a neural network to take our pixel values and predict whether or not there is a cat in the picture, our 2-layer network might have a problem. Just like in our last streetlight dataset, no individual pixel correlates with whether or not there is a cat in the picture. Only different *configurations of pixels* correlate with whether or not there is a cat.

This is the essence of Deep Learning. Deep Learning is all about creating intermediate layers (datasets) wherein each node in an intermediate layer represents the *presence* or *absence* of a different *configuration of inputs*. In this way, for our cat images dataset, no individual pixel has to correlate with whether or not there is a cat in the photo. Instead, our middle layer would attempt to identify different configurations of pixels that may or may not correlate with a cat (such as an ear, or cat eyes, or cat hair). The presence of many "cat like" configurations would then give the final layer the information (correlation) it needs to correctly predict the presence or absence of a cat!

Believe it or not, we can take our 3-layer network and continue to *stack more and more layers*. Some neural networks even have *hundreds* of layers, each neuron playing its part in detecting different configurations of input data. The rest of this book will be dedicated to studying different phenomena within these layers in an effort to explore the full power of deep neural networks.

To that end, I must issue the same challenge as I did in the previous chapter. Memorize the code on the previous page. You will need to be *very familiar* with each of the operations in the code in order for the following chapters to be readable. Do not progress past this page until you can build a 3 layer neural network from memory!!!



IN THIS CHAPTER •

- Correlation Summarization
- Simplified Visualization
- Seeing the Network Predict
- Visualizing Using Letters Instead of Pictures
- Linking our Variables
- The Importance of Visualization Tools

“Numbers have an important story to tell. They rely on you to give them a clear and convincing voice.”

— STEPHEN FEW

It's Time to Simplify

It's impractical to think about everything all the time. Mental tools can help.

At the end of the previous chapter, we finished with a code example that was really quite impressive. Just the neural network itself contained 35 lines of incredibly dense code. Reading through it, it's clear that there's a lot going on, and contained in that code is over 100 pages of concepts that when combined together can predict whether it's safe to cross the street.

I hope that you're continuing to rebuild these examples from memory in each chapter. As these examples get larger, this exercise becomes less about remembering specific letters of code and more about remembering the concepts and then re-building the code based on those concepts. In this chapter, this construction of efficient concepts in your mind is exactly what I want to set aside time to talk about. Even though it's not some architecture or experiment, it's perhaps the most important value I can give you. In this case, I want to show you how I summarize all of these little lessons in an efficient way in my mind so that I can do things like build new architectures, debug experiments, and leverage an architecture on new problems and new datasets.

Let's start by reviewing the concepts we've leareened so far.

In the beginning of this book, we started with very small lessons and then built layers of abstraction on top of them. We began by talking about the ideas behind Machine Learning in general. Then, we progressed to how individual linear neurons learned, followed by horizontal groups of neurons (layers) and then vertical groups (stacks of layers). Along the way, we discussed how learning is actually just reducing error down to 0, and we then leveraged calculus to learn how to change each weight in our network to help move our error in the direction of 0. In the last chapters, we then discussed how neural networks actually search for (and sometimes even create) correlation between the input and output datasets. This last idea allowed us to overlook the previous lessons on how individual neurons behaved because it concisely summarizes the previous lessons. The sum total of the neurons, gradients, stacks of layers, etc. all lead up to a single idea: neural networks find and create correlation. Holding onto this idea of correlation instead of the previous smaller ideas is important to learning Deep Learning. Otherwise, it would be easy to become overwhelmed with the complexity of neural networks. In fact, let's create a name for this idea. Let's call this our Correlation Summarization.

Correlation Summarization

This is the key to sanely moving forward to more advanced neural networks.

Correlation Summarization

Neural networks seek to find direct and indirect correlation between an input layer and an output layer, which are determined by the input and output datasets respectively.

At the 10,000 ft. level, this is what all neural networks do. So, given that a neural network is really just a series of matrices connected by layers, let's zoom in slightly and consider what any particular weight matrix is doing.

Local Correlation Summarization

Any given set of weights optimizes to learn how to correlate its input layer with what the output layer says it should be.

When we only have two layers (input and output), this means that our weight matrix knows what the output layer says it should be based on the output dataset. It looks for correlation between the input and output datasets because they are captured in the input and output layers. However, this becomes more nuanced when we have multiple layers, remember?

Global Correlation Summarization

What an earlier layer says it should be can be determined by taking what a later layer says it should be and multiplying it by the weights in between them. In this way, later layers can tell earlier layers what kind of signal they need to ultimately find correlation with the output. We call this cross-communication, backpropagation.

Thus, when global correlation teaches each layer what it should be, then local correlation can optimize weights locally. In other words, given that some neuron in the final layer says "I need to be a little higher", it then proceeds to tell all the neurons in the layer immediately preceeding it, "Hey previous layer, send me higher signal!". They then tell the neurons preceeding them "Hey! Send us higher signal!". It's like a giant game of telephone! At the end of the game, every layer knows which of its neurons need to be higher and lower, and the Local Correlation Summarization takes over, updating the weights accordingly.

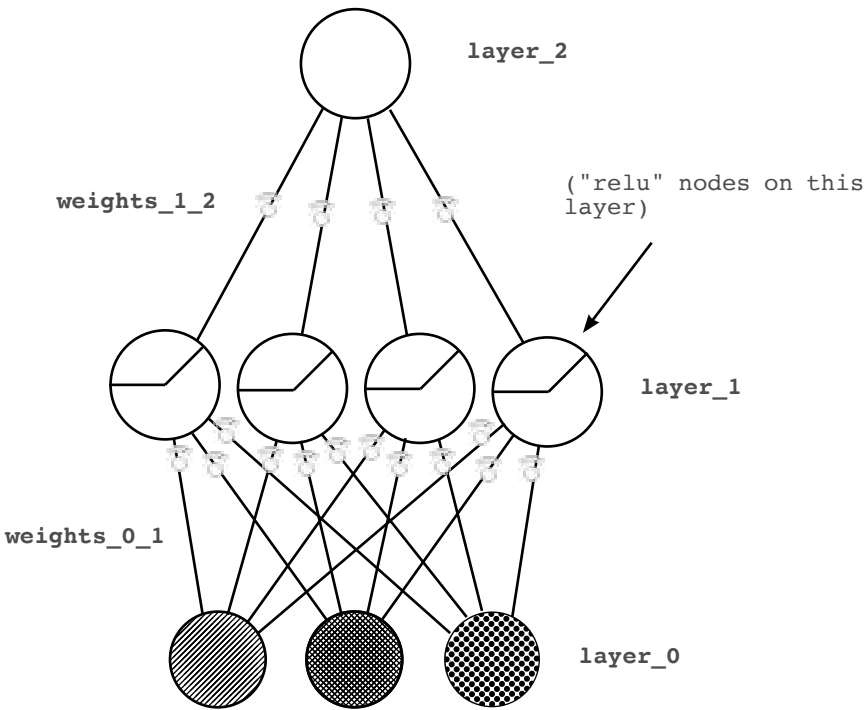
Our Previously Overcomplicated Visualization

While we simplify our mental picture, let's simplify our visualization as well.

At this point, I would expect that the visualization of neural networks in your head is something like the picture in the bottom right (because that's the one we used). We have our input dataset in `layer_0`, connected by a weight matrix (bunch of lines) to `layer_1` and so on. This was a very useful tool to learn the basics of how collections of weights and layers come together to learn a function.

However, moving forward, this picture simply has too much detail. Given our correlation summarization, we already know that we no longer need to worry about how individual weights are updated. Later layers already know how to communicate to earlier layers and tell them "hey, I need higher signal" or "hey, I need lower signal". Truth be told, we don't really care about the actual weight values anymore, only that they're behaving how they should, properly capturing correlation in a way that generalizes.

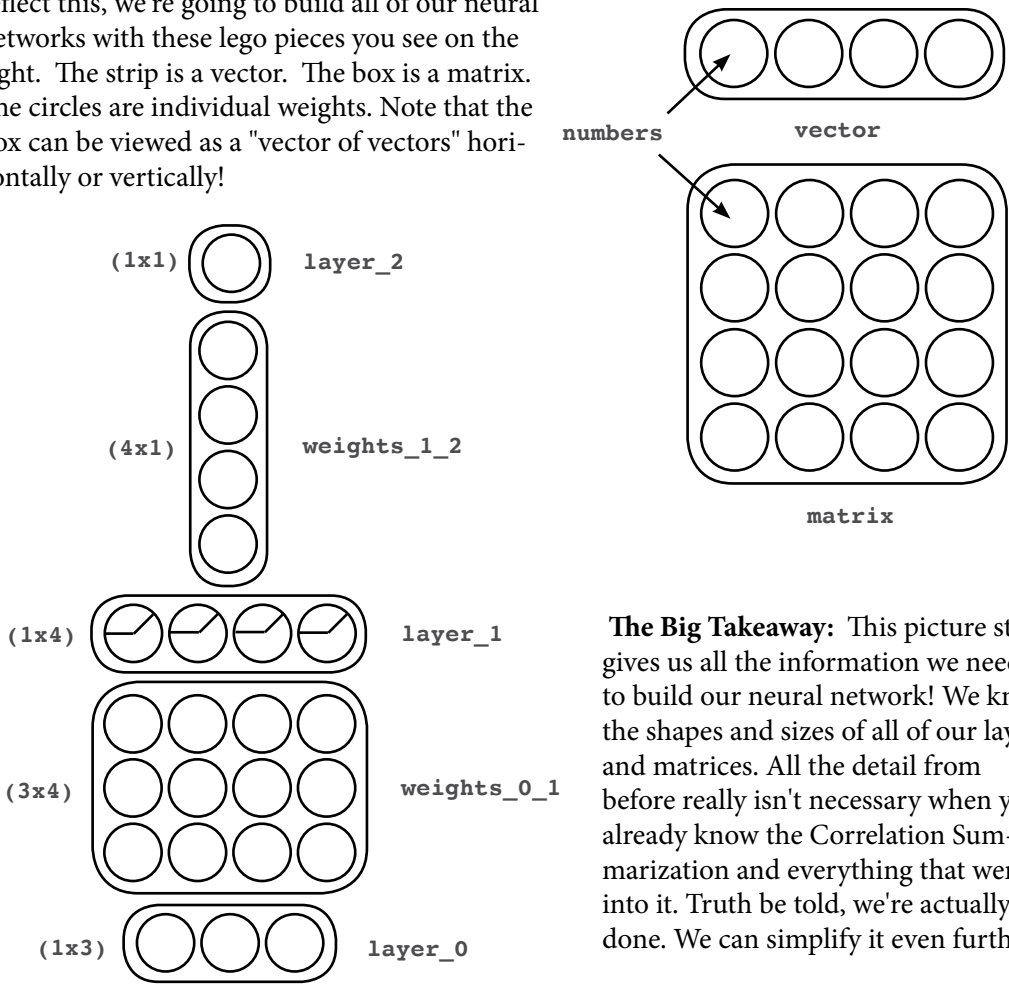
So, to reflect the change in our attention, we need to update our visualization on paper. We're also going to do a few more things which will make sense later. As you know, our neural network is really just a series of weight matrices. When we're using the network, we also end up creating vectors corresponding to each layer. In our picture, the weight matrices are the lines going from node to node, and the vectors are the strips of nodes themselves. For example, `weights_1_2` is a matrix, `weights_0_1` is a matrix, and `layer_1` is a vector. In later chapters, we're going to be arranging vectors and matrices in increasingly creative ways, so instead of all this detail showing each node connected by each weight (which gets hard to read if we have, say, 500 nodes in `layer_1`), let's instead just think in general terms. Let's think of them as just vectors and matrices of arbitrary size.



Our Simplified Visualization

Neural networks are like legos, and each block is a vector or matrix.

Moving forward, we're going to be building new neural network architectures in the same way that people build new structures with lego pieces. The great thing about our correlation summarization is that all the bits and pieces that lead to it (backpropagation, gradient descent, alpha, dropout, mini-batching, etc.) don't really depend on a particular configuration of our legos! No matter how we piece together our series of matrices, gluing them together with layers, our neural network will try to learn the pattern in our data by modifying the weights between wherever we put the input layer and the output layer. To reflect this, we're going to build all of our neural networks with these lego pieces you see on the right. The strip is a vector. The box is a matrix. The circles are individual weights. Note that the box can be viewed as a "vector of vectors" horizontally or vertically!



The Big Takeaway: This picture still gives us all the information we need to build our neural network! We know the shapes and sizes of all of our layers and matrices. All the detail from before really isn't necessary when you already know the Correlation Summarization and everything that went into it. Truth be told, we're actually not done. We can simplify it even further!

Simplifying Even Further

The dimensionality of our matrices are determined by the layers!

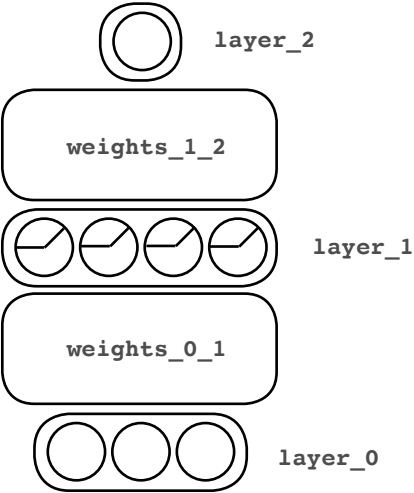
On the previous page, you might have noticed a pattern. Each matrix's dimensionality (number of rows and columns) has a direct relationship to the dimensionality of the layers before and after them! Thus, we can actually simplify our visualization even further! Consider the visualization on the right. We still have all of the information we need to build the neural network. We can infer that `weight_0_1` is a (3x4) matrix because the previous layer (`layer_0`) has 3 dimensions and the next layer (`layer_1`) has 4 dimensions. Thus, in order for the matrix to be big enough to have a single weight connecting each node in `layer_0` to each node in `layer_1`, it must be a (3x4) matrix!

This allows us to really start thinking about our neural networks using the correlation summarization. All this neural network is going to do is adjust the weights to find correlation between `layer_0` and `layer_2`. It's going to do this using all of the methods we've mentioned so far in this book. However, we will find that the different configurations of weights and layers between our input and output layers have a strong impact on whether or not the network is successful in finding correlation (and/or how fast it finds correlation).

The particular configuration of layers and weights in a neural network is called its **architecture**, and we will spend the majority of the rest of this book discussing the pros and cons of various architectures. As the correlation summarization reminds us, the neural network adjusts weights to find correlation between the input and output layers, sometimes even inventing correlation in the hidden layers. We will find that different architectures *channel signal to make correlation easier to discover*.

Good neural architectures channel signal so that correlation is easy to discover. Great architectures also filter noise to help prevent overfitting.

Much of research into neural networks is about finding new architectures that can find correlation faster and generalize better to unseen data. We will spend the vast majority of the rest of this book discussing new architectures.



Let's See This Network Predict

Let's picture data from our streetlight example flowing through the system.

1

A single datapoint from our "streetlight" dataset is selected. `layer_0` is set to the correct values. (Pictured right.)

2

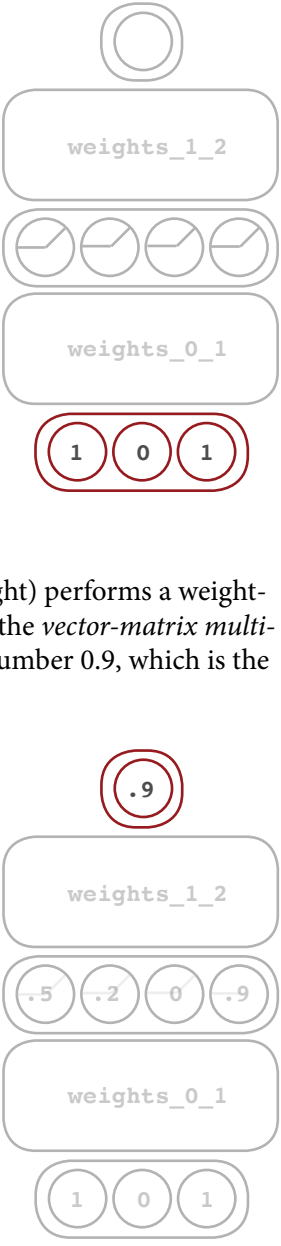
Four different weighted sums of `layer_0` are performed. The four weighted sums are performed by `weights_0_1`. As a reminder, this process is called *vector-matrix multiplication*. These four values are deposited into the four positions of `layer_1` and passed through the "relu" function (setting negative values to 0). To be clear, the 3rd value from the left in `layer_1` would have been negative, but the "relu" function sets it to 0. (Pictured left.)

3

The final step (pictured bottom right) performs a weighted average of `layer_1`, again using the *vector-matrix multiplication* process. This yields the number 0.9, which is the network's final prediction.

Review: Vector-Matrix Multiplication

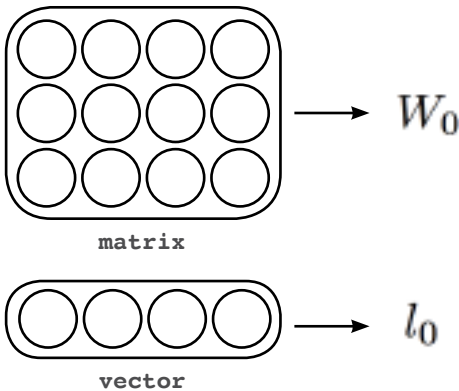
Vector-Matrix multiplication performs multiple *weighted sums* of a vector. The matrix must have the same number of rows as the vector has values, so that each column in the matrix performs a unique weighted sum. Thus, if the matrix has 4 columns, 4 weighted sums will be generated. The weightings of each sum are performed depending on the values of the matrix.



Visualizing Using Letters Instead of Pictures

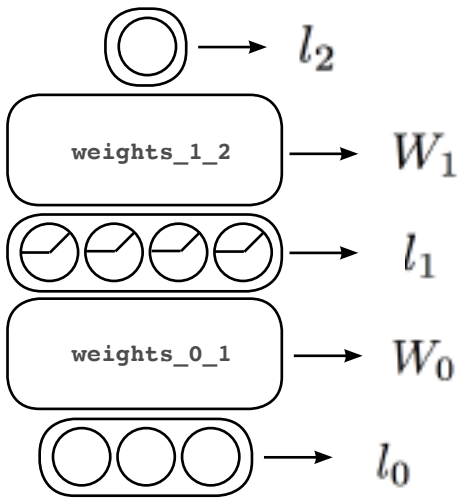
All these pictures and detailed explanations are actually a simple piece of algebra.

Just like we defined simpler pictures for our matrix and vector, we can actually perform the same visualization in the form of letters. Let me show you.



So, how do we visualize a matrix using math? Well, we simply pick a capital letter! I try to pick one that's easy to remember such as "W" for "weights". The little 0 down there? That just means it's probably one of several Ws. In our case, our network has 2. Perhaps surprisingly, I literally could have picked any capital letter, and the little 0 at the bottom is extra. It just lets me call all my weight matrices W so I can keep them apart. It's your visualization! Make it easy to remember!

How do we visualize a vector using math? Well, we pick a lowercase letter! Why did I choose the letter "L"? Well, since I have a bunch of vectors that are "layers", I thought "L" would be easy to remember! Why did I choose to call it "L"-zero? Well, Since I have multiple "layers", it seems nice to make all of them "L"s and just number them instead of having to think of new letters for every layer. There's no wrong answer here! How about that? So, if that's how we visualize matrices and vectors, in math, what do all of the pieces in our network look like? Well, to the right you can see a nice selection of variables pointing to their respective sections of the neural network. However, simply defining them doesn't show how they relate. Let's combine our variables together via vector-matrix multiplication!



Linking Our Variables

Our letters can be combined together to indicate functions and operations.

Vector-matrix multiplication is very simple. If you want to visualize that two letters are being multiplied by each other, you literally just put them next to each other. For example:

algebra	translation
$l_0 W_0$	"take the layer 0 vector and perform vector-matrix multiplication with the weight matrix 0"
$l_1 W_1$	"take the layer 1 vector and perform vector-matrix multiplication with the weight matrix 1"

We can even throw in arbitrary functions like "relu" using notation that looks almost exactly like the Python code! This really is crazy intuitive stuff!

$l_1 = \text{relu}(l_0 W_0)$	"to create the layer 1 vector, take the layer 0 vector and perform vector-matrix multiplication with the weight matrix 0, then perform the "relu" function on the output (setting all negative numbers to 0)."
$l_2 = l_1 W_1$	"to create the layer 2 vector, take the layer 1 vector and perform vector-matrix multiplication with the weight matrix 1"

If you notice, the layer 2 algebra actually contains layer 1 as an input variable. This means we can actually represent our *entire neural network* in one expression by just chaining them

together. Thus, all the logic in our forward propagation step can be contained in this one formula in the bottom left corner of this page. Note: baked into this formula is the assumption that our vectors and matrices have the right dimensions.

$$l_2 = \text{relu}(l_0 W_0) W_1$$

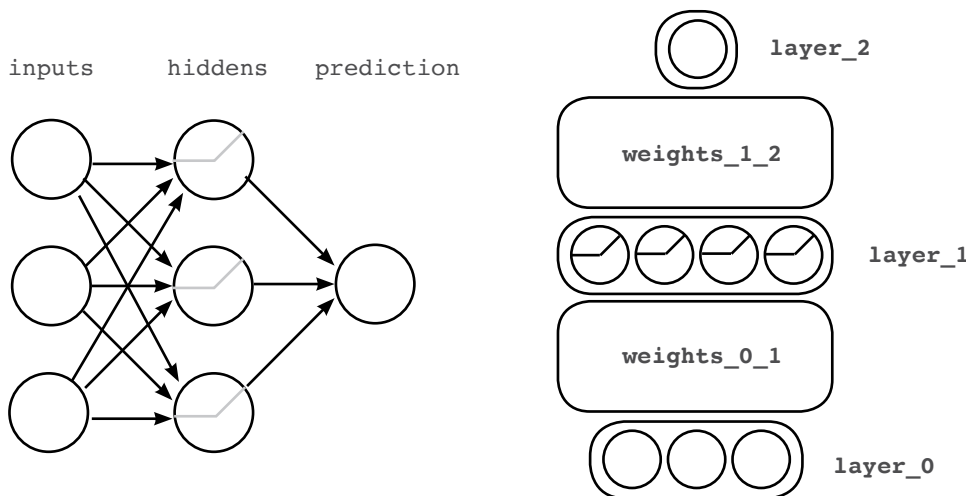
Everything Side-by-Side

Let's see our visualization, algebra formula, and python code in one place.

I don't think too much dialogue is necessary on this page. Just take a minute and look at each piece of forward propagation through these 4 different ways of seeing it. It is my hope that you would truly grok forward propagation, and understand the architecture by seeing it from different perspectives, all in one place.

```
layer_2 = relu(layer_0.dot(weights_0_1)).dot(weights_1_2)
```

$$l_2 = \text{relu}(l_0 W_0) W_1$$



The Importance of Visualization Tools

We're going to be studying new architectures.

In the following chapters, we're going to be taking these vectors and matrices and combining them in some very creative ways. My ability to describe each architecture for you is entirely dependent on our ability to have a mutually agreed upon language for describing them. Thus, please don't move beyond this chapter until you can clearly see how forward propagation manipulates these vectors and matrices, and how these various forms of describing them are articulated.

Good neural architectures channel signal so that correlation is easy to discover. Great architectures also filter noise to help prevent overfitting.

As mentioned previously, a neural architecture controls how signal flows through a network. How we create these architectures will affect the ways in which the network can detect correlation. We will find that we want to create architectures that maximize the network's ability to focus on the areas where meaningful correlation exist, and minimize the network's ability to focus on the areas that simply contain noise.

However, different datasets and domains have different characteristics. For example, image data has different kinds of signal and noise than text data. Thus, we will find that even though neural networks can be used in many situations, different architectures will be better suited to different problems because of their ability to locate certain types of correlations. So, for the next few chapters, we're going to explore how you can modify neural networks to specifically find the correlation you're looking for. See you there!

IN THIS CHAPTER •

Overfitting

Dropout

Batch Gradient Descent

“ With four parameters I can fit an elephant, and with five I can make him wiggle his trunk. ”

— JOHN VON NEUMANN

3 Layer Network on MNIST

Let's return to our MNIST dataset and attempt to classify it with our new network

In last several chapters, we have learned that neural networks model correlation. In fact, the hidden layers (the middle one in our 3 layer network) can even create "intermediate" correlation to help solve for a task (seemingly out of mid air). How do we know that our network is creating good correlation?

Back when we learned about Stochastic Gradient Descent with Multiple Inputs, we ran an experiment where we froze one weight and then asked our network to continue training. As it was training, we watch the "dots find the bottom of the bowls" as it were. We were watching the weights become adjusted to minimize the error.

When we froze the weight, however, we were surprised to see that the frozen weight still found the bottom of the bowl! For some reason, the bowl moved so that the frozen weight value became optimal. Furthermore, if we unfroze the weight after training to do some more training, it wouldn't learn! Why? Well, the error had already fallen to 0! As far as the network was concerned, there was nothing more to learn!

This begs the question, what if the input to the frozen weight was actually important to predicting baseball victory in the real world? What if the network had figured out a way to accurately predict the games in the training dataset (because that's what networks do, they minimize the error), but somehow forgot to include a valuable input?

This phenomenon is, unfortunately, extremely common in neural networks. In fact, one might say it's the "Arch Nemesis" of neural networks, **Overfitting**, and unfortunately, the more powerful your neural network's expressive power (more layers / weights), the more prone the network is to overfit. So, there's an everlasting battle going on in research where people continually find tasks that need more powerful layers but find themselves having to do lots of problem solving to make sure the network doesn't "overfit".

In this chapter, we're going to study the basics of Regularization, which are key to combatting overfitting in neural networks. In order to do this, we're going to first start with our most powerful neural network (3 layer network with relu hidden layer) on our most challenging task (MNIST digit classification).

So, to start, go ahead and train the network on the following page. You should see the same results as those listed below. Alas! Our network learned to perfectly predict the training data! Should we celebrate?

```
import sys, numpy as np
from keras.datasets import mnist

(x_train, y_train), (x_test, y_test) = mnist.load_data()

images, labels = (x_train[0:1000].reshape(1000,28*28) \
                  255, y_train[0:1000])
one_hot_labels = np.zeros((len(labels),10))

for i,l in enumerate(labels):
    one_hot_labels[i][l] = 1
labels = one_hot_labels

test_images = x_test.reshape(len(x_test),28*28) / 255
test_labels = np.zeros((len(y_test),10))
for i,l in enumerate(y_test):
    test_labels[i][l] = 1

np.random.seed(1)
relu = lambda x:(x>=0) * x # returns x if x > 0, return 0 otherwise
relu2deriv= lambda x: x>=0 # returns 1 for input > 0, return 0 otherwise
alpha, iterations, hidden_size, pixels_per_image, num_labels = \
    (0.005, 350, 40, 784, 10)
weights_0_1 = 0.2*np.random.random((pixels_per_image,hidden_size)) - 0.1
weights_1_2 = 0.2*np.random.random((hidden_size,num_labels)) - 0.1

for j in range(iterations):
    error, correct_cnt = (0.0, 0)

    for i in range(len(images)):
        layer_0 = images[i:i+1]
        layer_1 = relu(np.dot(layer_0,weights_0_1))
        layer_2 = np.dot(layer_1,weights_1_2)
        error += np.sum((labels[i:i+1] - layer_2) ** 2)
        correct_cnt += int(np.argmax(layer_2) == \
                                np.argmax(labels[i:i+1]))

        layer_2_delta = (labels[i:i+1] - layer_2)
        layer_1_delta = layer_2_delta.dot(weights_1_2.T)\
                        * relu2deriv(layer_1)
        weights_1_2 += alpha * layer_1.T.dot(layer_2_delta)
        weights_0_1 += alpha * layer_0.T.dot(layer_1_delta)

    sys.stdout.write("\r"+ \
        " I:"+str(j)+ \
        " Error:" + str(error/float(len(images)))[0:5] +\
        " Correct:" + str(correct_cnt/float(len(images))))
```

```
....
I:349 Error:0.108 Correct:1.0
```

Well... That was easy!

Our neural network perfectly learned to predict all 1000 images!

So, in some ways, this is a real victory. Our neural network was able to take a dataset of 1000 images and learn to correlate each input image with the correct label. How did it do this? Well, it simply iterated through each image, made a prediction, and then updated each weight ever so slightly so that the prediction was better next time. Doing this long enough on all the images eventually reached a state where the network could correctly predict on all of the images!

Perhaps a non-obvious question: how well will the neural network do on an image that it hasn't seen before? In other words, how well will it do on an image that wasn't part of the 1000 images it was trained on? Well, our MNIST dataset has many more images than just the 1000 we trained on; let's try it out! In the notebook from the previous page there are two variables called `test_images` and `test_labels`. If you execute the following code, it will run the neural network on these images and evaluate how well the network classifies them.

```
if(j % 10 == 0 or j == iterations-1):
    error, correct_cnt = (0.0, 0)

    for i in range(len(test_images)):

        layer_0 = test_images[i:i+1]
        layer_1 = relu(np.dot(layer_0,weights_0_1))
        layer_2 = np.dot(layer_1,weights_1_2)

        error += np.sum((test_labels[i:i+1] - layer_2) ** 2)
        correct_cnt += int(np.argmax(layer_2) == \
                                np.argmax(test_labels[i:i+1]))
    sys.stdout.write(" Test-Err:" + str(error/float(len(test_images)))[0:5] +\
        " Test-Acc:" + str(correct_cnt/float(len(test_images))))
    print()
```

Error:0.653 Correct:0.7073

The network did horribly! It only predicted with an accuracy of 70.7%. Why does it do so terribly on these new testing images when it learned to predict with 100% accuracy on the training data? How strange! We call this 70.7% number the *test accuracy*. It's the accuracy of the neural network on data that the network was NOT trained on. This number is important because it simulates how well our neural network will perform if we tried to use it in the real world (which only gives us images we haven't seen before). This is the score that matters!

Memorization vs Generalization

"Memorizing" 1000 images is easier than "Generalizing" to all images.

Let's consider again how the neural network learns. It adjusts each weight in each matrix so that the network is better able to take *specific inputs* and make a *specific prediction*. Perhaps a better question might be, "If we train it on 1000 images which it learns to predict perfectly, why does it work on other images at all?" As you might expect, when our fully trained neural network is applied to a new image, it is only guaranteed to work well if the new image is *nearly identical to an image from the training data*. Why? Well, the neural network only learned to transform input data to output data for *very specific input configurations*. If you give it something that doesn't look *familiar*, then it will predict randomly!

Well this makes neural networks kind of pointless! What's the point in a neural network only working on the data you trained it on? You already know the correct classifications for those datapoints! Neural networks are only useful if they work on data we don't already know the answer to! As it turns out, there's a way to combat this. For a clue, below I've printed out BOTH the training *and testing accuracy* of the neural network *as it was training* (every 10 iterations). Notice anything interesting? You should see a clue to better networks!

```
I:0 Train-Err:0.722 Train-Acc:0.537 Test-Err:0.601 Test-Acc:0.6488
I:10 Train-Err:0.312 Train-Acc:0.901 Test-Err:0.420 Test-Acc:0.8114
I:20 Train-Err:0.260 Train-Acc:0.93 Test-Err:0.414 Test-Acc:0.8111
I:30 Train-Err:0.232 Train-Acc:0.946 Test-Err:0.417 Test-Acc:0.8066
I:40 Train-Err:0.215 Train-Acc:0.956 Test-Err:0.426 Test-Acc:0.8019
I:50 Train-Err:0.204 Train-Acc:0.966 Test-Err:0.437 Test-Acc:0.7982
I:60 Train-Err:0.194 Train-Acc:0.967 Test-Err:0.448 Test-Acc:0.7921
I:70 Train-Err:0.186 Train-Acc:0.975 Test-Err:0.458 Test-Acc:0.7864
I:80 Train-Err:0.179 Train-Acc:0.979 Test-Err:0.466 Test-Acc:0.7817
I:90 Train-Err:0.172 Train-Acc:0.981 Test-Err:0.474 Test-Acc:0.7758
I:100 Train-Err:0.166 Train-Acc:0.984 Test-Err:0.482 Test-Acc:0.7706
I:110 Train-Err:0.161 Train-Acc:0.984 Test-Err:0.489 Test-Acc:0.7686
I:120 Train-Err:0.157 Train-Acc:0.986 Test-Err:0.496 Test-Acc:0.766
I:130 Train-Err:0.153 Train-Acc:0.99 Test-Err:0.502 Test-Acc:0.7622
I:140 Train-Err:0.149 Train-Acc:0.991 Test-Err:0.508 Test-Acc:0.758
.....
I:210 Train-Err:0.127 Train-Acc:0.998 Test-Err:0.544 Test-Acc:0.7446
I:220 Train-Err:0.125 Train-Acc:0.998 Test-Err:0.552 Test-Acc:0.7416
I:230 Train-Err:0.123 Train-Acc:0.998 Test-Err:0.560 Test-Acc:0.7372
I:240 Train-Err:0.121 Train-Acc:0.998 Test-Err:0.569 Test-Acc:0.7344
I:250 Train-Err:0.120 Train-Acc:0.999 Test-Err:0.577 Test-Acc:0.7316
I:260 Train-Err:0.118 Train-Acc:0.999 Test-Err:0.585 Test-Acc:0.729
I:270 Train-Err:0.117 Train-Acc:0.999 Test-Err:0.593 Test-Acc:0.7259
I:280 Train-Err:0.115 Train-Acc:0.999 Test-Err:0.600 Test-Acc:0.723
I:290 Train-Err:0.114 Train-Acc:0.999 Test-Err:0.607 Test-Acc:0.7196
I:300 Train-Err:0.113 Train-Acc:0.999 Test-Err:0.614 Test-Acc:0.7183
I:310 Train-Err:0.112 Train-Acc:0.999 Test-Err:0.622 Test-Acc:0.7165
I:320 Train-Err:0.111 Train-Acc:0.999 Test-Err:0.629 Test-Acc:0.7133
I:330 Train-Err:0.110 Train-Acc:0.999 Test-Err:0.637 Test-Acc:0.7125
I:340 Train-Err:0.109 Train-Acc:1.0 Test-Err:0.645 Test-Acc:0.71
I:349 Train-Err:0.108 Train-Acc:1.0 Test-Err:0.653 Test-Acc:0.7073
```


Overfitting in Neural Networks

Neural networks can get worse if we train them too much?

For some reason, our *test* accuracy went up for the first 20 iterations, and then slowly decreased as the network trained more and more (during which time the *training* accuracy was still improving). This is very common in neural networks. Let me explain the phenomenon via analogy.

Imagine that you are creating a *mold* for a common dinner fork, but instead of using it to create other forks you want to use it to *identify* if a particular utensil is in fact, a fork. If an object fits in the mold, you would conclude that the object is a fork, and if it does not, then you would conclude that it is *not* a fork.

Let's say you set out to make this mold, and you start with a wet piece of clay and a big bucket of 3-pronged forks, spoons, and knives. You then press each of the forks into the same place in the mold to create an outline, which sortof looks like a mushy fork. You repeatedly place all the forks in the clay over and over, hundreds of times. When you let the clay dry, you then find that none of the spoons or knives fit into this mold, but all of the forks fit into this mold. Awesome! You did it! You correctly made a mold that can only fit the shape of a fork!

However, what happens if someone hands you a 4-pronged fork? You look in your mold and notice that there is a specific outlines for three, thin prongs in your clay. Your 4-pronged fork doesn't fit! Why not? It's still a fork!

This is because the clay wasn't molded on any 4-pronged forks. It was only molded on the 3-pronged variety. In this way, the clay has **overfit** to only recognize the types of forks that it was "trained" to shape.

This is exactly the same phenomenon that we just witnessed in our neural network. It's actually an even closer parallel than you might think. In truth, one way to view the weights of a neural network is as a *high dimensional shape*. As you train, this shape *molds* around the shape of your data, learning to distinguish one pattern from another. Unfortunately, the images in our testing dataset were *slightly* different than the patterns in our training dataset. This caused our network to fail on many of our *testing* examples.

This phenomenon is known as **Overfitting**. A more official definition of a neural network that overfits is a neural network that has learned the *noise* in the dataset instead of only making decisions based on the *true signal*.

Where Overfitting Comes From

What causes our neural networks to overfit?

Let's alter this scenario just a bit. Picture the fresh clay in your head again (unmolded). What if you only pushed a single fork into it? Assuming the clay was very thick, it wouldn't have as much detail as the previous mold did (which was imprinted many times). Thus, it would be only a *very general shape of a fork*. This shape might, in fact, be compatible with both the 3 and 4 pronged variety of fork, because it's still a very *fuzzy* imprint.

Assuming this information, our mold actually got worse at our testing dataset as we imprinted more forks because it learned more *detailed information* about the training dataset that it was being molded to. This caused it to reject images that were even the slightest bit off from what it had repeatedly seen in the training data. So, what is this *detailed information* in our images that is incompatible with our test data? In our fork analogy, this was the number of prongs on the fork. In images, it's generally referred to as *noise*. In reality, it's a bit more nuanced. Consider these two dog pictures.



Everything that makes these pictures unique *beyond* what captures the essence of "dog" is included in this term *noise*. In the picture on the left, the pillow and the background are both noise. In the picture on the right, the empty, middle blackness of the dog is actually a form of *noise* as well. It's really the edges that tell us that it's a dog. The middle blackness doesn't really tell us anything. On the picture on the left, however, the middle of the dog has the furry texture and color of a dog, which could help the classifier correctly identify it.

So, how do we get our neural networks to train only on the *signal* (the essence of a dog) and ignore the *noise* (other stuff irrelevant to the classification)? Well, one way of doing this is by **early stopping**. It turns out that a large amount of noise comes in the fine grained detail of an image, and most of the signal (for objects) is found in the general *shape* and perhaps *color* of the image.

The Simplest Regularization: Early Stopping

Stop training the network when it starts getting worse.

So, how do we get our neural networks to train only on the *signal* (the essence of a dog) and ignore the *noise* (other stuff irrelevant to the classification)? Well, one way of doing this is by **early stopping**. It turns out that a large amount of noise often comes in the fine grained detail present in input data, and most of the signal is found in the more general characteristics of your input data (for images, this is things like big shapes and color).

So, how do we get our neural network to ignore the fine grained detail and only capture the general information present in our data (i.e. the general shape of dog or of an MNIST digit)? Well, we don't let the network train long enough to learn it! Just like in the "fork mold" example, it takes many forks imprinted many times to create the perfect outline of a 3 pronged fork. The first few imprints only generally capture the shallow outline of a fork. The same can be said for neural networks. As a result, early stopping is the cheapest form of "regularization", and if you're in a pinch, it can be quite effective.

This brings us to the subject that this chapter is all about, **Regularization**. Regularization is a subfield of methods for getting your model to *generalize* to new data points (instead of just memorize the training data). It's a subset of methods that help your neural network learn the *signal* and ignore the *noise*. In our case, it's a toolset at our disposal to create neural networks that have these properties.

Regularization

A subset of methods used to encourage generalization in learned models, often by increasing the difficulty for a model to learn the fine-grained details of training data.

So, the next question might be, how do we know when to stop? In truth, the only real way to know is to run the model on data that isn't in your training dataset. This is typically done using a *second* test dataset called a "validation set". In some circumstances, if we used our test set for knowing when to stop, we could actually *overfit to our test set*. So, as a general rule, we don't use it to control training. We use a validation set instead.

Industry Standard Regularization: Dropout

The Method: Randomly turning neurons off (setting to 0) during training.

So, this regularization technique really is as simple as it sounds. During training, you randomly set neurons in your network to zero (and usually the deltas on the same nodes during backpropagation, but you technically don't have to). This causes the neural network to train exclusively using *random subsections* of the neural network. Believe it or not, this regularization technique is generally accepted as the go-to, state-of-the-art regularization technique for the vast majority of networks. It's methodology is simple and inexpensive, although the intuitions behind *why* it works are a bit more complex.

Why does Dropout Work? (perhaps oversimplified)

Dropout makes our big network act like a little one by randomly training little subsections of the network at a time, and little networks don't overfit.

It turns out that the smaller a neural network is, the less it is able to overfit. Why? Well, small neural networks don't have very much expressive power. They can't latch onto the more granular details (i.e. noise) that tend to be the source of overfitting. They only have "room" to capture the big, obvious, high level features.

This notion of "room" or "capacity" is actually a really important one for you to keep in your mind. You can think of it like this. Remember our "clay" analogy from a few pages ago? Imagine if your clay was actually made of "sticky rocks" that were the size of dimes. Would that clay be able to make a very good imprint of a fork? Of course not! Why? Well, those stones are much like our weights. They form around our data, capturing the patterns we're interested in. If we only have a few, larger stones, then it can't capture nuanced detail. Each stone instead is pushed on by large parts of the fork, more or less *averaging* the shape (ignoring fine creases and corners).

Imagine again clay made up of very fine-grained sand. It's actually made up of millions and millions of small stones that can fit into every nook and cranny of a fork. This is what gives *big* neural networks the expressive power they often use to overfit to a dataset.

So, how do we have the power of a large neural network with the resistance to overfitting of the small neural network? We take our big neural network and turn off nodes randomly. What happens when you take a big neural network and only use a small part of it? Well, it behaves like a small neural network! However, when we do this randomly over potentially millions of different "sub-networks", the sum total of the entire network still maintains its expressive power! Neat, eh?

Why Dropout Works: Ensembling Works

Dropout is actually a form of training a bunch of networks and averaging them

Something to keep in mind: neural networks always start out randomly. Why does this matter? Well, since neural networks learn by trial and error, this ultimately means that every neural network learns just a little bit *differently*. It might learn equally effectively, but no two neural networks are ever exactly the same (unless they start out exactly the same for some random or intentional reason).

This has an interesting property. When you overfit two neural networks, no two neural networks overfit in exactly the same way. Why? Well, overfitting only occurs until every training image can be predicted perfectly, at which point the error == 0 and the network stops learning (even if you keep iterating). However, since each neural network starts by predicting randomly, then adjusting its weights to make better predictions, each network inevitably makes different mistakes, resulting in different updates. This culminates in a core concept:

While it is very likely for large, unregularized neural networks to overfit to noise, it is very *unlikely* for them to overfit to the **same** noise.

Why do they not overfit to the same noise? Well, they start randomly, and they stop training once they have learned enough noise to disambiguate between all the images in the training set. Truth be told, our MNIST network only needs to find a handful of random pixels that happen to correlate with our output labels to overfit. However, this is contrasted with, perhaps, an even more important concept:

Neural networks, even though they are randomly generated, still start by learning the biggest, most broadly sweeping features before learning much about the noise.

The takeaway is this; if you train 100 neural networks (all initialized randomly), they will each tend to latch onto different noise but similar broad *signal*. Thus, when they make mistakes, they often make *differing* mistakes. This means that if we allowed them to vote equally, their noise would tend to cancel out, revealing only what they have all learned in common, *the signal*.

Dropout In Code

Here's how you actually use Dropout in the real world

In our MNIST classification model, we're going to add Dropout to our hidden layer, such that 50% of the nodes are turned off (randomly) during training. Perhaps you will be surprised that this is actually only a 3 line change in our code. Below you can see a familiar snippet from our previous neural network logic with our dropout mask added:

```
i = 0
layer_0 = images[i:i+1]
dropout_mask = np.random.randint(2,size=layer_1.shape)

layer_1 *= dropout_mask * 2
layer_2 = np.dot(layer_1, weights_1_2)

error += np.sum((labels[i:i+1] - layer_2) ** 2)

correct_cnt += int(np.argmax(layer_2) == \
                    np.argmax(labels[i:i+1]))

layer_2_delta = (labels[i:i+1] - layer_2)
layer_1_delta = layer_2_delta.dot(weights_1_2.T)\
                * relu2deriv(layer_1)

layer_1_delta *= dropout_mask

weights_1_2 += alpha * layer_1.T.dot(layer_2_delta)
weights_0_1 += alpha * layer_0.T.dot(layer_1_delta)
```

Indeed, to implement Dropout on a layer (in this case layer_1), one simply multiplies the layer_1 values by a random matrix of 1s and 0s. This has the affect of randomly "turning off" nodes in layer_1 by setting them to equal 0. Note that our dropout_mask uses what's called a "50% bernoulli distribution" such that 50% of the time, each value in the dropout_mask is a 1, and (1 - 50% = 50%) of the time, it is a 0.

This is followed by something that may seem a bit peculiar. We multiply layer_1 by two! Why do we do this? Well, remember that layer_2 is going to perform a *weighted sum* of layer_1. Even though it's weighted, it's still a **sum** over the values of layer_1. Thus, if we turn off half the nodes in layer_1, then that sum is going to be cut in half! Thus, layer_2 would increase its sensitivity to layer_2, kind of like a person leaning closer to a radio when the volume is too low to better hear it. However, at test time, when we no longer use Dropout, the volume would be back up to normal! You may be surprised to find that this throws off layer_2's ability to *listen* to layer_1. Thus, we need to counter this by multiplying layer_1 by (1 / the percentage of turned on nodes). In this case, that's 1/0.5 which equals 2. In this way, the volume of layer_1 is the same between training and testing, despite Dropout.

```

import numpy, sys
np.random.seed(1)
def relu(x):
    return (x >= 0) * x # returns x if x > 0
                        # returns 0 otherwise

def relu2deriv(output):
    return output >= 0 #returns 1 for input > 0

alpha, iterations, hidden_size = (0.005, 300, 100)
pixels_per_image, num_labels = (784, 10)

weights_0_1 = 0.2*np.random.random((pixels_per_image,hidden_size)) - 0.1
weights_1_2 = 0.2*np.random.random((hidden_size,num_labels)) - 0.1

for j in range(iterations):
    error, correct_cnt = (0.0,0)
    for i in range(len(images)):
        layer_0 = images[i:i+1]
        layer_1 = relu(np.dot(layer_0,weights_0_1))
        dropout_mask = np.random.randint(2, size=layer_1.shape)
        layer_1 *= dropout_mask * 2
        layer_2 = np.dot(layer_1,weights_1_2)

        error += np.sum((labels[i:i+1] - layer_2) ** 2)
        correct_cnt += int(np.argmax(layer_2) == \
                               np.argmax(labels[i:i+1]))
        layer_2_delta = (labels[i:i+1] - layer_2)
        layer_1_delta = layer_2_delta.dot(weights_1_2.T) * relu2deriv(layer_1)
        layer_1_delta *= dropout_mask

        weights_1_2 += alpha * layer_1.T.dot(layer_2_delta)
        weights_0_1 += alpha * layer_0.T.dot(layer_1_delta)

    if(j%10 == 0):
        test_error = 0.0
        test_correct_cnt = 0

        for i in range(len(test_images)):
            layer_0 = test_images[i:i+1]
            layer_1 = relu(np.dot(layer_0,weights_0_1))
            layer_2 = np.dot(layer_1, weights_1_2)

            test_error += np.sum((test_labels[i:i+1] - layer_2) ** 2)
            test_correct_cnt += int(np.argmax(layer_2) == \
                                       np.argmax(test_labels[i:i+1]))

        sys.stdout.write("\n" + \
            "I:" + str(j) + \
            " Test-Err:" + str(test_error/ float(len(test_images)))[0:5] + \
            " Test-Acc:" + str(test_correct_cnt/ float(len(test_images)))+ \
            " Train-Err:" + str(error/ float(len(images)))[0:5] + \
            " Train-Acc:" + str(correct_cnt/ float(len(images))))

```

Dropout Evaluated on MNIST

Here's how you actually use Dropout in the real world.

If you remember from before, our neural network (without Dropout) previously reached a test accuracy of 81.14% before falling down to finish training at 70.73% accuracy. When we add dropout, our neural network instead behaves this way.

```

I:0 Test-Err:0.641 Test-Acc:0.6333 Train-Err:0.891 Train-Acc:0.413
I:10 Test-Err:0.458 Test-Acc:0.787 Train-Err:0.472 Train-Acc:0.764
I:20 Test-Err:0.415 Test-Acc:0.8133 Train-Err:0.430 Train-Acc:0.809
I:30 Test-Err:0.421 Test-Acc:0.8114 Train-Err:0.415 Train-Acc:0.811
I:40 Test-Err:0.419 Test-Acc:0.8112 Train-Err:0.413 Train-Acc:0.827
I:50 Test-Err:0.409 Test-Acc:0.8133 Train-Err:0.392 Train-Acc:0.836
I:60 Test-Err:0.412 Test-Acc:0.8236 Train-Err:0.402 Train-Acc:0.836
I:70 Test-Err:0.412 Test-Acc:0.8033 Train-Err:0.383 Train-Acc:0.857
I:80 Test-Err:0.410 Test-Acc:0.8054 Train-Err:0.386 Train-Acc:0.854
I:90 Test-Err:0.411 Test-Acc:0.8144 Train-Err:0.376 Train-Acc:0.868
I:100 Test-Err:0.411 Test-Acc:0.7903 Train-Err:0.369 Train-Acc:0.864
I:110 Test-Err:0.411 Test-Acc:0.8003 Train-Err:0.371 Train-Acc:0.868
I:120 Test-Err:0.402 Test-Acc:0.8046 Train-Err:0.353 Train-Acc:0.857
I:130 Test-Err:0.408 Test-Acc:0.8091 Train-Err:0.352 Train-Acc:0.867
I:140 Test-Err:0.405 Test-Acc:0.8083 Train-Err:0.355 Train-Acc:0.885
I:150 Test-Err:0.404 Test-Acc:0.8107 Train-Err:0.342 Train-Acc:0.883
I:160 Test-Err:0.399 Test-Acc:0.8146 Train-Err:0.361 Train-Acc:0.876
I:170 Test-Err:0.404 Test-Acc:0.8074 Train-Err:0.344 Train-Acc:0.889
I:180 Test-Err:0.399 Test-Acc:0.807 Train-Err:0.333 Train-Acc:0.892
I:190 Test-Err:0.407 Test-Acc:0.8066 Train-Err:0.335 Train-Acc:0.898
I:200 Test-Err:0.405 Test-Acc:0.8036 Train-Err:0.347 Train-Acc:0.893
I:210 Test-Err:0.405 Test-Acc:0.8034 Train-Err:0.336 Train-Acc:0.894
I:220 Test-Err:0.402 Test-Acc:0.8067 Train-Err:0.325 Train-Acc:0.896
I:230 Test-Err:0.404 Test-Acc:0.8091 Train-Err:0.321 Train-Acc:0.894
I:240 Test-Err:0.415 Test-Acc:0.8091 Train-Err:0.332 Train-Acc:0.898
I:250 Test-Err:0.395 Test-Acc:0.8182 Train-Err:0.320 Train-Acc:0.899
I:260 Test-Err:0.390 Test-Acc:0.8204 Train-Err:0.321 Train-Acc:0.899
I:270 Test-Err:0.382 Test-Acc:0.8194 Train-Err:0.312 Train-Acc:0.906
I:280 Test-Err:0.396 Test-Acc:0.8208 Train-Err:0.317 Train-Acc:0.9
I:290 Test-Err:0.399 Test-Acc:0.8181 Train-Err:0.301 Train-Acc:0.908

```

Not only does the network instead peak out at a score of 82.36%, it also doesn't over fit nearly as badly, finishing training with a testing accuracy of 81.81%. Notice that the dropout also slows down the Training-Acc, which previously went straight to 100% and then stayed there.

This should point to what Dropout really is. It's noise. It makes it more difficult for the network to train on the training data. It's like running a marathon with weights on your legs. It's harder to train, but when you take them off for the big race, you end up running quite a bit faster because you trained for something that was much more difficult.

Batch Gradient Descent

A method for increasing the speed of training and the rate of convergence.

In the context of this chapter, I would like to briefly apply a concept introduced several chapters ago, the concept of mini-batched stochastic gradient descent. I won't go into too much detail, as it's something that's largely taken for granted in neural network training. Furthermore, it's a very simple concept that doesn't really get more advanced even with the most state of the art neural networks. Simply stated, previously we trained one training example at a time, updating the weights after each example. Now, we're going to train 100 training examples at a time averaging the weight updates between all 100 examples. The code for this training logic is on the next page, and the training/testing output is below.

```
I:0 Test-Err:0.815 Test-Acc:0.3832 Train-Err:1.284 Train-Acc:0.165
I:10 Test-Err:0.568 Test-Acc:0.7173 Train-Err:0.591 Train-Acc:0.672
I:20 Test-Err:0.510 Test-Acc:0.7571 Train-Err:0.532 Train-Acc:0.729
I:30 Test-Err:0.485 Test-Acc:0.7793 Train-Err:0.498 Train-Acc:0.754
I:40 Test-Err:0.468 Test-Acc:0.7877 Train-Err:0.489 Train-Acc:0.749
I:50 Test-Err:0.458 Test-Acc:0.793 Train-Err:0.468 Train-Acc:0.775
I:60 Test-Err:0.452 Test-Acc:0.7995 Train-Err:0.452 Train-Acc:0.799
I:70 Test-Err:0.446 Test-Acc:0.803 Train-Err:0.453 Train-Acc:0.792
I:80 Test-Err:0.451 Test-Acc:0.7968 Train-Err:0.457 Train-Acc:0.786
I:90 Test-Err:0.447 Test-Acc:0.795 Train-Err:0.454 Train-Acc:0.799
I:100 Test-Err:0.448 Test-Acc:0.793 Train-Err:0.447 Train-Acc:0.796
I:110 Test-Err:0.441 Test-Acc:0.7943 Train-Err:0.426 Train-Acc:0.816
I:120 Test-Err:0.442 Test-Acc:0.7966 Train-Err:0.431 Train-Acc:0.813
I:130 Test-Err:0.441 Test-Acc:0.7906 Train-Err:0.434 Train-Acc:0.816
I:140 Test-Err:0.447 Test-Acc:0.7874 Train-Err:0.437 Train-Acc:0.822
I:150 Test-Err:0.443 Test-Acc:0.7899 Train-Err:0.414 Train-Acc:0.823
I:160 Test-Err:0.438 Test-Acc:0.797 Train-Err:0.427 Train-Acc:0.811
I:170 Test-Err:0.440 Test-Acc:0.7884 Train-Err:0.418 Train-Acc:0.828
I:180 Test-Err:0.436 Test-Acc:0.7935 Train-Err:0.407 Train-Acc:0.834
I:190 Test-Err:0.434 Test-Acc:0.7935 Train-Err:0.410 Train-Acc:0.831
I:200 Test-Err:0.435 Test-Acc:0.7972 Train-Err:0.416 Train-Acc:0.829
I:210 Test-Err:0.434 Test-Acc:0.7923 Train-Err:0.409 Train-Acc:0.83
I:220 Test-Err:0.433 Test-Acc:0.8032 Train-Err:0.396 Train-Acc:0.832
I:230 Test-Err:0.431 Test-Acc:0.8036 Train-Err:0.393 Train-Acc:0.853
I:240 Test-Err:0.430 Test-Acc:0.8047 Train-Err:0.397 Train-Acc:0.844
I:250 Test-Err:0.429 Test-Acc:0.8028 Train-Err:0.386 Train-Acc:0.843
I:260 Test-Err:0.431 Test-Acc:0.8038 Train-Err:0.394 Train-Acc:0.843
I:270 Test-Err:0.428 Test-Acc:0.8014 Train-Err:0.384 Train-Acc:0.845
I:280 Test-Err:0.430 Test-Acc:0.8067 Train-Err:0.401 Train-Acc:0.846
I:290 Test-Err:0.428 Test-Acc:0.7975 Train-Err:0.383 Train-Acc:0.851
```

Notice that our training accuracy has a bit of a smoother trend to it than it did before. Taking an average weight update consistently creates this kind of phenomenon during training. As it turns out, individual training examples are very noisy in terms of the weight updates they generate. Thus, averaging them makes for a smoother learning process.

```
import numpy as np
np.random.seed(1)

def relu(x):
    return (x >= 0) * x # returns x if x > 0

def relu2deriv(output):
    return output >= 0 # returns 1 for input > 0

batch_size = 100
alpha, iterations = (0.001, 300)
pixels_per_image, num_labels, hidden_size = (784, 10, 100)

weights_0_1 = 0.2*np.random.random((pixels_per_image,hidden_size)) - 0.1
weights_1_2 = 0.2*np.random.random((hidden_size,num_labels)) - 0.1

for j in range(iterations):
    error, correct_cnt = (0.0, 0)
    for i in range(int(len(images) / batch_size)):
        batch_start, batch_end = ((i * batch_size),((i+1)*batch_size))

        layer_0 = images[batch_start:batch_end]
        layer_1 = relu(np.dot(layer_0,weights_0_1))
        dropout_mask = np.random.randint(2,size=layer_1.shape)
        layer_1 *= dropout_mask * 2
        layer_2 = np.dot(layer_1,weights_1_2)

        error += np.sum((labels[batch_start:batch_end] - layer_2) ** 2)
        for k in range(batch_size):
            correct_cnt += int(np.argmax(layer_2[k:k+1]) == \
                np.argmax(labels[batch_start+k:batch_start+k+1]))

        layer_2_delta = (labels[batch_start:batch_end]-layer_2) \
            /batch_size
        layer_1_delta = layer_2_delta.dot(weights_1_2.T)* \
            relu2deriv(layer_1)
        layer_1_delta *= dropout_mask

        weights_1_2 += alpha * layer_1.T.dot(layer_2_delta)
        weights_0_1 += alpha * layer_0.T.dot(layer_1_delta)

    if(j%10 == 0):
        test_error = 0.0
        test_correct_cnt = 0

        for i in range(len(test_images)):
            layer_0 = test_images[i:i+1]
            layer_1 = relu(np.dot(layer_0,weights_0_1))
            layer_2 = np.dot(layer_1, weights_1_2)
```


Batch Gradient Descent (con't)

The first thing you'll notice when running this code is that it runs way faster. This is because each "np.dot" function is now performing 100 vector dot products at a time. As it turns out, CPU architectures are way faster at performing dot products batched in this way.

There's actually more going on here, however. Notice that our alpha is 20x larger than it was before. We can increase this for a rather fascinating reason. Consider if you were trying to find a city using a very wobbly compass. If you just looked down, got a heading, and then ran 2 miles, you'd likely be way off course! However, if you looked down, took 100 headings and then averaged them, running 2 miles would probably take you in the general right direction. Thus, because we're taking an average of a noisy signal (i.e.,

Conclusion

In this chapter we have addressed two of the most widely used methods for increasing the accuracy and training speed of almost any neural architecture. In the following chapters, we will pivot from sets of tools that are universally applicable to nearly all neural networks to special purpose architectures that are advantageous for modeling specific types of phenomenon in data. See you there!

Modeling Probabilities and Non-Linearities

Activation Functions

9

IN THIS CHAPTER

What is an Activation Function?

Standard Hidden Activation Functions

Sigmoid

Tanh

Standard Output Activation Functions

Softmax

Activation Function "Installation Instructions"

“ I know that two and two make four, and should be glad to prove it too, if I could, though I must say if by any sort of process I could convert two and two into five it would give me much greater pleasure. ”

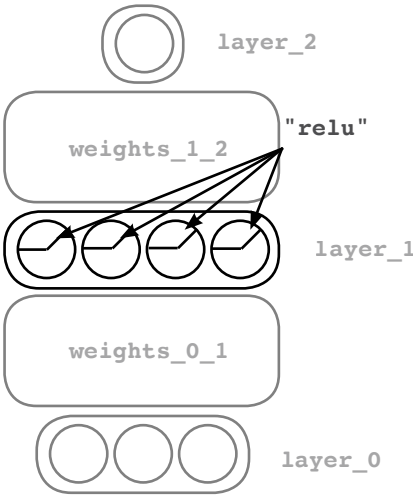
— GEORGE GORDON BYRON

What is an Activation Function?

Definition: a function applied to the neurons in a layer during prediction.

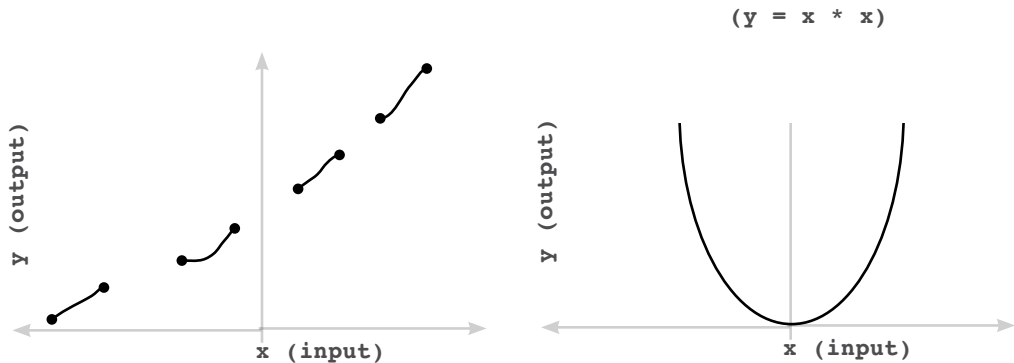
An activation function is a function applied to the neurons in a layer during prediction. This should seem very familiar to you, as we have already been using an activation function called "relu" (pictured right in our 3 layer neural network). The relu function had the affect of turning all negative numbers to 0.

Oversimplified, an activation function is any function that can take one number and return another number. However, there are an infinite number of functions in the universe, and not all of them are useful as activation functions. There are several constraints on what makes a function an "Activation Function". Using functions outside of these constraints is usually a bad idea, as we will see.



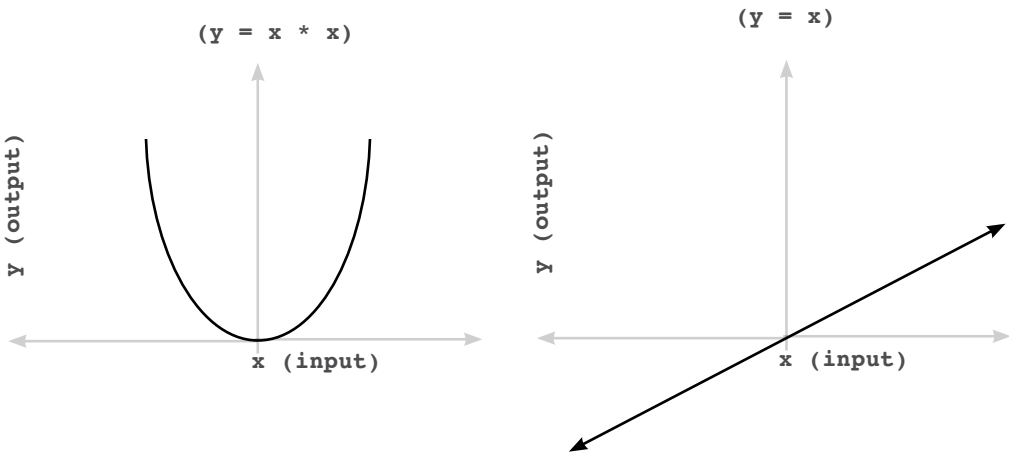
Constraint 1: The function must be continuous and infinite in domain.

The first constraint on what makes a proper "Activation Function" is that it must have an output number for *any* input. In other words, we shouldn't be able to put in a number that doesn't have an output for some reason. A bit overkill, but see how the function on the left (4 distinct lines) doesn't have y values for every x value? It's only defined in 4 spots. This would make for a horrible activation function. The function on the right, however, is continuous and infinite in domain. There is no input (x) for which we can't compute an output(y).



Constraint 2: Good activation functions are "monotonic", never changing direction.

The second constraint is that the function is 1:1. It must never "change direction". In other words, it must either be "always increasing" or "always decreasing". As an example, look at the two functions below. These shapes answer, "Given x as input, what value of y does the function describe?" The function on the left ($y = x * x$) is not an ideal activation function because it is not either "always increasing" or "always decreasing". How can we tell? Well, notice that there are many cases in which two values of x have a single value of y (this is true for every value except 0, actually). The function on the right, however, is always increasing! There is no point at which 2 values of x have the same value of y.



Now, this particular constraint is not technically a "requirement". Unlike functions that have missing values (non-continuous), we can optimize functions that aren't monotonic. However, consider the implication of having multiple input values map to the same output value. When we're learning in our neural networks, we're "searching" for the right weight configurations to give us a specific output. This problem can get a lot harder if we have multiple "right answers". In other words, if there are multiple ways to get the same output, then the network has multiple possible "perfect" configurations.

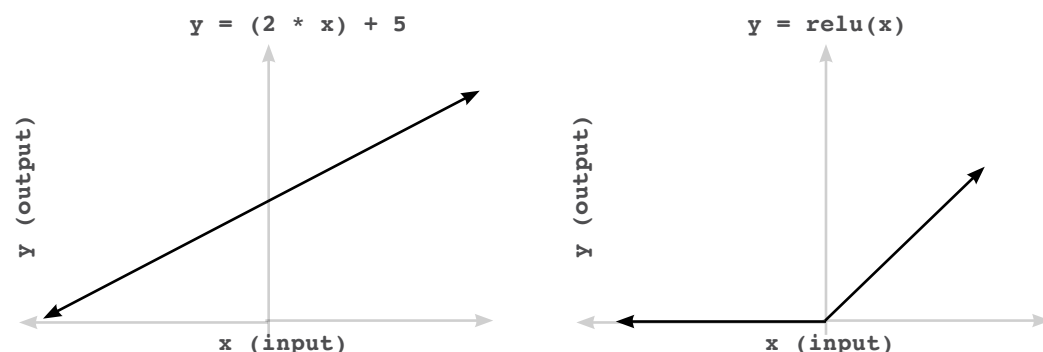
An optimist might say, "Hey, this is great! We're more likely to find the right answer if it can be found in multiple places!" A pessimist would say, "This is terrible! Now we don't have a 'correct direction to go to reduce the error' since we can go in either direction and theoretically make progress." Unfortunately, the phenomenon the pessimist identified is more important. For an advanced study of this subject, look more into "Convex vs Non-Convex Optimization". Many universities (and online courses) will have entire courses dedicated to these kinds of questions.

Constraint 3: Good activation functions are "nonlinear" (i.e., they squiggle or turn).

The third constraint requires a bit of recollection back to our chapter on deep neural networks. Remember how we needed to have "Sometimes Correlation"? In order to create this "Sometimes Correlation", we needed to allow our neurons to "selectively correlate to input neurons" such that a very negative signal from one input into a neuron could reduce how much it correlated to any input at all (by forcing the neuron to simply drop to 0, in the case of "relu").

As it turns out, this phenomenon is facilitated by *any function that curves*. Functions that look like "straight lines", on the other hand, simply scale the weighted average coming in. Simply scaling something (multiplying it by a constant... like "2") doesn't affect how correlated a neuron is to its various inputs. It simply makes the collective correlation that is represented louder or softer. However, the activation doesn't allow one weight to affect how correlated the neuron is to the other weights. What we *really* want is *selective* correlation. Given a neuron with an activation function, we want one incoming signal to be able to increase or decrease how correlated the neuron is to all of the other incoming signals. All curved lines do this (to varying degrees, as we will see).

Thus, the function on the left is considered a "linear" function, whereas the one on the right is considered "non-linear" and will usually make for a better activation function (there are exceptions that we'll discuss later).



Constraint 4: Good activation functions (and their derivatives) should be efficiently computable.

This one is pretty simple. You're going to be calling this function a lot (sometimes billions of times), so you don't want it to be too slow to compute. In fact, many recent activation functions have become popular because they are so easy to compute at the expense of their expressiveness ("relu" is a great example of this).

Standard Hidden Layer Activation Functions

Out of the infinite possible functions, which ones are most commonly used?

Even with these constraints, it should be clear that there are an infinite (possibly transfinite?) number of functions that could be used as activation functions. In fact, the last few years have seen a lot of progress in state-of-the-art activations. However, there are still a relatively small list of activations that account for the vast majority of activation needs, and improvements on them have been relatively minute in most cases.

Sigmoid: The Bread and Butter

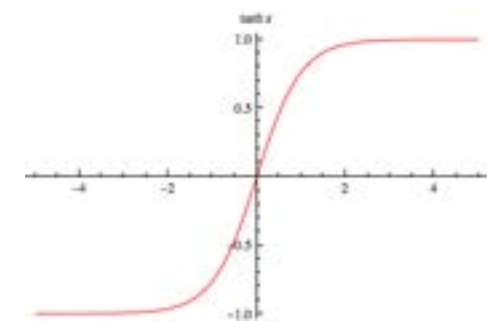
Sigmoid is great because it smoothly squishes the infinite amount of input to an output between 0 and 1. In many circumstances, this lets you interpret the output of any individual neuron as a "probability". Thus, people use this non-linearity both in hidden layers and output layers. (image: wikipedia)



Tanh: Better than Sigmoid for Hidden

Here's the cool thing about Tanh. Remember how we wanted to model "selective correlation?" Well, Sigmoid gives you "varying degrees of positive correlation". That's nice. Tanh is the same as sigmoid *except* it's between -1 and 1! This means it can also throw in some *negative correlation*. While not that useful for output layers (unless the data you're predicting goes between -1 and 1), this aspect of *negative correlation* is very powerful for hidden layers and on many problems, tanh will outperform sigmoid in hidden layers.

(image: wolfram alpha)



Standard Output Layer Activation Functions

Choosing the best one depends on what you're trying to predict.

So, it turns out that what is best for Hidden Layer activation functions can be quite different than what is best for Output Layer activation functions, especially when it comes to classification. Broadly speaking, there are 3 major types of output layer.

Configuration 1: Predicting Raw Data Values (No Activation Function)

This is perhaps the most straightforward but least common type of output layer. In some cases, people want to train a neural network to transform "one matrix of numbers" into "another matrix of numbers", where the range of the output (difference between lowest and highest value) is something other than a probability. One example might be predicting the average temperature in Colorado given the temperature in the surrounding states.

The main thing to focus on here is to ensure that your output non-linearity can actually predict the right answers. In the case above, a sigmoid or tanh would be inappropriate since they force every prediction to be between 0 and 1 (but we want it to predict any temperature, not just between 0 and 1). If I were training a network to do this prediction, I would very likely just train the network without an activation function on the output at all.

Configuration 2: Predicting Unrelated Yes/No Probabilities (Sigmoid)

There will often be times where you want to make multiple binary probabilities in one neural network. We actually already did this in our "Networks with Multiple Inputs and Outputs" subchapter, where we predicted whether the team would win, whether there would be injuries, and the morale of the team (happy or sad) based on the input data. As an aside, when you have neural networks with hidden layers, predicting multiple things at once can be beneficial. Often times the network will learn something in predicting one label that will be useful to one of your other labels. For example, if the network got really good at predicting whether or not the team would win ballgames, the same hidden layer would likely be very useful for predicting whether the team would be happy or sad. However, the network might have a harder time simply predicting happiness or sadness without this extra signal. This tends to vary greatly from problem to problem, but it's good to keep in mind.

In these instances, it's best to use the Sigmoid activation function as it models individual probabilities separately for each output node.

Configuration 3: Predicting "Which One" Probabilities (Softmax)

By far the most common usecase in neural networks is when you want to predict a single label out of many. For example, in our MNIST digit classifier, we want to predict *which* number is in the image. We know ahead of time that the image cannot be more than one number at a time. We could train this network with a sigmoid activation function and simply declare that the highest output probability is the most likely! This will work reasonably well. However, it is far better to have an activation function that models the idea that "The more likely it's one label, the less likely it is any of the other labels."

Why do we like this phenomenon? Well, consider how weight updates are performed. Let's say our MNIST digit classifier should predict that the image is a 9. Let's say that the raw weighted sums going into our final layer (before we apply an activation function) are the following values:



In other words, the network's raw input to the last layer predicts a "0" for every node but 9, where it predicts 100. We might call this "perfect". Let's see what happens when we run these numbers through a sigmoid activation function:



Strangely, the network seems less sure now! I mean, 9 is still the highest, but it seems to think that there's a 50% chance that it could have been any one of the other numbers as well. Weird! Softmax, on the other hand, would interpret the input very differently:



This looks great! Not only is 9 the highest, but the network doesn't even suspect it's any of the other possible MNIST digits! This might seem like just a theoretical flaw of sigmoid, but it can have serious consequences when we backpropagate. Consider how our "Mean Squared Error" is calculated on the *sigmoid* output. In theory, the network is predicting nearly perfectly, right? Surely it won't backprop hardly any error. Not so for sigmoid!



Look at all of the error! These weights are in for a *massive* weight update even though the network predicted perfectly! Why? Well, for sigmoid to reach 0 error, it doesn't just have to predict the highest positive number for the true output, but it has to predict a 0 everywhere else. In other words, where softmax asks, "Which digit seems like the best fit for this input?" Sigmoid says, "You better believe that it's only digit '9' and doesn't have anything in common with the other MNIST digits!!"

The Core Issue: Inputs Have Similarity

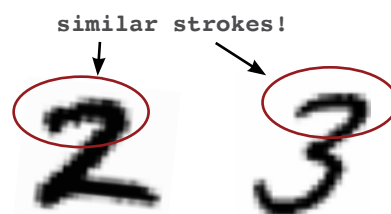
Different numbers share characteristics. It's good to let the network believe that.

So, it turns out that MNIST digits aren't all completely different. In other words, they have overlapping pixel values! The average 2 shares quite a bit in common with the average 3. Why is this important? Well, as a general rule, similar inputs create similar outputs. When you take some numbers and multiply them by a matrix, if the starting numbers are pretty similar, the ending numbers are going to be pretty similar.

Consider the 2 and 3 above. If you forward propagate the 2 and a small amount of probability accidentally goes to the label "3", what does it mean for the network to consider this a *big mistake* and respond with a *big weight update*? It means that it's going to penalize the network for recognizing a 2 by anything other than features that are exclusively related to 2s! It penalizes the network for recognizing a 2 based on, say, the top curve! Why? Well, 2 and 3 share the same curve at the top of the image. Training with a sigmoid would penalize the network for trying to predict a 2 based on this input, because by doing so it would be looking for the same input as it does for 3s. Thus, when a 3 came along, the "2" label would get some probability (because part of the image looks "2ish").

What's the side affect? Well, since most images share lots of pixels in the middle of images, the network will start trying to focus on the edges of the images. Consider this "2 detector" node's weights. See how muddy the middle of the image is? The heaviest weights are the end points of the 2 that are toward the edge of the image. On one hand, these are probably the *best individual* indicators of a 2, but the best overall is going to be a network that sees the entire shape for what it is. These individual indicators can be accidentally triggered by a 3 that is slightly off center or tilted in the wrong way. It's not learning the true essence of a 2 because it needs to learn "2 and NOT 1 NOT 3, NOT 4, etc."

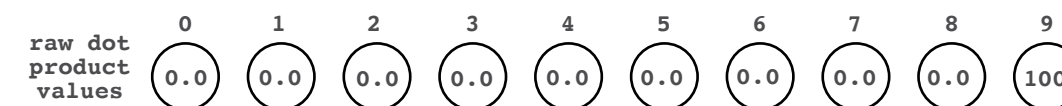
What we really want is an output activation that won't penalize labels that are similar. Instead, we want it to pay attention to all of the information that can be indicative of any potential input. Furthermore, it's also quite nice that a softmax's probabilities always sum to 1. We can interpret any individual prediction as a global probability that the prediction is a particular label. Softmax works better in both theory and practice.



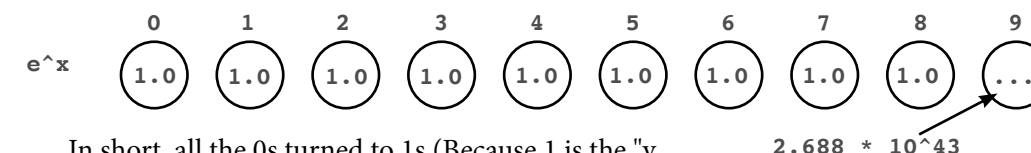
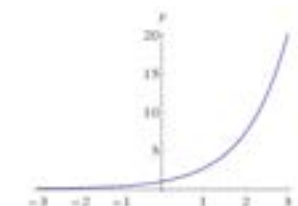
Softmax Computation

Softmax raises each input value exponentially and then divides by the layer's sum.

Let's see a softmax computation on our neural network's hypothetical output values from a few pages ago. I'll picture them here again so we can see the input to softmax:



The first thing that you must do to compute a softmax on the whole layer is to raise each value exponentially. So for each value x , we are going to compute e to the power of x (e is a special number $\sim 2.71828...$). You can see the value of e^x on the right. Notice that it turns every prediction into a positive number, where negative numbers turn into *very small* positive numbers and big numbers turn into *very big* numbers. (If you've heard of "exponential growth" it was likely talking about this function or one very similar to it.)



In short, all the 0s turned to 1s (Because 1 is the "y intercept" of e^x) and the 100 turned into a massive number (2 followed by 43 zeros!). If there were any negative numbers, they would have turned into something between 0 and 1. The next step is to sum all of the nodes in the layer, and divide each value in the layer by that sum. This effectively makes every number 0 except for the value for label 9.



So, what is it about softmax that we like? The nice thing about softmax is that the higher the network predicts one value, the lower it predicts all the others. It increases what is called "sharpness of attenuation". It encourages the network to predict one output with very high probability.

If you wanted to adjust how aggressively it does this, you could simply use numbers slightly higher or lower than "e" when you're exponentiating. Lower numbers will result in lower attenuation, and higher numbers will result in higher attenuation. However, most people just stick with "e".

Activation Installation Instructions

How do you add your favorite activation function to any layer?

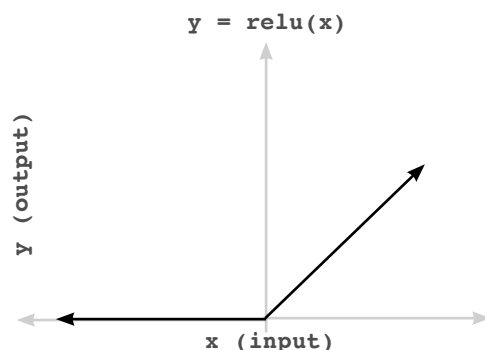
Now that we've covered a wide variety of activation functions and explained their usefulness in hidden and output layers of neural networks, let's talk about the proper way to install one into a neural network. Fortunately, we've already seen an example of how to use a non-linearity in our first deep neural network. In this case, we added a "relu" activation function to our hidden layer. Adding this to forward propagation was relatively straightforward. We just took what layer_1 would have been (without an activation) and applied the "relu" function to each value.

```
layer_0 = images[i:i+1]
layer_1 = relu(np.dot(layer_0, weights_0_1))
layer_2 = np.dot(layer_1, weights_1_2)
```

There's a bit of lingo here to remember. The "input" to a layer refers to the value before the non-linearity. In this case, the input to layer_1 is `np.dot(layer_0, weights_0_1)`. This is not to be confused with the "previous layer" itself, layer_0.

Adding an activation function to a layer in forward propagation is relatively straightforward. However, properly compensating for the activation function in backpropagation is a bit more nuanced. If you remember back to Chapter 6, we performed an interesting operation to create the layer_1_delta variable. Wherever relu had forced a layer_1 value to be 0, we also multiplied the delta by 0. Why did we do this? The reasoning at the time was: "Because a layer_1 value of 0 had no affect on the output prediction, it shouldn't have any impact on the weight update either. It wasn't responsible for the error." This is the extreme form of a more nuanced property. Consider the shape of the "relu" function.

The "slope" of relu for positive numbers is exactly 1. The "slope" of relu for negative numbers is exactly 0. What this really means is that modifying the input to this function (by a tiny amount) will have a 1:1 affect if it was predicting positively, and will have a 0:1 (none) affect if predicting negatively. This "slope" is a measure of how much the output of relu will change given a change in its input. Since the point of the "delta" at this point is to tell earlier layers "make my input higher or lower next time", this delta is very useful! It modifies the delta backpropagated from the following layer to take into account whether this node contributed to the error.



Thus, when we backpropagate, in order to generate layer_1_delta, we multiply the backpropagated delta from layer_2 (`layer_2_delta.dot(weights_1_2.T)`) by the slope of relu *at the point predicted in forward propagation*. For some deltas the slope is 1 (positive numbers) and for others it is 0 (negative numbers).

```
error += np.sum((labels[i:i+1] - layer_2) ** 2)

correct_cnt += int(np.argmax(layer_2) == \
                    np.argmax(labels[i:i+1]))

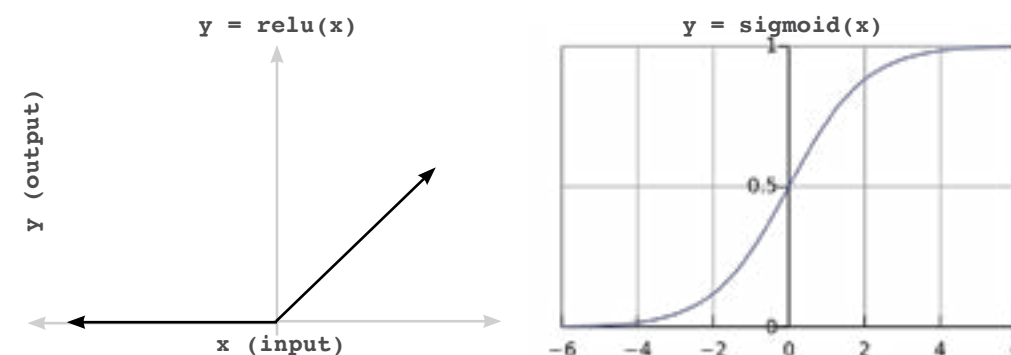
layer_2_delta = (labels[i:i+1] - layer_2)
layer_1_delta = layer_2_delta.dot(weights_1_2.T) \
                * relu2deriv(layer_1)

weights_1_2 += alpha * layer_1.T.dot(layer_2_delta)
weights_0_1 += alpha * layer_0.T.dot(layer_1_delta)

def relu(x):
    return (x >= 0) * x # returns x if x > 0
                        # return 0 otherwise

def relu2deriv(output):
    return output >= 0 # returns 1 for input > 0
                        # return 0 otherwise
```

"relu2deriv" is a special function that can take the output of "relu" and calculate the slope of "relu" at that point (and it does this for all of the values in the "output" vector). This begs the question, how do we make similar adjustments for all of the other non-linearities that aren't "relu"? Consider "relu" and "sigmoid":

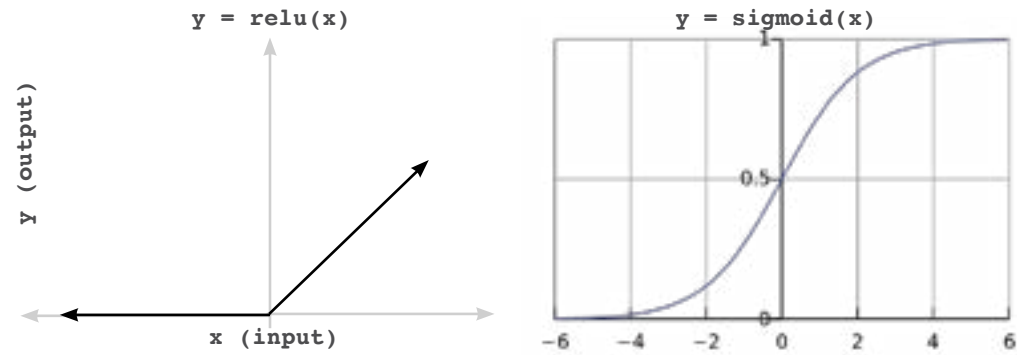
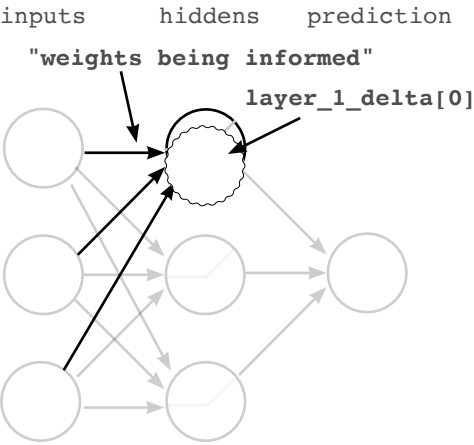


The important point to see in these pictures is that the "slope" is an indicator of how much a *tiny* change to the input would affect the output. We want to modify our incoming delta (from the following layer) to take into account whether or not a weight update before this node would have any affect. Remember, the end goal is to adjust weights to reduce our error. This step encourages the network to leave weights alone if adjusting them will have little to no affect. It does this by multiplying it by the slope. This is no different for sigmoid.

Multiplying Delta By The Slope

To compute `layer_delta`, we multiply the backpropagated delta by the layer's slope.

`layer_1_delta[0]` represents how much higher or lower the first hidden node of layer 1 should be in order to reduce the error of the network (for a particular training example). When there is no non-linearity, this is simply the weighted average delta of `layer_2`. However, the end goal of the delta on a neuron is to inform the weights before it about whether or not they should move. If moving them would have no affect, they (as a group) should be left alone. This is obvious for "relu" which is either "on" or "off". Sigmoid is, perhaps, more nuanced.



Consider a single sigmoid neuron. Sigmoid's "sensitivity to change in the input" slowly increases as the input approaches 0 from either direction. However, very positive and very negative inputs approach a slope of very near 0. Thus, as the input becomes very positive or very negative, small changes to the incoming weights become less relevant to the neuron's error at this training example. In broader terms, many hidden nodes are irrelevant to the accurate prediction of a "2" (perhaps they're only used for "8"s). Thus, we shouldn't mess with their weights too much because we could corrupt their usefulness elsewhere.

Inversely, this also creates a notion of "stickiness". Weights that have previously been updated a lot in one direction (for similar training examples) are very "confidently" predicting a high value or low value. These non-linearities help make it harder for occasional erroneous training examples to corrupt intelligence that has been reinforced many times.

Converting Output to Slope (derivative)

Most great activations can convert their output to their slope. (Efficiency win!)

So, now that we see that adding an activation to a layer changes how we compute the delta for that layer, we should discuss how the industry does this efficiently. The new operation that is necessary is the computation of the derivative of whatever nonlinearity was used.

Most nonlinearities (all of the popular ones) utilize a method of computing a derivative that will seem surprising to those of you who are already familiar with calculus. Instead of computing the derivative at a certain point on its curve the "normal" way, most great activation functions have a means by which the *output* of the layer (at forward propagation) can be used to compute the derivative. This has become the standard practice for computing derivatives in neural networks and it is quite handy. Below is a small table for the functions we've seen so far, paired with their derivatives. `input` is a NumPy vector (corresponding to the input to a layer). `output` is the prediction of the layer. `deriv` is the derivative of the vector of activation derivatives corresponding to the derivative of the activation at each node. `true` is the vector of true values (typically a 1 for the correct label position, 0 everywhere else).

function	forward prop	back prop delta
relu	<pre>ones_and_zeros = (input > 0) output = input*ones_and_zeros</pre>	<pre>mask = output > 0 deriv = output * mask</pre>
sigmoid	<pre>output = 1/(1 + np.exp(-input))</pre>	<pre>deriv = output*(1-output)</pre>
tanh	<pre>output = np.tanh(input)</pre>	<pre>deriv = 1 - (output**2)</pre>
softmax	<pre>temp = np.exp(input) output /= np.sum(temp)</pre>	<pre>temp = (output - true) output = temp/len(true)</pre>

Note that the delta computation for softmax is special because it's only used for the last layer. There's a bit more going on (theoretically) than we have time to discuss here. We'll discuss this more in a later chapter. For now, let's install some better activation functions in our MNIST classification network.

Upgrading our MNIST Network

Let's upgrade our MNIST network to reflect what we've learned.

Theoretically, our tanh function should make for a better hidden layer activation, and our softmax should make for a better output layer activation function. When we test them, they do in fact reach a higher score. However, things are not always as simple as they seem.

I had to make a couple of adjustments in order to have the network "tuned" properly with these new activations. For tanh, I had to reduce the standard deviation of the incoming weights. Remember that we initialize our weights randomly. `np.random.random` simply creates a random matrix with numbers randomly spread between 0 and 1. By multiplying by 0.2 and subtracting by 0.1, we rescale this random range to be between -0.1 and 0.1. This worked great for relu, but is less optimal for tanh. Tanh likes to have a narrower random initialization, so I adjusted it to be between -0.01 and 0.01.

Furthermore, I removed our "error" calculation because we're not ready for that yet. Technically, softmax is best used with an error function called Cross Entropy. This network properly computes the `layer_2_delta` for this error measure, but since we haven't analyzed why this error function is advantageous, I removed the lines to compute it.

Finally, as with almost all changes you make to a neural network, I had to revisit our "alpha" tuning. I found that a much higher alpha was required to reach a good score within 300 iterations. And voilà! As expected, we reached a higher testing accuracy of 87%!

```
import numpy as np, sys
np.random.seed(1)

from keras.datasets import mnist
(x_train, y_train), (x_test, y_test) = mnist.load_data()

images, labels = (x_train[0:1000].reshape(1000,28*28)\
                  / 255, y_train[0:1000])

one_hot_labels = np.zeros((len(labels),10))
for i,l in enumerate(labels):
    one_hot_labels[i][l] = 1
labels = one_hot_labels

test_images = x_test.reshape(len(x_test),28*28) / 255
test_labels = np.zeros((len(y_test),10))
for i,l in enumerate(y_test):
    test_labels[i][l] = 1

def tanh(x):
    return np.tanh(x)
def tanh2deriv(output):
    return 1 - (output ** 2)
def softmax(x):
    temp = np.exp(x)
    return temp / np.sum(temp, axis=1, keepdims=True)
```

```
alpha, iterations, hidden_size = (2, 300, 100)
pixels_per_image, num_labels = (784, 10)
batch_size = 100

weights_0_1 = 0.02*np.random.random((pixels_per_image,hidden_size))-0.01
weights_1_2 = 0.2*np.random.random((hidden_size,num_labels))- 0.1

for j in range(iterations):
    correct_cnt = 0
    for i in range(int(len(images) / batch_size)):
        batch_start, batch_end=((i * batch_size),((i+1)*batch_size))
        layer_0 = images[batch_start:batch_end]
        layer_1 = tanh(np.dot(layer_0,weights_0_1))
        dropout_mask = np.random.randint(2,size=layer_1.shape)
        layer_1 *= dropout_mask * 2
        layer_2 = softmax(np.dot(layer_1,weights_1_2))

        for k in range(batch_size):
            correct_cnt += int(np.argmax(layer_2[k:k+1]) == \
                                     np.argmax(labels[batch_start+k:batch_start+k+1]))
        layer_2_delta = (labels[batch_start:batch_end]-layer_2)\
                        / (batch_size * layer_2.shape[0])
        layer_1_delta = layer_2_delta.dot(weights_1_2.T) \
                        * tanh2deriv(layer_1)

        layer_1_delta *= dropout_mask

        weights_1_2 += alpha * layer_1.T.dot(layer_2_delta)
        weights_0_1 += alpha * layer_0.T.dot(layer_1_delta)
    test_correct_cnt = 0

    for i in range(len(test_images)):

        layer_0 = test_images[i:i+1]
        layer_1 = tanh(np.dot(layer_0,weights_0_1))
        layer_2 = np.dot(layer_1,weights_1_2)
        test_correct_cnt += int(np.argmax(layer_2) == \
                                     np.argmax(test_labels[i:i+1]))

    if(j % 10 == 0):
        sys.stdout.write("\n" + "I:" + str(j) + \
            " Test-Acc:" + str(test_correct_cnt/float(len(test_images))) + \
            " Train-Acc:" + str(correct_cnt/float(len(images))))

I:0 Test-Acc:0.394 Train-Acc:0.156 I:150 Test-Acc:0.8555 Train-Acc:0.914
I:10 Test-Acc:0.6867 Train-Acc:0.723 I:160 Test-Acc:0.8577 Train-Acc:0.925
I:20 Test-Acc:0.7025 Train-Acc:0.732 I:170 Test-Acc:0.8596 Train-Acc:0.918
I:30 Test-Acc:0.734 Train-Acc:0.763 I:180 Test-Acc:0.8619 Train-Acc:0.933
I:40 Test-Acc:0.7663 Train-Acc:0.794 I:190 Test-Acc:0.863 Train-Acc:0.933
I:50 Test-Acc:0.7913 Train-Acc:0.819 I:200 Test-Acc:0.8642 Train-Acc:0.926
I:60 Test-Acc:0.8102 Train-Acc:0.849 I:210 Test-Acc:0.8653 Train-Acc:0.931
I:70 Test-Acc:0.8228 Train-Acc:0.864 I:220 Test-Acc:0.8668 Train-Acc:0.93
I:80 Test-Acc:0.831 Train-Acc:0.867 I:230 Test-Acc:0.8672 Train-Acc:0.937
I:90 Test-Acc:0.8364 Train-Acc:0.885 I:240 Test-Acc:0.8681 Train-Acc:0.938
I:100 Test-Acc:0.8407 Train-Acc:0.88 I:250 Test-Acc:0.8687 Train-Acc:0.937
I:110 Test-Acc:0.845 Train-Acc:0.891 I:260 Test-Acc:0.8684 Train-Acc:0.945
I:120 Test-Acc:0.8481 Train-Acc:0.90 I:270 Test-Acc:0.8703 Train-Acc:0.951
I:130 Test-Acc:0.8505 Train-Acc:0.90 I:280 Test-Acc:0.8699 Train-Acc:0.949
I:140 Test-Acc:0.8526 Train-Acc:0.90 I:290 Test-Acc:0.8701 Train-Acc:0.94
```

IN THIS CHAPTER •

- Natural Language Processing (NLP)
- Supervised NLP
- Capturing Word Correlation in Input Data
- Intro to an Embedding Layer
- Neural Architecture
- Comparing Word Embeddings
- Filling in the Blank
- Meaning is Derived from Loss
- Word Analogies

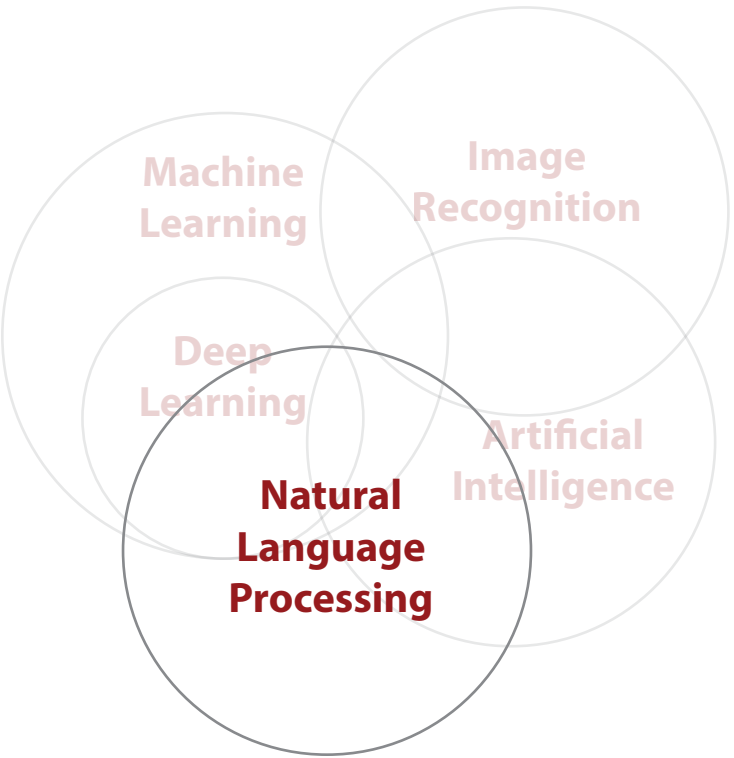
“Computers are incredibly fast, accurate and stupid; humans are incredibly slow, inaccurate and brilliant; together they are powerful beyond imagination.”

— ALBERT EINSTEIN

What does it mean to Understand Language?

What kinds of predictions do people do on language?

Up until now, we've been using neural networks to model image data. However, neural networks can be used to understand a much wider variety of datasets. Exploring new datasets also teaches us a lot about neural networks in general, as different datasets often justify different styles of neural network training according the challenges hidden in the data.



We begin by exploring a much older field that overlaps with Deep Learning called Natural Language Processing (NLP). This field is dedicated exclusively to the automated understanding of human language (previously not using Deep Learning). In this chapter, we're going to discuss the basics of Deep Learning's approach to this field.

Natural Language Processing (NLP)

Natural Language Processing is divided into a collection of tasks or challenges.

Perhaps the best way to quickly get to know Natural Language Processing is by considering a few of the many challenges that the NLP community seeks to solve. Here are a few types of classification problem that are common to NLP.

- Using the **characters** of a document to predict *where words start and end*.
- Using the **words** of a document to predict *where sentences start and end*.
- Using the **words in a sentence** to predict *the part of speech for each word*.
- Using **words in a sentence** to predict *where phrases start and end*.
- Using **words in a sentence** to predict *where named entity (person, place, thing) references start and end*.
- Using **sentences in a document** to predict *which pronouns refer to the same person / place / thing*.
- Using **words in a sentence** to predict the *sentiment* of a sentence.
- and more...

Generally speaking, NLP tasks seek to do one of three things. A task is either labeling a region of text (such as Part-of-Speech Tagging, Sentiment Classification, or Named Entity Recognition), linking two or more regions of text (such as Coreference, which tries to answer whether two mentions of a "real world thing" are in fact referencing the same "real world thing", where "real world thing" is generally a person, place, or some other named entity), or trying to fill in missing information (missing words) based on context.

Perhaps it is also apparent how Machine Learning and Natural Language Processing are deeply intertwined. Until recently, most state-of-the-art NLP algorithms were advanced, probabilistic, non-parametric models (not Deep Learning). However, the recent development and popularization of two major neural algorithms have swept the field of NLP, namely neural word embeddings and recurrent neural networks.

In this chapter, we're going to build a word embedding algorithm and demonstrate why it increases the accuracy of NLP algorithms. In the next chapter, we're doing to create a recurrent neural network and demonstrate why it is so effective at predicting across sequences.

It is also worth mentioning the key role that NLP (perhaps leveraging Deep Learning) plays in the advancement of Artificial Intelligence. A.I. seeks to create machines that can think and engage with the world as humans do (and beyond). NLP plays a very special role in this endeavor, as language is the bedrock of conscious logic and communication in humans. As such, the methods by which machines can leverage and understand language form the foundation of "human-like" logic in machines - the foundation of "thought".

Supervised NLP

Words go in and predictions come out.

Perhaps you will remember the picture below from Chapter 2. Supervised learning is all about taking "what we do know" and transforming it into "what we want to know." Up until now "what we know" has always been comprised of numbers in one way or another. However, NLP leverages text as input. How do we process this?



Well, since neural networks only map input numbers to output numbers, our first step is to convert our text into numerical form. Much like we converted our streetlight dataset before, we need to convert our real world data (in this case text) into a *matrix* that the neural network can consume. As it turns out, how we do this is extremely important!



How should we convert text to numbers? Answering that question requires some thought regarding our problem. Remember, neural networks look for correlation between their input and output layers. Thus, we want to convert our text into numbers in such a way that the correlation between input and output is *most obvious* to the network. This will make for faster training and better generalization.

In order to know what input format makes input/output correlation the most obvious to the network, we need to know what our input/output dataset looks like. So, to explore this topic, we're going to take on the challenge of Topic Classification.

IMDB Movie Reviews Dataset

Predicting whether people post positive or negative reviews.

The IMDB movie reviews dataset is a collection of review -> rating pairs that often look like the following (note: this is an imitation not actually pulled from IMDB).

"This movie was terrible! The plot was dry, the acting unconvincing, and I spilled popcorn on my shirt.

Rating: 1 (stars)

The entire dataset is comprised of around 50,000 of these pairs, where the input reviews are usually a few sentences and the output ratings are between 1 and 5 stars. People consider it a "sentiment dataset" because the stars are very indicative of the overall sentiment of the movie review. However, it should be quite obvious that this "sentiment" dataset might be very different from other sentiment datasets, such as product reviews, or hospital patient reviews.

What we want to do is train a neural network that can leverage the input text to make accurate predictions of the output score. To accomplish this, we must first decide how to turn our input and output datasets into matrices. Interestingly, our output dataset is a number, which perhaps makes it an easier place to start. We'll adjust the range of stars to be between 0 and 1 instead of 1 and 5, so that we can use binary softmax. This is really all we need to do to the output. I'll show an example on the next page.

The input data, however, is a bit trickier. To begin, let's just consider the raw data itself. It's a list of characters. This presents a few problems, not only is the input data text instead of numbers, but it's *variable length* text. So far, our neural networks always take an input of a fixed size. We'll need to overcome this.

So, the raw input simply won't work. The next question to ask yourself is, "What about this data will have correlation with the output?" Simply representing that property might work quite well. For starters, I wouldn't expect any characters (in our list of characters) to have any correlation with the sentiment. We need to think about it differently.

What about the words? There are several words in this dataset that would have quite a bit of correlation! I would bet that "terrible" and "unconvincing" have significant *negative* correlation with the rating. By *negative*, I mean that as they *increase* in frequency in any input datapoint (any review), the rating tends to *decrease*.

Perhaps this property is more general! Perhaps words by themselves (even out of context) would have significant correlation with sentiment. Let's explore this further.

Capturing Word Correlation in Input Data

"Bag of Words": Given a review's vocabulary, predict the sentiment.

If we observe that there is correlation between the vocabulary of an IMDB review and its rating, then we can proceed to the next step, creating an input matrix that represents the vocabulary of a movie review.

What is commonly done in this case is to create a matrix where each row (vector) corresponds to each movie review, and each column represents whether a review contains a particular word in our vocabulary. So, to create the vector for a review, we calculate the vocabulary of the review and then put a "1" in each corresponding column for that review, and "0"s everywhere else. How big are these vectors? Well, if there are 2000 words, and we need to have a place in each vector for each word, then each vector will have 2000 dimensions. This form of storage is called "one-hot encoding" and is the most common format for encoding binary data (the "binary" presence or absence of an input datapoint amongst a vocabulary of possible input datapoints). So, if our vocabulary was only 4 words, our one-hot encoding might look like the following.

```
import numpy as np

onehots = {}
onehots['cat'] = np.array([1,0,0,0])
onehots['the'] = np.array([0,1,0,0])
onehots['dog'] = np.array([0,0,1,0])
onehots['sat'] = np.array([0,0,0,1])

sentence = ['the','cat','sat']
x = word2hot[sentence[0]] + \
    word2hot[sentence[1]] + \
    word2hot[sentence[2]]

print("Sent Encoding:" + str(x))
```

cat

the

dog

sat

As you can see, we create a vector for each term in the vocabulary and this allows us to use simple vector addition to create a vector representing a subset of the total vocabulary (such as a subset corresponding to the words in a sentence).

```
Code Output: Sent Encoding:[1 1 0 1]
```

"the cat sat"

Note that when we create an embedding for several (such as "the cat sat"), we have multiple options in the case that words occur multiple times. If our phrase was "cat cat cat", we could either sum the vector for "cat" three times (resulting in [3,0,0,0]) or just take the unique "cat" a single time (resulting in [1,0,0,0]). The latter typically works better for language.

Predicting Movie Reviews

With our encoding strategy and our previous network, we can predict sentiment.

So, using the strategy identified on the previous page, we can build a vector for each word in our sentiment dataset and use our previous 2 layer network to predict sentiment. While I'm going to give you the code below, I strongly recommend attempting this from memory. Open up a new Jupyter notebook, load in the dataset, build your one-hot vectors, and then build a neural network to predict the rating of each movie review (positive or negative). Below is how I would do the pre-processing step.

```
import sys

f = open('reviews.txt')
raw_reviews = f.readlines()
f.close()

f = open('labels.txt')
raw_labels = f.readlines()
f.close()

tokens = list(map(lambda x:set(x.split(" ")),raw_reviews))

vocab = set()
for sent in tokens:
    for word in sent:
        if(len(word)>0):
            vocab.add(word)
vocab = list(vocab)

word2index = {}
for i,word in enumerate(vocab):
    word2index[word]=i

input_dataset = list()
for sent in tokens:
    sent_indices = list()
    for word in sent:
        try:
            sent_indices.append(word2index[word])
        except:
            ""
    input_dataset.append(list(set(sent_indices)))

target_dataset = list()
for label in raw_labels:
    if label == 'positive\n':
        target_dataset.append(1)
    else:
        target_dataset.append(0)
```

Intro to an Embedding Layer

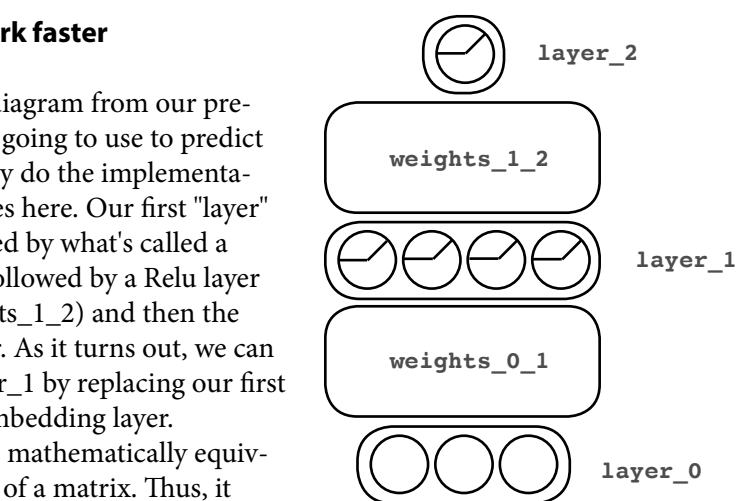
One more trick to make our network faster

On the right I've written the diagram from our previous neural network that we are now going to use to predict sentiment. However, before we actually do the implementation, I want to describe the layer names here. Our first "layer" is our dataset (layer_0). This is followed by what's called a "Linear" layer (weights_0_1). This is followed by a Relu layer (layer_1), another Linear layer (weights_1_2) and then the output, which is our "prediction" layer. As it turns out, we can actually take a bit of a shortcut to layer_1 by replacing our first Linear layer (weights_0_1) with an Embedding layer.

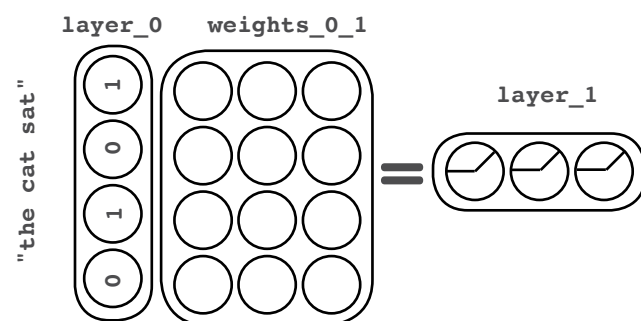
Taking a vector of 1s and 0s is mathematically equivalent to simply summing several rows of a matrix. Thus, it is actually much more efficient to simply select the relevant rows of `weights_0_1` and sum

them as opposed to doing a big vector-matrix multiplication. Since our Sentiment vocabulary is on the order of 70,000 words, most of the vector-matrix multiplication is spent multiplying 0s in the input vector by different rows of the matrix before summing them. Simply selecting the rows corresponding to each word in a matrix and summing them is much more efficient.

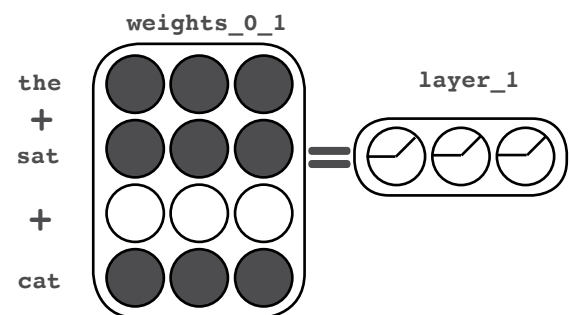
Using this process of selecting rows and performing a sum (or average) means that we're treating our first Linear layer (weights_0_1) as an Embedding layer. Structurally, they are identical (layer_1 is exactly the same using either method for forward propagation). The only difference is that summing a small number of rows is much faster.



One-Hot Vector-Matrix Multiplication



Matrix Row Sum



Predicting Movie Reviews

After running the code from two pages ago, run this code.

```
import numpy as np
np.random.seed(1)

def sigmoid(x):
    return 1/(1 + np.exp(-x))

alpha, iterations = (0.01, 2)
hidden_size = 100

weights_0_1 = 0.2*np.random.random((len(vocab),hidden_size)) - 0.1
weights_1_2 = 0.2*np.random.random((hidden_size,1)) - 0.1

correct,total = (0,0)
for iter in range(iterations):

    # train on first 24,000
    for i in range(len(input_dataset)-1000):

        x,y = (input_dataset[i],target_dataset[i])
        layer_1 = sigmoid(np.sum(weights_0_1[x],axis=0)) #embed + sigmoid
        layer_2 = sigmoid(np.dot(layer_1,weights_1_2)) # linear + softmax

        layer_2_delta = layer_2 - y # compare pred with truth
        layer_1_delta = layer_2_delta.dot(weights_1_2.T) #backprop

        weights_0_1[x] -= layer_1_delta * alpha
        weights_1_2 -= np.outer(layer_1,layer_2_delta) * alpha

        if(np.abs(layer_2_delta) < 0.5):
            correct += 1
            total += 1
            if(i % 10 == 9):
                progress = str(i/float(len(input_dataset)))
                sys.stdout.write('\rIter:' + str(iter) + '\n'
                                + 'Progress:' + progress[2:4] + '\n'
                                + '.' + progress[4:6] + '\n'
                                + '% Training Accuracy:' + '\n'
                                + str(correct/float(total)) + '%')

    print()
    correct,total = (0,0)
    for i in range(len(input_dataset)-1000,len(input_dataset)):

        x = input_dataset[i]
        y = target_dataset[i]

        layer_1 = sigmoid(np.sum(weights_0_1[x],axis=0))
        layer_2 = sigmoid(np.dot(layer_1,weights_1_2))

        if(np.abs(layer_2 - y) < 0.5):
            correct += 1
            total += 1
    print("Test Accuracy:" + str(correct / float(total)))
```

Interpreting the Output

What did our neural network learn along the way?

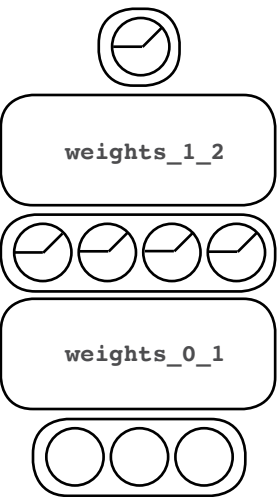
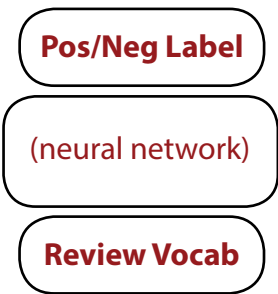
Below, you can see the output of our Movie Reviews neural network. From one perspective, this is simply the same Correlation Summarization that we have already discussed.

```
Iter:0 Progress:95.99% Training Accuracy:0.832%
Iter:1 Progress:95.99% Training Accuracy:0.8663333333333333%
Test Accuracy:0.849
```

From this perspective, the neural network on the previous page was simply looking for correlation between the input datapoints and the output datapoints. However, those datapoints have characteristics that we are quite familiar with (notably those of language). Furthermore, it is extremely beneficial to consider what patterns of language would be detected by the Correlation Summarization, and more importantly, which ones would not. After all, just because the network is able to find correlation between our input and output datasets does not mean that it understands every useful pattern of language. Furthermore, understanding the difference between what the network (in its current configuration) is capable of learning relative to what it needs to know to properly understand language is an incredibly fruitful line of thinking. This is what researchers on the front lines of state-of-the-art research consider, and it's what we're going to consider today.

So, what about language did our movie reviews network actually learn? Well, let's first start by considering what we presented to the network. As displayed in the diagram on the top right, we presented each review's vocabulary as input, and asked it to predict one of two labels (positive or negative). So, given that our Correlation Summarization says the network will look for correlation between our input and output datasets, at a very minimum, we would expect our network to identify words that have either a positive or negative correlation (by themselves).

This follows naturally from the Correlation Summarization itself. We present the presence or absence of a word. As such, the Correlation Summarization will find direct correlation between this presence/absence and each of our two labels. However, this isn't the whole story.



Neural Architecture

How did our choice of architecture affect what the network learned?

On the last page, we discussed the first, most trivial type of information that our neural network learned: direct correlation between our input and target datasets. This observation is largely the "clean slate" of neural intelligence. (If your network cannot discover direct correlation between input and output data, something is probably broken.) The development of more sophisticated architectures is based on the need to find more complex patterns than direct correlation, and this network is no exception.

The minimal architecture needed to identify direct correlation is a 2 layer network, where a network has a single weight matrix which connects directly from the input layer to the output layer. However, we used a network that has a hidden layer. This begs the question, what does this hidden layer do?

Fundamentally, hidden layers are about grouping datapoints from a previous layer into "n" groups (where "n" is the number of neurons in the hidden layer). Each hidden neuron takes in a datapoint and answers the question: "Is this datapoint in my group?" As the hidden layer learns, it searches for useful groupings of its input. What are useful groupings?

An input datapoint "grouping" is useful if it does two things. First and foremost, the grouping must be useful to the prediction of an output label. If it's not useful to the output prediction, the Correlation Summarization will never lead the network to find the group. This is a hugely valuable realization. Much of neural network research is about finding training data (or some other manufactured "signal" for the network to artificially predict) so that it finds groupings that are useful for a task (such as predicting movie review stars). We'll discuss this more in a moment.

Secondly, a "grouping" is useful if it is an actual phenomenon in the data that we care about. Bad groupings just memorize the data. Good groupings pick up on phenomena that are useful linguistically. For example, when predicting whether a movie review is positive or negative, understanding the difference between "terrible" and "NOT terrible" is a powerful grouping. We would love to have a neuron that turned OFF when it saw "awful" and turned ON when it saw "not awful". This would be a powerful grouping for the next layer to use to make the final prediction. However, as the input to our neural network is simply the vocabulary of a review, "it was great, not terrible" creates exactly the same "layer_1" value as "it was terrible, not great". For this reason, we know that our network is very unlikely to create a hidden neuron that understands negation. In fact, this means that testing whether a layer is the same or different based on a certain language pattern is a great first step for knowing whether an architecture is likely to find said pattern using the correlation summarization. If you can construct two examples with an identical hidden layer, one with the pattern you find interesting and one without, the network is unlikely to find that pattern.

Neural Architecture (cont.)

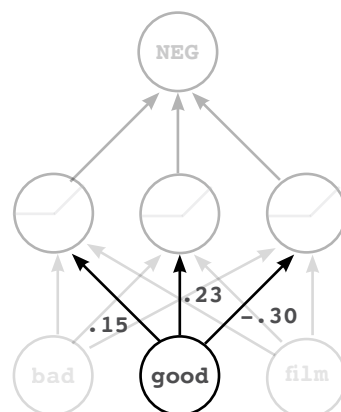
How did our choice of architecture affect what the network learned?

As we learned on the last page, a hidden layer fundamentally groups the previous layer's data. At a granular level, each neuron classifies a datapoint as either subscribing or not subscribing to its group. At a higher level, two datapoints (i.e., "movie reviews") are similar if they subscribe to many of the same groups. Finally, two inputs (i.e., words) are similar if the weights linking them to various hidden neurons (a measure of each word's "group affinity") are similar. So, given this knowledge, in our previous neural network, what should we observe in the weights going into the hidden neurons from the words?

What should we see in the weights connecting words and hidden neurons?

Here's a hint. Words that have a similar predictive power should subscribe to similar groups (hidden neuron configurations). So, what does this mean for the weights connecting each word to each hidden neuron?

Here's the answer. Words that correlate with similar labels (positive or negative) will have similar weights connecting them to various hidden neurons. This is because the neural network learns to bucket them into similar hidden neurons so that the final layer (weights_1_2) can make the correct positive or negative predictions. We can see this phenomenon by taking a particularly positive or negative word and searching for the other words with the most similar weight values. In other words, we can take each word and see which other words have the most similar "weight values" connecting them to each hidden neuron (to each "group"). Words that subscribe to similar groups will have similar predictive power for positive or negative labels. As such, words that subscribe to similar groups, having similar weight values, will also have similar meaning. Abstractly, in terms of neural networks, a neuron has similar meaning to other neurons in the same layer if and only if it has similar weights connecting it to the next and/or previous layers.



The 3 bold weights for "good" form the "embedding" for good. They reflect how much the term "good" is a member of each group (i.e. hidden neuron). Words with similar predictive power will have similar word embeddings (weight values).

Comparing Word Embeddings

How can we visualize weight similarity?

For each input word, we can select the list of weights proceeding out of it to the various hidden neurons by simply selecting the corresponding row of weights_0_1. Each entry in the row represents each weight proceeding from that row's word to each hidden neuron. Thus, to figure out which words are most similar to a target term, we simply compare each word's vector (row of the matrix) to that of the target term. Our comparison of choice is called "Euclidian Distance" which we can perform as seen in the code below.

```
from collections import Counter
import math

def similar(target='beautiful'):
    target_index = word2index[target]
    scores = Counter()
    for word, index in word2index.items():
        raw_difference = weights_0_1[index] - (weights_0_1[target_index])
        squared_difference = raw_difference * raw_difference
        scores[word] = -math.sqrt(sum(squared_difference))

    return scores.most_common(10)
```

This allows us to easily query for the most similar word (neuron) according to the network.

<code>print(similar('beautiful'))</code>	<code>print(similar('terrible'))</code>
<pre>[('beautiful', -0.0), ('atmosphere', -0.70542101298), ('heart', -0.7339429768542354), ('tight', -0.7470388145765346), ('fascinating', -0.7549291974), ('expecting', -0.759886970744), ('beautifully', -0.7603669338), ('awesome', -0.76647368382398), ('masterpiece', -0.7708280057), ('outstanding', -0.7740642167)]</pre>	<pre>[('terrible', -0.0), ('dull', -0.760788602671491), ('lacks', -0.76706470275372), ('boring', -0.7682894961694), ('disappointing', -0.768657), ('annoying', -0.78786389931), ('poor', -0.825784172378292), ('horrible', -0.83154121717), ('laughable', -0.8340279599), ('badly', -0.84165373783678)]</pre>

As you might expect, the most similar term to every word is simply itself, followed by words which had similar usefulness as the target term. Again, as you might expect, since the network really only has two labels (positive/negative), the input terms are grouped according to which label they tend to predict. This is a standard phenomenon of the correlation summarization. It seeks to create similar representations (layer_1 values) within the network based on the label being predicted so that it can predict the right label. In this case, the side effect is that the weights feeding into layer_1 get grouped according to output label. The key takeaway is a gut instinct on this phenomenon of the correlation summarization. It consistently attempts to convince the hidden layers to be similar based on which label should be predicted.

What is the Meaning of a Neuron?

Meaning is entirely based on the target labels being predicted.

Note that the meanings of different words didn't totally reflect how we might group them. Notice that the most similar term to "beautiful" is "atmosphere". This is actually a very valuable lesson. For the purposes of predicting whether or not a movie review is positive or negative, these words have nearly identical meaning. However, in the real world, their meaning is quite different (one is an adjective and another a noun, for example).

```
print(similar('beautiful'))
[('beautiful', -0.0),
 ('atmosphere', -0.70542101298),
 ('heart', -0.7339429768542354),
 ('tight', -0.7470388145765346),
 ('fascinating', -0.7549291974),
 ('expecting', -0.759886970744),
 ('beautifully', -0.7603669338),
 ('awesome', -0.76647368382398),
 ('masterpiece', -0.7708280057),
 ('outstanding', -0.7740642167)]

print(similar('terrible'))
[('terrible', -0.0),
 ('dull', -0.760788602671491),
 ('lacks', -0.76706470275372),
 ('boring', -0.7682894961694),
 ('disappointing', -0.768657),
 ('annoying', -0.78786389931),
 ('poor', -0.825784172378292),
 ('horrible', -0.83154121717),
 ('laughable', -0.8340279599),
 ('badly', -0.84165373783678)]
```

This realization is incredibly important. The "meaning" (of a neuron) in our network is defined based on our target labels. Everything in the neural network is contextualized based on the correlation summarization trying to correctly make predictions. Thus, even though you and I know a great deal about these words, the neural network is entirely ignorant of all information outside of the task at hand. So, how can we convince our network to learn more nuanced information about neurons (in this case, word neurons)? Well, if we give it input and target data that requires a more nuanced understanding of language, it will have reason to learn more nuanced interpretations of various terms. So, what should we use our neural network to predict so that it learns more interesting weight values for our word neurons?

Mary had a little lamb whose ???? was white as snow .

The task that we're going to use to learn more interesting weight values for our word neurons is a glorified "fill in the blank" task. Why are we going to use this? Well, first off, we have nearly infinite training data (the internet), which means we have nearly infinite "signal" for the neural network to use to learn more nuanced information about words. Furthermore, being able to accurately fill in the blank requires at least some notion of context about the real world. For example, in the blank above, is it more likely that the blank is correctly filled by the word "anvil" or "wool"? Let's see if our neural network can figure it out.

Filling in The Blank

Learning richer meanings for words by having a richer signal to learn.

In this next neural network, we're going to use almost exactly the same network as the previous example with only a few modifications. First, instead of predicting a single label given a movie review, we're going to take each (5 word) phrase, remove one word (a focus term), and attempt to train a network to figure out the identity of the word we removed using the rest of the phrase. Secondly, we're going to use a trick called "Negative Sampling" to make our network train a bit faster. Consider that in order to predict which term is missing, we have to have one label for each possible word. This means that we have several thousand labels, which would cause our network to train quite slowly. To overcome this, we're going to randomly ignore most of our labels for each forward propagation step (as in, pretend they don't exist at all). While this might seem like a crude approximation, it's a technique that works quite well in practice. Below is our pre-processing code for this example.

```
import sys, random, math
from collections import Counter
import numpy as np

np.random.seed(1)
random.seed(1)
f = open('reviews.txt')
raw_reviews = f.readlines()
f.close()

tokens = list(map(lambda x: (x.split(" ")), raw_reviews))
wordcnt = Counter()
for sent in tokens:
    for word in sent:
        wordcnt[word] += 1
vocab = list(set(map(lambda x: x[0], wordcnt.most_common()))))

word2index = {}
for i, word in enumerate(vocab):
    word2index[word] = i

concatenated = list()
input_dataset = list()
for sent in tokens:
    sent_indices = list()
    for word in sent:
        try:
            sent_indices.append(word2index[word])
            concatenated.append(word2index[word])
        except:
            ""

    input_dataset.append(sent_indices)
concatenated = np.array(concatenated)
random.shuffle(input_dataset)
```

Filling in The Blank (con't)

Learning richer meanings for words by having a richer signal to learn.

```
alpha, iterations = (0.05, 2)
hidden_size, window, negative = (50, 2, 5)

weights_0_1 = (np.random.rand(len(vocab), hidden_size) - 0.5) * 0.2
weights_1_2 = np.random.rand(len(vocab), hidden_size) * 0

layer_2_target = np.zeros(negative + 1)
layer_2_target[0] = 1

def similar(target='beautiful'):
    target_index = word2index[target]

    scores = Counter()
    for word, index in word2index.items():
        raw_difference = weights_0_1[index] - (weights_0_1[target_index])
        squared_difference = raw_difference * raw_difference
        scores[word] = -math.sqrt(sum(squared_difference))
    return scores.most_common(10)

def sigmoid(x):
    return 1 / (1 + np.exp(-x))

for rev_i, review in enumerate(input_dataset * iterations):
    for target_i in range(len(review)):

        # since it's really expensive to predict every vocabulary
        # we're only going to predict a random subset
        target_samples = [review[target_i]] + list(concatenated\
            [(np.random.rand(negative) * len(concatenated)).astype('int').tolist()])

        left_context = review[max(0, target_i - window):target_i]
        right_context = review[target_i + 1:min(len(review), target_i + window)]

        layer_1 = np.mean(weights_0_1[left_context + right_context], axis=0)
        layer_2 = sigmoid(layer_1.dot(weights_1_2[target_samples].T))
        layer_2_delta = layer_2 - layer_2_target
        layer_1_delta = layer_2_delta.dot(weights_1_2[target_samples])

        weights_0_1[left_context + right_context] -= layer_1_delta * alpha
        weights_1_2[target_samples] -= np.outer(layer_2_delta, layer_1) * alpha

    if (rev_i % 250 == 0):
        sys.stdout.write('\rProgress: ' + str(rev_i / float(len(input_dataset)
            * iterations)) + " " + str(similar('terrible')))
        sys.stdout.write('\rProgress: ' + str(rev_i / float(len(input_dataset)
            * iterations)))
        print(similar('terrible'))

Progress:0.99998 [('terrible', -0.0), ('horrible', -2.846300248788519),
('brilliant', -3.039932544396419), ('pathetic', -3.4868595532695967),
('superb', -3.6092947961276645), ('phenomenal', -3.660172529098085),
('masterful', -3.6856112636664564), ('marvelous', -3.9306620801551664),
```

Meaning is Derived from Loss

Learning richer meanings for words by having a richer signal to learn.

With this new neural network, we can subjectively see that our word embeddings cluster in a rather different way. Where before it clustered words according to their likelihood to predict a Positive or Negative label, now they cluster based on their likelihood to occur within the same phrase (sometimes regardless of sentiment). See for yourself!

Predicting POS/NEG

```
print(similar('terrible'))

[('terrible', -0.0),
 ('dull', -0.760788602671491),
 ('lacks', -0.76706470275372),
 ('boring', -0.7682894961694),
 ('disappointing', -0.768657),
 ('annoying', -0.78786389931),
 ('poor', -0.825784172378292),
 ('horrible', -0.83154121717),
 ('laughable', -0.8340279599),
 ('badly', -0.84165373783678)]
```

Fill In The Blank

```
print(similar('terrible'))

[('terrible', -0.0),
 ('horrible', -2.79600898781),
 ('brilliant', -3.3336178881),
 ('pathetic', -3.49393193646),
 ('phenomenal', -3.773268963),
 ('masterful', -3.8376122586),
 ('superb', -3.9043150978490),
 ('bad', -3.9141673639585237),
 ('marvelous', -4.0470804427),
 ('dire', -4.178749691835959)]
```

Predicting POS/NEG

```
print(similar('beautiful'))

[('beautiful', -0.0),
 ('atmosphere', -0.70542101298),
 ('heart', -0.7339429768542354),
 ('tight', -0.7470388145765346),
 ('fascinating', -0.7549291974),
 ('expecting', -0.759886970744),
 ('beautifully', -0.7603669338),
 ('awesome', -0.76647368382398),
 ('masterpiece', -0.7708280057),
 ('outstanding', -0.7740642167)]
```

Fill In The Blank

```
print(similar('beautiful'))

[('beautiful', -0.0),
 ('lovely', -3.0145597243116),
 ('creepy', -3.1975363066322),
 ('fantastic', -3.2551041418),
 ('glamorous', -3.3050812101),
 ('spooky', -3.4881261617587),
 ('cute', -3.592955888181448),
 ('nightmarish', -3.60063813),
 ('heartwarming', -3.6348147),
 ('phenomenal', -3.645669007)]
```

The key takeaway here is that, even though we were training over the same dataset with a very similar architecture (3 layers, cross entropy, sigmoid nonlinear), we can influence what the network learns within its weights by changing what we tell the network to predict. Thus, even though it's "looking at" the same statistical information, we can target what it learns based on what we select as the input and target values. For the moment, let's call this process of choosing what we want our network to learn "intelligence targeting".

Controlling the input/target values is not the only way we can perform "intelligence targeting". We can also adjust how our network measures error, the size and types of layers that it has, and the types of regularization we apply. In Deep Learning research, all of these techniques fall under the umbrella of constructing what's called a **loss function**.

Meaning is Derived from Loss (con't)

Neural networks don't really learn data, they minimize the loss function

In Chapter 4, we learned that learning is really about adjusting each weight in our neural network to bring our error down to 0. On this page, we're going to explain the exact same phenomena from a different perspective, choosing our error so that our neural network learns the patterns we're interested in. Remember these lessons from Chapter 4?

The Golden Method for Learning	The Secret
Adjusting each <code>weight</code> in the correct <i>direction</i> and by the correct <i>amount</i> so that our <code>error</code> reduces to 0.	For any <code>input</code> and <code>goal_pred</code> , there is an <i>exact relationship</i> defined between our <code>error</code> and <code>weight</code> , found by combining our <code>prediction</code> and <code>error</code> formulas.
<pre>error = ((0.5 * weight) - 0.8) ** 2</pre>	

Perhaps you remember the formula above from our 1 weight neural network. In that network, we could evaluate our error by first forward propagating $(0.5 * \text{weight})$ and then comparing to our target (0.8). I would encourage you to think about this not from the perspective of two different steps (forward propagation, then error evaluation), but to instead consider the entire formula (including forward prop) to be the evaluation of an error value. This context will reveal the true cause of our different word embedding clusterings. Even though the network and datasets were similar, the error function was fundamentally different, leading to different word clusterings within each network.

Predicting POS/NEG	Fill In The Blank
<pre>print(similar('terrible'))</pre> <pre>[('terrible', -0.0), ('dull', -0.760788602671491), ('lacks', -0.76706470275372), ('boring', -0.7682894961694), ('disappointing', -0.768657), ('annoying', -0.78786389931), ('poor', -0.825784172378292), ('horrible', -0.83154121717), ('laughable', -0.8340279599),</pre>	<pre>print(similar('terrible'))</pre> <pre>[('terrible', -0.0), ('horrible', -2.79600898781), ('brilliant', -3.3336178881), ('pathetic', -3.49393193646), ('phenomenal', -3.773268963), ('masterful', -3.8376122586), ('superb', -3.9043150978490), ('bad', -3.9141673639585237), ('marvelous', -4.0470804427),</pre>

Meaning is Derived from Loss (cont.)

Our choice of loss function determines our neural network's knowledge.

The more formal term for an **error function** is a **loss function** or **objective function** (all three phrases are interchangeable). Considering learning to be all about minimizing a loss function (which includes our forward propagation) gives us a far broader perspective on how neural networks learn. We can have two neural networks with identical starting weights, trained over identical datasets that ultimately learn very different patterns because we choose a different loss function. In the case of our two movie review neural networks, our loss function was different because we chose two different target values (POS/NEG vs Fill In the Blank).

Different kinds of architectures, layers, regularization techniques, datasets, and non-linearities aren't really that different. These are simply the different ways we can choose to construct a loss function. Furthermore, if our network isn't learning properly, the solution can often come from any of these possible categories.

For example, if your network is overfitting, you can augment your loss function by choosing simpler nonlinearities, smaller layer sizes, shallower architectures, larger datasets, or more aggressive regularization techniques. All of these choices will have a fundamentally similar affect on your loss function and a similar consequence on the behavior of your network. They all interplay together, and over time you will learn how changing one can affect the performance of another, but for now, the important takeaway is that learning is about constructing a loss function and then minimizing it.

So, whenever you want a neural network to learn a pattern, everything you need to know to do so will be contained in your loss function. When we only had a single weight, this allowed our loss function to be quite simple, as you remember:

```
error = ((0.5 * weight) - 0.8) ** 2
```

However, as you chain large numbers of complex layers together, your loss function will become more complicated (and that's ok). Just remember, if something is going wrong, the solution is in your loss function, which includes both your forward prediction and your raw error evaluation (such as mean squared error or cross entropy).

King - Man + Woman ~ = Queen

Word analogies are an interesting consequence of the previously built network.

Before closing out this chapter, we should discuss what is, at the time of writing, still one of the most famous properties of neural word embeddings (word vectors like those we just created). The task of "filling in the blank" creates word embeddings with some interesting phenomena known as "Word Analogies", wherein you can take the vectors for different words and perform basic algebraic operations on them. For example, if you train the previous network on a large enough corpus, you'll be able to take the vector for "king", subtract from it the vector for "man" add in the vector for "woman" and then search for the most similar vector (other than those in your query). As it turns out, the most similar vector will often be the word "queen". We can even see similar phenomena in the "Fill In the Blank" network trained over movie reviews.

```
def analogy(positive=['terrible','good'],negative=['bad']):  
    norms = np.sum(weights_0_1 * weights_0_1,axis=1)  
    norms.resize(norms.shape[0],1)  
    normed_weights = weights_0_1 * norms  
    query_vect = np.zeros(len(weights_0_1[0]))  
    for word in positive:  
        query_vect += normed_weights[word2index[word]]  
    for word in negative:  
        query_vect -= normed_weights[word2index[word]]  
    scores = Counter()  
    for word,index in word2index.items():  
        raw_difference = weights_0_1[index] - query_vect  
        squared_difference = raw_difference * raw_difference  
        scores[word] = -math.sqrt(sum(squared_difference))  
    return scores.most_common(10)[1:]
```

terrible - bad + good ~ =	elizabeth - she + he ~ =
analogy(['terrible','good'],['bad'])	analogy(['elizabeth','he'],['she'])
<div>[('superb', -223.3926217861), ('terrific', -223.690648739), ('decent', -223.7045545791), ('fine', -223.9233021831882), ('worth', -224.03031703075), ('perfect', -224.125194533), ('brilliant', -224.2138041), ('nice', -224.244182032763), ('great', -224.29115420564)]</div>	<div>[('christopher', -192.7003), ('it', -193.3250398279812), ('him', -193.459063887477), ('this', -193.59240614759), ('william', -193.63049856), ('mr', -193.6426152274126), ('bruce', -193.6689279548), ('fred', -193.69940566948), ('there', -193.7189421836)]</div>

Word Analogies

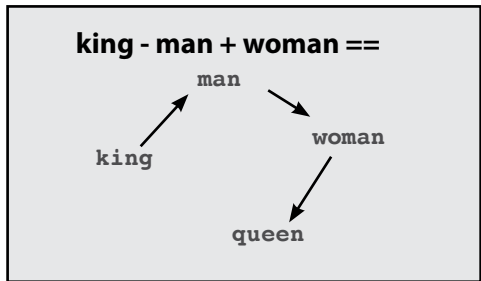
Linear compression of an existing property in the data.

When this property was first discovered, it created a rather large flurry of excitement as people extrapolated many possible applications of such a technology. It is a rather amazing property in its own right, and it did create a veritable cottage industry around generating word embeddings of one variety or another. However, the word analogy property in and of itself hasn't grown that much since then, and most of the current work in language focuses instead on recurrent architectures (which we'll get to in a later chapter).

That being said, getting a good intuition for what's going on with word embeddings as a result of your chosen loss function is extremely valuable. We've already learned that our choice of loss function can affect how words are grouped, but this word analogy phenomenon is something different. What about our new loss function causes it to happen?

If you consider a word embedding having 2 dimensions, it's perhaps easier to envision exactly what it means for these word analogies to work.

```
king = [0.6 , 0.1]  
man = [0.5 , 0.0]  
woman = [0.0 , 0.8]  
queen = [0.1 , 1.0]  
  
king - man = [0.1 , 0.1]  
queen - woman = [0.1 , 0.2]
```



What this means it that the relative usefulness to the final prediction between king/man and queen/woman is similar. Why? The difference between king and man leaves a vector of "royalty". So, there's a bunch of male and female related words in one grouping, and then there's another grouping of "royal words" that's grouped in the "royal" direction.

This can be traced back to the loss that we chose. When the word "king" shows up in a phrase, it changes the probability of other words showing up in a certain way. It increases the probability of words related to "man" and the probability of words related to "royalty". "queen" showing up in a phrase increases the probability of words related to "woman" and the probability of words related to "royalty" (as a group). Thus, because they have this sort of ven-diagram of impact on the output probability, they end up subscribing to similar combinations of groupings. Oversimplified, "king" subscribes to the "male" and the "royal" dimensions of the hidden layer, while "queen" subscribes to the "female" and "royal" dimensions of the hidden layer. Thus, when you take the vector for "king" and subtract out some approximation of the "male" dimensions and add in the "female" ones, you get something close to "queen". The most important takeaway is that this is more about the properties of language than Deep Learning. Any linear compression of these co-occurrence statistics will behave similarly.

Conclusion

Neural word embeddings and the impact of loss on learning.

In this chapter, we've unpacked the fundamental principles of using neural networks to study language. We started with an overview of the primary problems in Natural Language Processing, then proceeded to learn how neural networks model language at the word level using word embeddings. We further learned how our choice of loss function can change the kinds of properties that are captured by our word embeddings, and we finished with a discussion of perhaps the most magical of neural phenomena in this space: word analogies.

As with the other chapters, I would like to encourage you to build the examples in this chapter from scratch. While it might seem like this chapter is self standing, the lessons in loss function creation and tuning are invaluable and will be extremely important as we tackle increasingly more complicated strategies in future chapters. Good luck!

Neural Networks that Write like Shakesphere

Recurrent Layers for Variable Length Data

12

IN THIS CHAPTER

- The Challenge of Arbitrary Length
- The Surprising Power of Averaged Word Vectors
- The Limitations of Bag-of-Words Vectors
- Using Identity Vectors to Sum Word Embeddings
- Learning the Transition Matrices
- Learning to Create Useful Sentence Vectors
- Forward Propagation in Python
- Forward Propagation and Backpropagation with Arbitrary Length
- Weight Update with Arbitrary Length

“ There’s something magical about Recurrent Neural Networks. ”

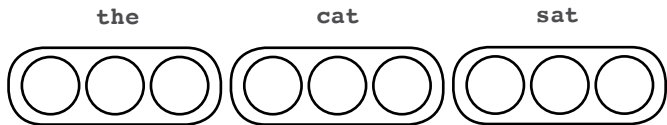
— ANDREJ KARPATHY

The Challenge of Arbitrary Length

Modeling arbitrarily long sequences of data with neural networks

This chapter and the previous chapter are quite intertwined, and I would encourage you to ensure that you have mastered the concepts and techniques of the previous chapter before you dive into this one. In the previous chapter, we learned about Natural Language Processing. This included learning how to modify a loss function to learn a specific pattern of information within the weights of our neural network. Futhermore, we developed an intuition for what a word embedding is and how it can represent shades of similarity with other word embeddings. In this chapter, we will expand upon this intuition of an embedding conveying the meaning of a single word by creating embeddings which convey the meaning of variable length phrases and sentences.

Let us first consider this challenge. If you wanted to create a vector that held an entire sequence of symbols within its contents in the same way that a word embedding stores information about a word, how would you accomplish this? Well, let's start with the simplest option. In theory, if we concatenated or stacked our word embeddings, we'd have a vector of sorts that held an entire sequence of symbols.



However, this approach leaves something to be desired, as different sentences will have different length vectors. This makes comparing two vectors together quite tricky, since one vector will "stick out" of the side. Consider the following second sentence:



In theory, these two sentences should be very similar, and comparing their vectors should indicate a high degree of similarity. However, since "the cat sat" is a shorter vector, we have to choose which part of "the cat sat still"s vector to compare to. If we align left, the vectors will appear to be exactly identical (ignoring the fact that "the cat sat still" is in fact a different sentence). However, if we align right, then the vectors will appear to be extraordinarily different, despite the fact that 3/4 words are the same, in the same order. So, while this naive approach shows some promise, it is far from ideal in terms of representing the meaning of a sentence in a useful way (a way that can be compared with other vectors).

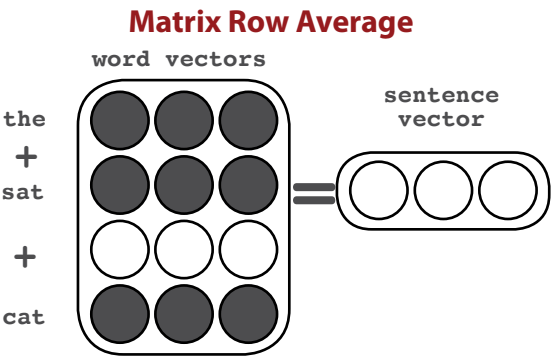
Do Comparisons Really Matter?

Why should we care about whether we can compare two sentence vectors?

The act of comparing two vectors is useful because it gives us an approximation of what the neural network sees. Even though we can't read two vectors, we can still tell when they're similar or different (using the function from the last chapter). If our method for generating sentence vectors doesn't reflect the similarity we observe between two sentences, then the network will also have difficulty recognizing when two sentences are similar. All it has to work with are the vectors!

So, as we continue to iterate and evaluate various methods for computing sentence vectors, I want you to remember why we're doing this. We're trying to take the perspective of a neural network. We're asking ourselves: "Will the correlation summarization find correlation between sentence vectors similar to this one and a desirable label, or will two nearly identical sentences instead generate wildly different vectors such that there is very little correlation between sentence vectors and the corresponding labels we're trying to predict?" We want to create sentence vectors that are useful for predicting things about the sentence, which, at a very minimum, means similar sentences need to create similar vectors.

So, if the previous way of creating our sentence vectors (concatenation) had some issues because of the rather arbitrary way of aligning them, let's explore what is perhaps the next simplest approach. What if we simply take the vector for each word in a sentence, and average them? Well, right off the bat, we don't have to worry about alignment because each sentence vector is of the same length!



Furthermore, the sentences "the cat sat" and "the cat sat still" would have quite similar sentence vectors because the words going into them are similar. Even better, it is quite likely that "a dog walked" would also be quite similar to "the cat sat" even though no words overlapped simply because the words used are also quite similar.

As it turns out, averaging word embeddings is a surprisingly effective way to create word embeddings. It's not perfect (as we'll see), but it does a rather amazingly strong job at capturing what you might perceive to be complex relationships between words. Before moving on, I think it's extremely beneficial to take our word embeddings from the last chapter and play around a little bit with our "average" strategy.

The Surprising Power of Averaged Word Vectors

It's the amazingly powerful go-to tool for neural prediction.

On the last page, we proposed our second method for creating vectors that convey the meaning of a sequence of words. This method simply took the average of the vectors corresponding to the words in a sentence, and intuitively we expect these new "average" sentence vectors to behave in several desirable ways. In the next couple of pages, we're going to play around with sentence vectors generated using the embeddings from the previous chapter. So, break out the code from the previous chapter, train the embeddings on the IMDB corpus as we did before, and let's experiment with "average" sentence embeddings!

On the right, you'll see the same normalization we performed when we were comparing word embeddings before. However, this time we'll pre-normalize all of the word embeddings into a matrix called `normed_weights`. Then, we create a function called `make_sent_vect` and use it to convert each review (list of words) into embeddings using our "average" approach. This is stored in the matrix `reviews2vectors`. After this, we create a function that queries for the most similar reviews given an input review, by performing a dot product between the input review's

```
import numpy as np
norms = np.sum(weights_0_1 * weights_0_1,axis=1)
norms.resize(norms.shape[0],1)
normed_weights = weights_0_1 * norms

def make_sent_vect(words):
    indices = list(map(lambda x:word2index[x],\
        filter(lambda x:x in word2index,words)))
    return np.mean(normed_weights[indices],axis=0)

reviews2vectors = list()
for review in tokens: # tokenized reviews
    reviews2vectors.append(make_sent_vect(review))
reviews2vectors = np.array(reviews2vectors)

def most_similar_reviews(review):
    v = make_sent_vect(review)
    scores = Counter()
    for i,val in enumerate(reviews2vectors.dot(v)):
        scores[i] = val
    most_similar = list()

    for idx,score in scores.most_common(3):
        most_similar.append(raw_reviews[idx][0:40])
    return most_similar
most_similar_reviews(['boring','awful'])

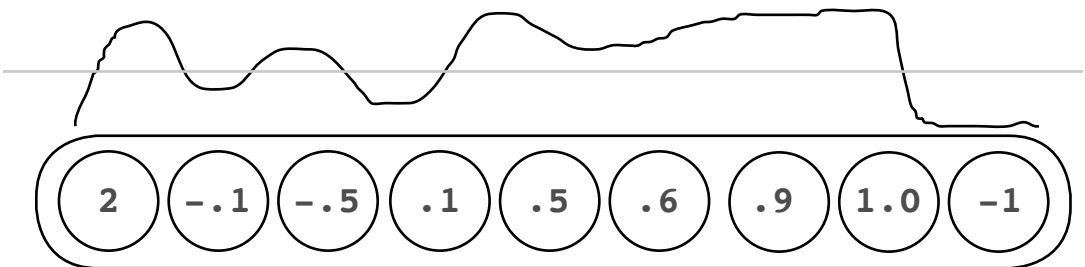
['I am amazed at how boring this film',
 'This is truly one of the worst dep',
 'It just seemed to go on and on and.]
```

vector and the vector of every other review in the corpus. This dot product similarity metric is the same one that we briefly discussed back in Chapter 4 when we were learning to predict with multiple inputs. Perhaps surprisingly, when we query for the most similar reviews to the average vector between the two words 'boring' and 'awful', we receive back three very negative reviews. It does appear that there is likely to be interesting statistical information within these vectors, such that negative and positive embeddings cluster together.

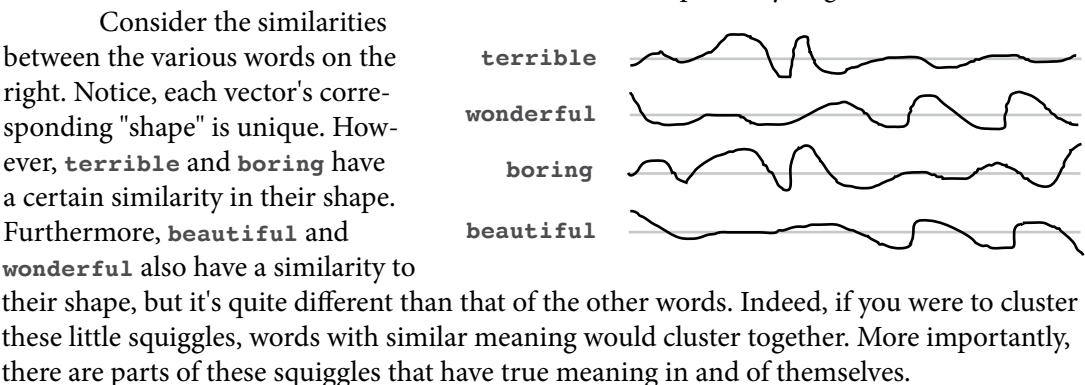
How is Information Stored in These Embeddings?

When we average word embeddings, average shapes remain.

Now, to consider what's going on here requires a little bit of abstract thought, and I'd recommend digesting this kind of information over a period of time, as it is perhaps a different kind of lesson than the one we are used to. For a moment, I would like for you to consider that a word vector can be visually observed as a *squiggly line* like this one.



Instead of thinking of a vector as a list of numbers, I want you to think about it as if it's a line with high and low points corresponding to high and low values at different places in the vector. So, if we selected several words from our corpus, they might look like this:



For example, for the negative words, there's a downward and then upward spike about 40% from the left. If I were to continue drawing lines corresponding to words, this spike would continue to be distinctive. Now, there's nothing magical about that spike meaning "negativity" and if I re-trained the network it would likely show up somewhere else. Indeed, the spike only indicates "negativity" because all of the negative words have it! Thus, during the course of training, these shapes are molded such that different curves in different positions convey meaning (as discussed in the previous chapter). Thus, when we take an "average curve" over the words in a sentence, the most dominant meanings of the sentence hold true, and the noise created by any particular word gets averaged away.

How does a Neural Network Use Embeddings?

Neural networks detect the "curves" that have correlation with a target label.

On the last page, we learned about a new way to "view" word embeddings as a squiggly line with distinctive properties (curves). We also learned that these curves are developed throughout the course of training to accomplish the target objective. Furthermore, words with similar meaning in one way or another will often share a distinctive "bend" in the curve, some combination of high-low pattern amongst the weights. In this page, we want to consider how the correlation summarization processes these curves as input. What does it mean for a layer to consume these curves as input?

Well, truth be told, a neural network consumes embeddings just like it consumed our streetlight dataset in the very early chapters. It looks for correlation between the various bumps and curves in the hidden layer and the target label it's trying to predict. Furthermore, this is why words with one particular aspect of similarity will share similar bumps and curves. At some point, during training, a neural network started developing unique characteristics between the shapes of different words so that it could tell them apart, and grouping them (giving them similar bumps/curves) to help make accurate predictions. However, this is simply another way of summarizing the lessons we learned at the end of the last chapter. We want to press in further.

In this chapter, we want to consider what it means for us to sum these embeddings into a sentence embedding. What kinds of classifications would this summed vector be useful for? At the end of the last page, we identified that taking an average across the word embeddings of a sentence would result in a vector that had an average of the characteristics of the words therein. Thus, if there were many positive words, the final embedding would look somewhat positive (with other noise from the words generally cancelling out). However, one should note that this approach is a bit "mushy". If we have enough words, these different wavy lines should simply all average together to generally just be a straight line.

This brings us to the first weakness of this approach: we're attempting to store arbitrarily long sequences (a sentence) of information into a fixed length vector, and if we try to store too much eventually the sentence vector (being an average of a multitude of word vectors) will simply average out to a straight line (a vector of near-0s).

In short, this process of storing the information of a sentence doesn't decay nicely. If we try to store too many words into a single vector, we end up storing almost nothing. That being said, a sentence is often not that many words, and if a sentence has repeating patterns then these sentence vectors can in fact be quite useful, as the sentence vector will retain the most dominant patterns present across the word vectors being summed (such as the "negative spike" on the last page).

The Limitations of Bag-of-Words Vectors

Order becomes irrelevant when you average word embeddings.

However, the biggest issue with "average" embeddings is that they have no concept of order. For example, let's say we have two sentences "Yankees defeat Red Socks" and "Red Socks defeat Yankees". Generating sentence vectors for each of these two sentences using our "average" approach will yield exactly identical vectors! However, these sentences are conveying the exact opposite information! Furthermore, this approach entirely ignores grammar and syntax, as "Socks Red Yankees defeat" will also yield an exactly identical sentence embedding.

This approach of summing or averaging word embeddings to form the embedding for a phrase or sentence is classically known as a "Bag of Words" approach because, much like throwing a bunch of words into a bag, order is not preserved. The key limitation here is that you can take any sentence, scramble all the words around, generate a sentence vector, and no matter how you scrambled the words, the vector will be the same (because addition is associative, i.e. $a + b == b + a$).

The *real topic* of this chapter is about generating sentence vectors in a way where order *does* matter. We want to create vectors in a way where scrambling them around actually changes the resulting vector. More importantly, we want the *way in which order matters* (otherwise known as *the way in which order changes the vector*) to be **learned**. In this way, our neural network's representation of order can be based around trying to solve a task in language and, by extension, hopefully capture the essence of order in language itself. Now, I'm using "language" as an example here, but we can generalize these statements to any sequence. Language is just a particularly challenging yet universally known domain.

One of the most famous and successful ways of generating vectors for sequences (such as a sentence) is called a Recurrent Neural Network (RNN), and in order to show you how it works, we're first going to start by coming up with a new, and seemingly wasteful, way of doing our "average" word embeddings using something called an **Identity Matrix**. So, before we jump into Recurrent Neural Networks, let's first learn about the object known as an Identity Matrix. An identity matrix is just an arbitrarily large, square matrix (num rows == num columns) of 0s with 1s stretching from the top left corner to the bottom right corner like in the examples on the right.

All three of these matrices are "identity" matrices, and they have one purpose. Performing vector-matrix multiplication with ANY vector will simply return the original vector. So, if I take the vector [3,5] and multiply it by the top identity matrix, the result will be [3,5].

[1, 0]
[0, 1]

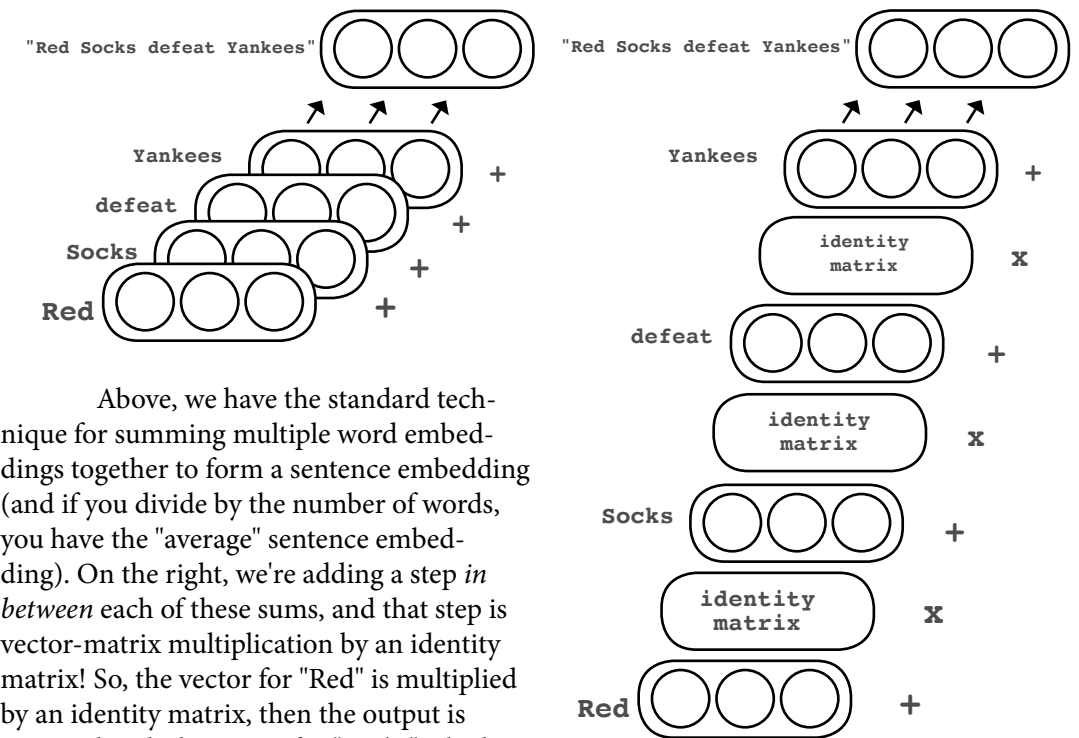
[1, 0, 0]
[0, 1, 0]
[0, 0, 1]

[1, 0, 0, 0]
[0, 1, 0, 0]
[0, 0, 1, 0]
[0, 0, 0, 1]

Using Identity Vectors to Sum Word Embeddings

We're going to implement the exact same logic using a different approach.

Perhaps you will think that identity matrices are actually quite useless. For what purpose is a matrix that takes a vector and outputs that same vector? In our case, we're going to use it as a teaching tool to show how we can setup a more complicated way of summing our word embeddings so that the neural network can take order into account when generating the final sentence embedding. So, let's explore another way of summing embeddings.



Above, we have the standard technique for summing multiple word embeddings together to form a sentence embedding (and if you divide by the number of words, you have the "average" sentence embedding). On the right, we're adding a step *in between* each of these sums, and that step is vector-matrix multiplication by an identity matrix! So, the vector for "Red" is multiplied by an identity matrix, then the output is summed with the vector for "Socks" which is then vector-matrix multiplied by the identity matrix and added to the vector "defeat" and so on throughout the sentence. Note that because the vector-matrix multiplication by the identity matrix returns the same vector that goes into it, the process on the right yields *exactly the same sentence embedding* as the process on the top left. Yes, we're doing wasteful computation, but that's all about to change. The main thing to consider here is that if the matrices we used were any matrix other than the identity matrix, changing the order of the words would change the resulting embedding. Let's see this in Python!

Matrices That Change Absolutely Nothing

Let's create sentence embeddings using Identity Matrices in Python.

On the previous page, we outlined a new technique for summing word embeddings using a new tool, the Identity Matrix. On this page, we're going to demonstrate how to play with Identity Matrices in Python and ultimately how to implement the new sentence vector technique from the last page (proving that it produces identical sentence embeddings).

On the right, we first initialize 4 vectors (a, b, c, and d) of length 3 as well as an identity matrix with 3 rows and 3 columns (identity matrices are always square). Notice that the identity matrix has the characteristic set of 1s running diagonally from the top left to the bottom right (which, by the way, is called "the diagonal" in Linear Algebra). Any square matrix with 1s along "the diagonal" and 0s everywhere else is an identity matrix.

We then proceed to perform vector->matrix multiplication with each of our vectors and the identity matrix (using NumPy's "dot" function). As you can see, the output of this process simply yields a new vector that is identical to the input vector.

Since vector->matrix multiplication by an identity matrix returns the exact same vector that you started with, incorporating this process into our sentence embedding should seem quite trivial, and in fact it is, as we can see below.

```
import numpy as np

a = np.array([1,2,3])
b = np.array([0.1,0.2,0.3])
c = np.array([-1,-0.5,0])
d = np.array([0,0,0])

identity = np.eye(3)
print(identity)

[[ 1.  0.  0.]
 [ 0.  1.  0.]
 [ 0.  0.  1.]]

print(a.dot(identity))
print(b.dot(identity))
print(c.dot(identity))
print(d.dot(identity))

[ 1.  2.  3.]
[ 0.1  0.2  0.3]
[-1. -0.5  0.]
[ 0.  0.  0.]
```

```
this = np.array([2,4,6])
movie = np.array([10,10,10])
rocks = np.array([1,1,1])

print(this + movie + rocks)
print((this.dot(identity) + movie).dot(identity) + rocks)
```

```
[13 15 17]
[13. 15. 17.]
```

As you can see, both ways of creating sentence embeddings generate the exact same vector. This is only because the identity matrix is a very special kind of matrix. However, what would happen if we didn't use the identity matrix? What if, instead, we used a different matrix? In fact, the identity matrix is the **ONLY** matrix that is guaranteed to return the same vector that it is vector->matrix multiplied with. No other matrix has this guarantee.

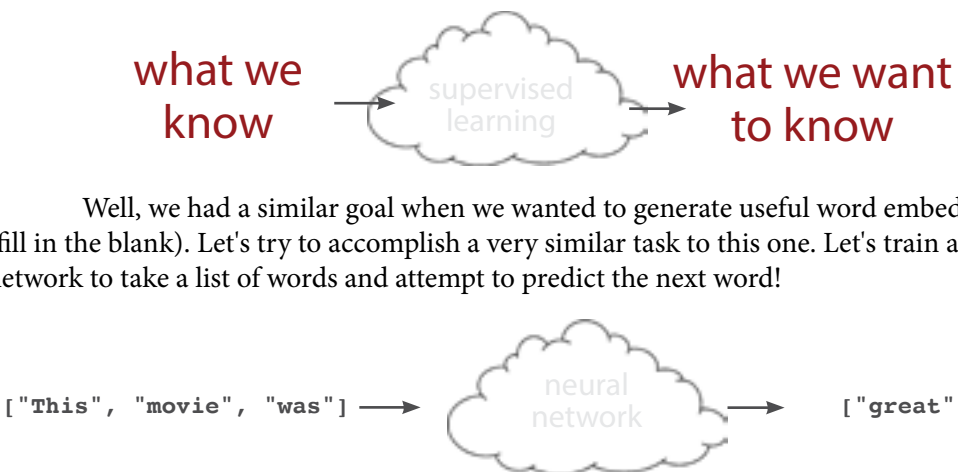
Learning the Transition Matrices

What if we allowed the identity matrices to change to minimize the loss?

Before we begin, let's remember our goal. We want to generate sentence embeddings that cluster according to the meaning of the sentence, such that when we take a sentence, we can use the vector to find sentences with a similar meaning. More specifically, we want these sentence embeddings to care about the order of words. Previously, we simply tried summing word embeddings. However, this meant that "Red Socks defeat Yankees" had an identical vector to the sentence "Yankees defeat Red Socks" despite the fact that these two sentences have opposite meanings. Instead, we want to form sentence embeddings where these two sentences generate *different* embeddings (yet still cluster in a meaningful way). Our theory is that if we use our Identity Matrix way of creating sentence embeddings, but used *any other matrix other than the identity matrix*, then sentence embeddings would be different depending on the order.

Now the obvious question: what matrix do we use instead of the Identity Matrix? Well, there are an infinite number of choices. However, in Deep Learning the standard answer to this kind of question is, simply, "We'll learn the matrix just like we learn any other matrix in a neural network!" Ok, so we'll just learn this matrix. How will we do that?

Well, whenever we want to train a neural network to learn something, we always need to have a task for it to learn. In this case, we want that task to require it to generate interesting sentence embeddings by learning both useful word vectors AND useful modifications to our Identity Matrices. So, what task should we use?



Well, we had a similar goal when we wanted to generate useful word embeddings (fill in the blank). Let's try to accomplish a very similar task to this one. Let's train a neural network to take a list of words and attempt to predict the next word!

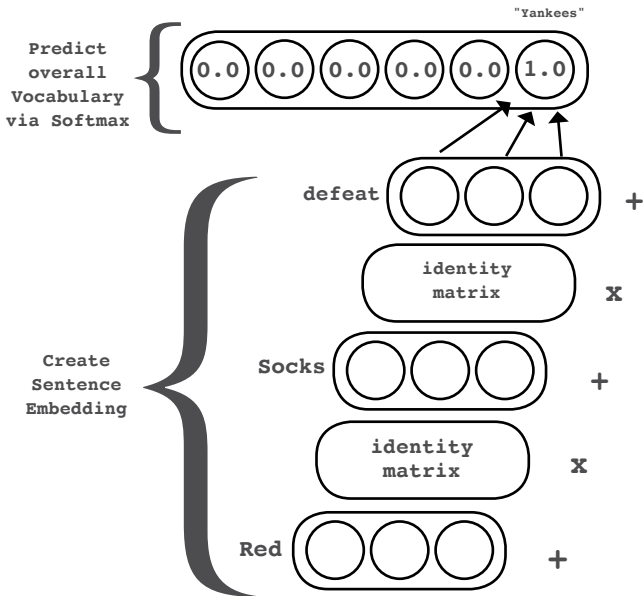
Learning To Create Useful Sentence Vectors

Create the sentence vector, make prediction, modify sentence vector via its parts.

So, in this next experiment, I don't want you to think about our network like you would previous neural networks. In this experiment, I want you to think about creating a sentence embedding, using it to predict the next word, and then modifying the respective parts that formed our sentence embedding to make this prediction more accurate. Since we're predicting the next word, our sentence embedding is going to be made from the "parts of the sentence we've seen so far". So, our neural network will look something like this:

On the right, we see our new neural network. It is composed of two steps. First, we create our sentence embedding, and then we use that embedding to predict which word comes next. Thus, the input to this network is the text "Red Socks defeat" and the word to be predicted is the word "Yankees".

Furthermore, while I've written "identity matrix" in the boxes between our word vectors, this is actually a matrix that will only *start* as an identity matrix. During training, we'll backpropagate gradients into these matrices and update them to help the network better make predictions (just like we do for the rest of the weights in this network). In this way, our network will *learn how to incorporate more information than just a sum of word embeddings*. By allowing our (initially "identity") matrices to change (and become NOT identity matrices anymore), we allow our neural network to learn how to create embeddings where the order in which the words are presented changes the sentence embedding. However, this change is not simply arbitrary. The network will learn how to incorporate the order of words in such a way that is *useful for the task of predicting the next word*. Furthermore, we will constrain our transition matrices (the matrices that are originally identity matrices) to all be the same matrix. In other words, the matrix from Red -> Socks will be re-used to transition from Socks -> defeat. In this way, whatever logic it learns in one transition will be reused in the next and only logic that is useful at every predictive step will be allowed to be learned in the network.



Forward Propagation in Python

Let's take this idea and show how to perform a simple forward propagation.

Now that we have the conceptual idea of what we're trying to build, let's check out a toy version in Python. First, let's setup our weights (I'm using a limited vocab of 9 words).

```
import numpy as np

def softmax(x_):
    x = np.atleast_2d(x_)
    temp = np.exp(x)
    return temp / np.sum(temp, axis=1, keepdims=True)

word_vects = {}
word_vects['yankees'] = np.array([[0.,0.,0.]])
word_vects['bears'] = np.array([[0.,0.,0.]])
word_vects['braves'] = np.array([[0.,0.,0.]])
word_vects['red'] = np.array([[0.,0.,0.]])
word_vects['socks'] = np.array([[0.,0.,0.]])
word_vects['lose'] = np.array([[0.,0.,0.]])
word_vects['defeat'] = np.array([[0.,0.,0.]])
word_vects['beat'] = np.array([[0.,0.,0.]])
word_vects['tie'] = np.array([[0.,0.,0.]])

sent2output = np.random.rand(3,len(word_vects))

identity = np.eye(3)
```

Word
Embeddings

Sentence Embedding
to
Output Classification
Weights

Transition
Weights

}

In the code above, we created three sets of weights. We created a Python dictionary of word embeddings, and our identity matrix (transition matrix), as well as a classification layer. This classification layer "sent2output" is just a weight matrix to predict the next word given a sentence vector of length 3. With these tools, forward propagation is quite trivial. Below, we can see how forward propagation works with the sentence "red socks defeat" -> "yankees".

```
layer_0 = word_vects['red']
layer_1 = layer_0.dot(identity) + word_vects['socks']
layer_2 = layer_1.dot(identity) + word_vects['defeat']

pred = softmax(layer_2.dot(sent2output))
print(pred)
```

Create
Sentence
Embedding

Predict
over all
Vocabulary

}

```
[[ 0.11111111  0.11111111  0.11111111  0.11111111  0.11111111  0.11111111
   0.11111111  0.11111111  0.11111111]]
```

How do we Backpropagate into this?

It might seem a bit trickier, but it's actually the same steps that we already learned.

On the bottom of the previous page, we showed how to perform forward prediction for this network. At first, it might not be clear how backpropagation can be performed. However, it's really quite simple. Perhaps this is what you see:

Normal Neural Network
(Chapters 1-5)

layer_0 = word_vects['red']
layer_1 = layer_0.dot(identity) + word_vects['socks']
layer_2 = layer_1.dot(identity) + word_vects['defeat']

some sort of strange
additional piece

```
pred = softmax(layer_2.dot(sent2output))
print(pred)
```

Normal Neural
Network Again
(Chapter 9 Stuff)

Based on previous chapters, you should feel comfortable with computing a loss and backpropagating until you get to the gradients at layer_2, which we'll call layer_2_delta. At this point, you might be wondering, "which direction do I backprop in?" Gradients could go back to layer_1 by going backwards through the identity matrix multiplication, or they could go into word_vects['defeat']. The truth is, when you add two vectors together during forward propagation, then you back propagate the same gradient into BOTH sides of the addition. So, when we generate layer_2_delta, we will backpropagate it twice, once across the identity matrix to create layer_1_delta, and again to word_vects['defeat'].

```
y = np.array([1,0,0,0,0,0,0,0,0]) # target one-hot vector for "yankees"

pred_delta = pred - y
layer_2_delta = pred_delta.dot(sent2output.T)
defeat_delta = layer_2_delta * 1 # can ignore the "1" like prev. chapter
layer_1_delta = layer_2_delta.dot(identity.T)
socks_delta = layer_1_delta * 1 # again... can ignore the "1"
layer_0_delta = layer_1_delta.dot(identity.T)
alpha = 0.01
word_vects['red'] -= layer_0_delta * alpha
word_vects['socks'] -= socks_delta * alpha
word_vects['defeat'] -= defeat_delta * alpha
identity -= np.outer(layer_0,layer_1_delta) * alpha
identity -= np.outer(layer_1,layer_2_delta) * alpha
sent2output -= np.outer(layer_2,pred_delta) * alpha
```

Let's Train it!

We've got all the tools, now let's train it on a toy corpus.

So that we can get an intuition for what's going on, we're first going to train our new network on a toy task called the Babi dataset. This dataset is a synthetically generated Question-Answer corpus to teach machines how to answer simple questions about an environment. While we're not using it for QA (yet), the simplicity of the task will help us better see the impact made by learning our identity matrix. First, let's download the Babi dataset.

Bash Commands

```
wget http://www.thespermwhale.com/jaseweston/babi/tasks_1-20_v1-1.tar.gz
tar -xvf tasks_1-20_v1-1.tar.gz
```

And with some simple python, we can open and clean a small dataset to train our network.

```
import sys, random, math
from collections import Counter
import numpy as np

f = open('tasksv11/en/qal_single-supporting-fact_train.txt', 'r')
raw = f.readlines()
f.close()

tokens = list()
for line in raw[0:1000]:
    tokens.append(line.lower().replace("\n", "").split(" ")[1:])

print(tokens[0:3])

[['Mary', 'moved', 'to', 'the', 'bathroom'],
 ['John', 'went', 'to', 'the', 'hallway'],
 ['Where', 'is', 'Mary', 'bathroom'],
```

As you can see, this dataset contains a variety of simple statements and questions (with punctuation removed). Each question is simply followed by the correct answer. When used in the context of QA, a neural network simply reads the statements in order and answers questions (either correctly or incorrectly) based on information in the recently read statements. For now, we're just going to train our network to attempt to finish each sentence when given one or more starting words. Along the way, we're going to see the importance of allowing our "recurrent" matrix (previously the "identity" matrix) to learn.

Setting Things Up

Before we can create matrices, we need to learn how many parameters we have.

As with our word embedding neural network, we first need to create a few useful counts, lists, and utility functions that we will use during the predict, compare, learn process. These utility functions and objects can be seen below. They should look familiar.

```
vocab = set()
for sent in tokens:
    for word in sent:
        vocab.add(word)

vocab = list(vocab)

word2index = {}
for i, word in enumerate(vocab):
    word2index[word] = i

def words2indices(sentence):
    idx = list()
    for word in sentence:
        idx.append(word2index[word])
    return idx

def softmax(x):
    e_x = np.exp(x - np.max(x))
    return e_x / e_x.sum(axis=0)
```

On the left, we create a simple list of our vocabulary words as well as a lookup dictionary allowing us to go back and forth between a word's text and its index. We will use its index in the vocabulary list to pick which row and column of our embedding and prediction matrices correspond to which word. On the right, we have a utility function for converting a list of words to a list of indices, as well as the function for softmax, which we'll use for predicting the next word.

In the code block below, we initialize our random seed (to get consistent results), then set our embedding size to 10. We create a matrix of word embeddings, recurrent embeddings, and an initial `start` embedding. This is the embedding modeling an "empty" phrase, which is key to the network modeling how sentences tend to start. Finally, we have a `decoder` weight matrix (just like from embeddings), and a `one_hot` utility matrix.

```
np.random.seed(1)
embed_size = 10

# word embeddings
embed = (np.random.rand(len(vocab), embed_size) - 0.5) * 0.1

# embedding -> embedding (initially the identity matrix)
recurrent = np.eye(embed_size)

# sentence embedding for empty sentence
start = np.zeros(embed_size)

# embedding -> output weights
decoder = (np.random.rand(embed_size, len(vocab)) - 0.5) * 0.1

# one hot lookups (for loss function)
one_hot = np.eye(len(vocab))
```

Forward Propagation with Arbitrary Length

We will forward propagate using the exact same logic described several pages ago.

Below you can see the logic by which we forward propagate and predict the next word. Note that while the construction might feel unfamiliar, it follows the exact same procedure as we had before for summing embeddings while using the identity matrix. However, we've replaced the identity matrix with a matrix called "recurrent", which is initialized to be all zeros (and will be learned through training). Furthermore, instead of only predicting at the last word, we make a prediction (`layer['pred']`) at every timestep, which is based on the embedding generated by the previous words. This is more efficient than doing a new forward propagation from the very beginning of the phrase each time we want to predict a new term.

```
def predict(sent):
    layers = list()
    layer = {}
    layer['hidden'] = start
    layers.append(layer)

    loss = 0

    # forward propagate
    preds = list()
    for target_i in range(len(sent)):

        layer = {}

        # try to predict the next term
        layer['pred'] = softmax(layers[-1]['hidden'].dot(decoder))

        loss += -np.log(layer['pred'][sent[target_i]])

        # generate the next hidden state
        layer['hidden'] = layers[-1]['hidden'].dot(recurrent) + \
            embed[sent[target_i]]

        layers.append(layer)
    return layers, loss
```

While there's nothing particularly new about this bit of code relative to what we've learned in the past, there's a particular piece you want to make sure you're familiar with before we move forward. The list called `layers` is a new way for us to forward propagate. Notice that we end up doing more forward propagations if the length of `sent` is larger. As a result, we can't just use static layer variables like we did before. This time, we need to simply keep appending new layers to the list based on the number we need. Make sure you're comfortable with what's going on in each part of this list, because if it's unfamiliar to you in the forward propagation pass, it will be very difficult to know what's going on during the backpropagation and weight update steps.

Backpropagation with Arbitrary Length

We will backpropagate using the exact same logic as described several pages ago.

As described several pages ago with the example "Red Socks defeat Yankees", we're now going to implement backpropagation over arbitrary length sequences assuming we have access to the forward propagation objects returned from the function on the previous page. The most important object is our `layers` list, which has two vectors (`layer['state']` and `layer['previous->hidden']`). In order to backpropagate, we're going to take our output gradient and add a new object to each list called `layer['state_delta']` which is going to represent the gradient at that layer. This corresponds to variables like `socks_delta`, `layer_0_delta`, and `defeat_delta` from our "Red Socks defeat Yankees" example previously described. We're just building the same logic in a way that it can consume the variable length sequences from our forward propagation logic.

```
# forward
for iter in range(30000):
    alpha = 0.001
    sent = words2indices(tokens[iter%len(tokens)][1:])
    layers, loss = predict(sent)

    # back propagate
    for layer_idx in reversed(range(len(layers))):
        layer = layers[layer_idx]
        target = sent[layer_idx-1]

        if(layer_idx > 0): # if not the first layer
            layer['output_delta'] = layer['pred'] - one_hot[target]
            new_hidden_delta = layer['output_delta'] \
                .dot(decoder.transpose())

            # if the last layer - don't pull from a
            # later one because it doesn't exist
            if(layer_idx == len(layers)-1):
                layer['hidden_delta'] = new_hidden_delta
            else:
                layer['hidden_delta'] = new_hidden_delta + \
                    layers[layer_idx+1]['hidden_delta'] \
                        .dot(recurrent.transpose())

        else: # if the first layer
            layer['hidden_delta'] = layers[layer_idx+1]['hidden_delta'] \
                .dot(recurrent.transpose())
```

Before moving on to the next section, make sure you can read along this code and explain it to a friend (or at least to yourself). There are no new concepts in this code, but the construction of it can make it seem a bit foreign at first. Just spend some time linking what's written in this code back to each line of the "Red Socks defeat Yankees" example, and you should be ready for the next section where we update our weights using the gradients we backpropagated.

Weight Update with Arbitrary Length

We will update weights using the exact same logic as previously described.

As with the forward and backprop logic, this weight update logic is not new. However, I'm presenting it after having explained it so that you can now focus on the engineering complexity having (hopefully) already grokked (ha!) the theory complexity.

```
# forward
for iter in range(30000):
    alpha = 0.001
    sent = words2indices(tokens[iter%len(tokens)][1:])

    layers, loss = predict(sent)

# back propagate
for layer_idx in reversed(range(len(layers))):
    layer = layers[layer_idx]
    target = sent[layer_idx-1]

    if(layer_idx > 0):
        layer['output_delta'] = layer['pred'] - one_hot[target]
        new_hidden_delta = layer['output_delta']\
            .dot(decoder.transpose())

        # if the last layer - don't pull from a
        # later one because it doesn't exist
        if(layer_idx == len(layers)-1):
            layer['hidden_delta'] = new_hidden_delta
        else:
            layer['hidden_delta'] = new_hidden_delta + \
                layers[layer_idx+1]['hidden_delta']\
                    .dot(recurrent.transpose())
    else:
        layer['hidden_delta'] = layers[layer_idx+1]['hidden_delta']\
            .dot(recurrent.transpose())

# update weights
start -= layers[0]['hidden_delta'] * alpha / float(len(sent))
for layer_idx, layer in enumerate(layers[1:]):

    decoder -= np.outer(layers[layer_idx]['hidden'],\
        layer['output_delta']) * alpha / float(len(sent))

    embed_idx = sent[layer_idx]
    embed[embed_idx] -= layers[layer_idx]['hidden_delta'] * \
        alpha / float(len(sent))

    recurrent -= np.outer(layers[layer_idx]['hidden'],\
        layer['hidden_delta']) * alpha / float(len(sent))

if(iter % 1000 == 0):
    print("Perplexity:" + str(np.exp(loss/len(sent))))
```

Execution and Output Analysis

We will update weights using the exact same logic as previously described.

Now the moment of truth: what happens when we run it? Well, when I run this code, I get a relatively steady downtrend in a metric called "Perplexity". Technically, the perplexity is just the probability of the correct label (word), passed through a **log** function, **negated**, and **exponentiated** (e^x). However, what it represents theoretically is the difference between two probability distributions. In our case, the perfect probability distribution would be 100% probability allocated to the correct term and 0% everywhere else. Perplexity is high when two probability distributions do not match, and it is low (approaching 1) when they do match. Thus, a decreasing perplexity, like all loss functions used with Stochastic Gradient Descent, is a good thing! It means our network is learning to predict probabilities that match the data.

```
Perplexity:82.09227500075585
Perplexity:81.87615610433569
Perplexity:81.53705034457951
.....
Perplexity:4.132556753967558
Perplexity:4.071667181580819
Perplexity:4.0167814473718435
```

However, this hardly tells us what's going on in the weights. In fact, perplexity has faced some criticism over the years (particularly in the Language Modeling community) for being over-used as a metric. We want to look a little more closely at the predictions.

```
sent_index = 4

l,_ = predict(words2indices(tokens[sent_index]))

print(tokens[sent_index])

for i,each_layer in enumerate(l[1:-1]):
    input = tokens[sent_index][i]
    true = tokens[sent_index][i+1]
    pred = vocab[each_layer['pred']].argmax()
    print("Prev Input:" + input + (' ' * (12 - len(input))) + \
        "True:" + true + (" " * (15 - len(true))) + "Pred:" + pred)
```

In this code, we take a sentence and simply predict the word that the model thinks is most likely. This is useful because we can get a sense for the kinds of characteristics the model takes on. What kinds of things does it get right? What kinds of mistakes does it make? We'll take a look on the next page.

Execution and Output Analysis (cont.)

Looking at predictions can help us understand what's going on.

We can look at the output predictions of our neural network as it trains to not only learn what kinds of patterns it picks up, but also the order in which it learns them. So, after 100 training steps, the output looks like so:

```
[ 'sandra', 'moved', 'to', 'the', 'garden.' ]
Prev Input:sandra      True:moved      Pred:is
Prev Input:moved       True:to        Pred:kitchen
Prev Input:to          True:the       Pred:bedroom
Prev Input:the         True:garden.   Pred:office
```

You'll find that neural networks tend to start off quite random. In the case above, the neural network is likely only biased toward whatever words it started with in its first random state. We should keep training:

```
[ 'sandra', 'moved', 'to', 'the', 'garden.' ]
Prev Input:sandra      True:moved      Pred:the
Prev Input:moved       True:to        Pred:the
Prev Input:to          True:the       Pred:the
Prev Input:the         True:garden.   Pred:the
```

After 10,000 training steps, the neural network simply picks out the most common word "the" and predicts it at every timestep. This is an extremely common error in recurrent neural networks. It takes lots of training to learn finer grained detail in a highly skewed dataset.

```
[ 'sandra', 'moved', 'to', 'the', 'garden.' ]
Prev Input:sandra      True:moved      Pred:is
Prev Input:moved       True:to        Pred:to
Prev Input:to          True:the       Pred:the
Prev Input:the         True:garden.   Pred:bedroom.
```

Now, these mistakes are really interesting. After only seeing the word "sandra", the network predicts "is" which, while not exactly the same as "moved", isn't a bad guess. It simply picked the wrong verb. Next, notice that the words "to" and "the" were correct, which isn't as surprising since these are some of the more common words in the dataset, and presumably the network has been trained to predict the phrase "to the" after the verb "moved" many times. The final mistake is also quite compelling, mistaking the room "bedroom" for the word "garden".

Something important to note here, there's almost no way this neural network could perfectly learn this task. After all, if I gave you the words "sandra moved to the", could you tell me the correct next word? More context is needed to solve this task, but the fact that it's unsolvable, in my opinion, creates educational analysis for the ways in which it fails.

Conclusion and Review

Recurrent Neural Networks Predict Over Arbitrary Length Sequences.

In this chapter, we learned how to create vector representations for arbitrary length sequences. In the last exercise, we trained a linear recurrent neural network to predict the next term given a previous phrase of terms. In order to do this, it needed to learn how to create embeddings that accurately represented variable length strings of terms into a fixed size vector.

Now, this last sentence should drive home a question: how does a neural network fit a variable amount of information into a fixed size box? Well, the truth is that these sentence vectors don't encode everything in the sentence. The name of the game in recurrent neural networks is not just about what these vectors remember, but also what they forget. In the case of predicting the next word, most RNNs learn that only the last couple of words are really necessary (<https://arxiv.org/abs/1702.04521>), and they learn to forget (aka, not make unique patterns in their vectors for) words farther back in the history.

However, one should note that we didn't have any non-linearities in the generation of these representations. What kinds of limitations do you think that could create? In the next chapter, we're going to explore these questions and more using nonlinearities and gates to form a neural network called a Long-Short Term Memory network. However, before we do that, make sure you can sit down and (from memory) code a working Linear RNN that converges. The dynamics and control flow of these networks can be a bit daunting, and the complexity is about to jump by quite a bit. So, before moving on, become comfortable with what we've made in this chapter.

And with that, let's dive into LSTMs!