

COMPUTER SCIENCE TRIPOS - PART II PROJECT

Language Modelling for Text Prediction

April 14, 2017

supervised by
Dr Marek Rei & Dr Ekaterina Shutova

Proforma

Name: **Devan Kuleindiren**
College: **Robinson College**
Project Title: **Language Modelling for Text Prediction**
Examination: **Computer Science Tripos – Part II, June 2017**
Word Count: **?**
Project Originator: **Devan Kuleindiren & Dr Marek Rei**
Supervisors: **Dr Marek Rei & Dr Ekaterina Shutova**

Original Aims of the Project

The primary aim of the project was to implement and benchmark a variety of language models, comparing the quality of their predictions as well as the time and space that they consume. More specifically, I aimed to build an n -gram language model along with several smoothing techniques, and a variety of recurrent neural network-based language models. An additional aim was to investigate ways to improve the performance of existing language models on error-prone text.

Work Completed

All of the project aims set out in the proposal have been met, resulting in a series of language model implementations and a generic benchmarking framework for comparing their performance. I have also proposed and evaluated a novel extension to an existing language model which improves its performance on error-prone text. Additionally, as an extension, I implemented a mobile keyboard on iOS that uses my language model implementations as a library.

Special Difficulties

None.

Declaration

I, Devan Kuleindiren of Robinson College, being a candidate for Part II of the Computer Science Tripos, hereby declare that this dissertation and the work described in it are my own work, unaided except as may be specified below, and that the dissertation does not contain material that has already been used to any substantial extent for a comparable purpose.

SIGNED

DATE

Contents

1	Introduction	5
1.1	Language Models	5
1.2	Motivation	6
1.3	Related Work	6
2	Preparation	7
2.1	n -gram Models	7
2.1.1	An Overview of n -gram Models	7
2.1.2	Smoothing Techniques	8
2.2	Recurrent Neural Network Models	11
2.2.1	An Overview of Neural Networks	11
2.2.2	Recurrent Neural Networks	15
2.2.3	Word Embeddings	17
2.3	Software Engineering	17
2.3.1	Requirements	17
2.3.2	Tools and Technologies Used	18
2.3.3	Starting Point	19
3	Implementation	20
3.1	System Overview	20
3.1.1	The Language Model Hierarchy	21
3.1.2	Testing	22
3.2	n -gram Models	23
3.2.1	Computing the Vocabulary	23
3.2.2	Counting n -grams Efficiently	23
3.2.3	Precomputing Smoothing Coefficients	25
3.3	Recurrent Neural Network Models	25
3.3.1	TensorFlow	26
3.3.2	Network Structure	27
3.3.3	Long Short-Term Memory and Gated Recurrent Units	27
3.3.4	Word Embeddings	31
3.3.5	Putting Theory into Practice	32
3.4	Mobile Keyboard	33
3.4.1	Updating Language Model Predictions On the Fly	34
3.5	Extending Models to Tackle Error-Prone Text	35
3.5.1	The Approach	35

4	Evaluation	37
4.1	Evaluation Methodology	37
4.1.1	Metrics	37
4.1.2	Datasets	39
4.2	Results	40
4.2.1	Existing Models	40
4.2.2	On a Mobile Device	43
4.2.3	On Error-Prone Text	44
5	Conclusions	48
	Bibliography	48
A	Backpropagation Recurrence Relation	51
B	Project Proposal	53

Chapter 1

Introduction

My project investigates the performance of various language models in the context of text prediction. I started by implementing a series of well-established models and comparing their performance, before assessing the tradeoffs that occur when you attempt to apply them in a practical context, such as in a mobile keyboard. Finally, I proposed a novel extension to an existing model which aims to improve its performance on error-prone text.

1.1 Language Models

Language models (LMs) produce a probability distribution over a sequence of words, which can be used to estimate the relative likelihood of words or phrases occurring in various contexts. This predictive power is useful in a variety of applications. For example, in speech recognition, if the speech recogniser has estimated two candidate word sequences from an acoustic signal; ‘*it’s not easy to wreck a nice beach*’ and ‘*it’s not easy to recognise speech*’, then a language model can be used to determine that the second candidate is more probable than the first. Language models are also used in machine translation, handwriting recognition, part-of-speech tagging and information retrieval.

$$\begin{array}{rcl} \overbrace{\text{Do you want to grab a}}^{w_1^k} \overbrace{\text{drink}}^{w_{k+1}} & \mathbb{P}(w_{k+1}|w_1^k) & \\ & (0.327) & \\ & \text{coffee} & (0.211) \\ & \text{bite} & (0.190) \\ & \text{spot} & (0.084) \\ & \vdots & \vdots \end{array}$$

My project focuses on language modelling in the context of text prediction. That is, given a sequence of words $w_1 w_2 \dots w_k = w_1^k$, I want to estimate $\mathbb{P}(w_{k+1}|w_1^k)$. For instance, if a user has typed ‘*do you want to grab a*’, then a language model could be used to suggest probable next words such as ‘*coffee*’, ‘*drink*’ or ‘*bite*’, and these predictions could further be narrowed down as the user continues typing.

1.2 Motivation

Benchmarking

Language models are central to a wide range of applications, but there are so many different ways of implementing them. Before using a language model in a particular context, it is important to understand how it performs in comparison to other methods available. In this project I focused in depth on the two most prominent types of language model: n -gram models and recurrent neural network-based models. n -gram models are typically coupled with smoothing techniques, which are explained in section 2.1. I investigated 5 different smoothing techniques and 3 different recurrent neural network architectures on a variety of datasets.

Error-prone Text

One problem with existing language models is that their next-word predictions tend to be less accurate when they are presented with error-prone text. This is not surprising, because they are only ever trained on sentences that do not contain any errors. Unfortunately, the assumption that humans will not make any mistakes when typing text is almost never valid. For this reason, I also investigated ways to narrow the gap in performance between language model predictions on error-prone text and language model predictions on error-free text.

1.3 Related Work

Chelba et al. [1] from Google explore the performance of a variety of language models on a huge, one billion word dataset. Their work presents the limits of language modelling, when vast quantities of data and computational resources are available. Chen and Goodman [2] compare the performance of a series of smoothing techniques for n -gram models, and later use their results to propose an extension to Kneser-Ney smoothing [3] which is implemented in this project.

In recent years, there have been joint efforts from Ng et al. at CoNLL to improve and compare the performance of grammatical error correction [4] [5]. Language modelling on error-prone text, however, has an important distinction: A test time, a language model cannot use words ahead of its current position to make predictions. In other words, if a language model is predicting the probability of the next word in the middle of an error-prone sentence, then it can only use the first half of the sentence for making corrections and predictions.

Chapter 2

Preparation

My preparation consisted of thoroughly understanding n -gram and RNN-based language models, as well as planning how to tie them all together in an efficient implementation.

2.1 n -gram Models

This section describes n -gram language models and the various smoothing techniques implemented in this project.

2.1.1 An Overview of n -gram Models

Language models are concerned with the task of computing $\mathbb{P}(w_1^N)$, the probability of a sequence of words $w_1 w_2 \dots w_N = w_1^N$, where $w_i \in V$ and V is some predefined vocabulary¹. By repeated application of the product rule, it follows that:

$$\mathbb{P}(w_1^N) = \prod_{i=1}^N \mathbb{P}(w_i | w_1^{i-1}) \quad (2.1)$$

n -gram language models make the Markov assumption that w_i only depends on the previous $(n-1)$ words. That is, $\mathbb{P}(w_i | w_1^{i-1}) \approx \mathbb{P}(w_i | w_{i-n+1}^{i-1})$:

$$\mathbb{P}(w_1^N) \approx \prod_{i=1}^N \mathbb{P}(w_i | w_{i-n+1}^{i-1}) \quad (2.2)$$

Using the maximum likelihood estimation, $\mathbb{P}(w_i | w_{i-n+1}^{i-1})$ can be estimated as follows:

$$\mathbb{P}(w_i | w_{i-n+1}^{i-1})_{MLE} = \frac{c(w_{i-n+1}^i)}{\sum_w c(w_{i-n+1}^{i-1} w)} \quad (2.3)$$

where $c(W)$ denotes the number of times that the word sequence W was seen in the training set.

To a first approximation, the aforementioned n -gram language model provides reasonable results and is simple to compute. However, it does have one major issue: if, as an example, a 3-gram (trigram) model does not encounter the trigram ‘*the cat sat*’ in the data

¹The vocabulary is typically taken as all of the words that occur at least k times in the training set. k is typically around 2 or 3.

it is trained upon, then it will assign a probability of 0 to that word sequence. This is problematic, because ‘the cat sat’ and many other plausible sequences of words might not occur in the training data. In fact, there are $|V|^n$ possible n -grams for a language model with vocabulary V , which is exponential in the value of n . This means that as the value of n is increased, the chances of encountering a given n -gram in the training data becomes exponentially less likely.

One way to get around this problem is to exponentially increase the size of the training set. This does, however, require significantly more memory and computation, and assumes that additional training data is available in the first place. An alternative solution is to adopt a technique called *smoothing*. The idea behind smoothing is to ‘smooth’ the probability distribution over the words in the vocabulary such that rare or unseen n -grams are given a non-zero probability. There are a variety of methods that achieve this. The ones which I have implemented are described in the next section.

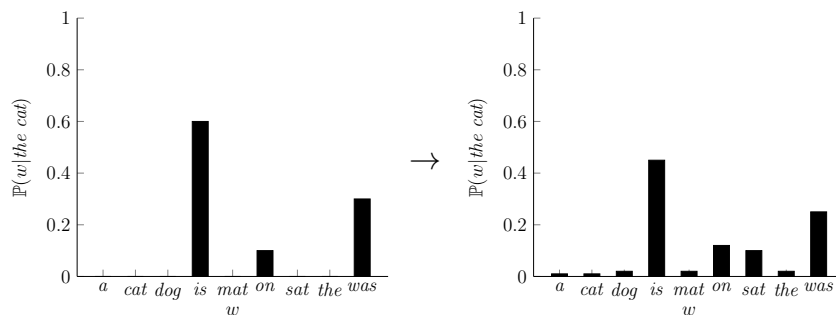


Figure 2.1: A toy example of smoothing. The probabilities of words that frequently follow ‘the cat’ are distributed to other less frequently occurring words in the vocabulary. The vocabulary is $V = \{a, cat, dog, is, mat, on, sat, the, was\}$.

2.1.2 Smoothing Techniques

Add-One Smoothing

Add-one smoothing [6] simply involves adding 1 to each of the n -gram counts, and dividing by $|V|$ to ensure the probabilities sum to 1:

$$\mathbb{P}(w_i | w_{i-n+1}^{i-1})_{\text{ADD-ONE}} = \frac{c(w_{i-n+1}^i) + 1}{\sum_w c(w_{i-n+1}^{i-1}w) + |V|} \quad (2.4)$$

One issue with add-one smoothing is that it gives an equal amount of probability to all n -grams, regardless of how likely they actually are. As an example, if both ‘the cat’ and ‘pizza cat’ are unseen in the training data of a bigram model, then $\mathbb{P}(cat | the)_{\text{ADD-ONE}} = \mathbb{P}(cat | pizza)_{\text{ADD-ONE}}$, despite the fact that ‘the’ is much more likely to precede ‘cat’ than ‘pizza’. This problem can be reduced by employing *backoff*, a technique whereby you recurse on the probability calculated by the $(n - 1)$ -gram model. In this case, it is likely that $\mathbb{P}(the)_{\text{ADD-ONE}} > \mathbb{P}(pizza)_{\text{ADD-ONE}}$, which could be used to deduce that ‘the cat’ is more likely than ‘pizza cat’.

Absolute Discounting

Absolute discounting employs backoff by interpolating higher and lower order n -gram models. It does this by subtracting a fixed discount $0 \leq D \leq 1$ from each non-zero count:

$$\begin{aligned} \mathbb{P}(w_i | w_{i-n+1}^{i-1})_{\text{ABS}} &= \frac{\max\{c(w_{i-n+1}^i) - D, 0\}}{\sum_w c(w_{i-n+1}^{i-1}w)} \\ &+ \frac{D}{\sum_w c(w_{i-n+1}^{i-1}w)} N_{1+}(w_{i-n+1}^{i-1} \bullet) \mathbb{P}(w_i | w_{i-n+2}^{i-1})_{\text{ABS}} \end{aligned} \quad (2.5)$$

where

$$N_{1+}(w_{i-n+1}^{i-1} \bullet) = |\{w \mid c(w_{i-n+1}^{i-1}w) \geq 1\}| \quad (2.6)$$

and the base case of recursion $\mathbb{P}(w)_{\text{ABS}}$ is given by the maximum likelihood unigram model. $N_{1+}(w_{i-n+1}^{i-1} \bullet)$ is the number of unique words that follow the sequence w_{i-n+1}^{i-1} , which is the number of n -grams that D is subtracted from. It is not difficult to show that the coefficient attached to the $\mathbb{P}(w_i | w_{i-n+2}^{i-1})_{\text{ABS}}$ term ensures that the probabilities sum to 1.

Ney, Essen and Kneser [7] suggested setting D to the value:

$$D = \frac{n_1}{n_1 + 2n_2} \quad (2.7)$$

where n_1 and n_2 are the total number of n -grams with 1 and 2 counts respectively.

Kneser-Ney Smoothing

Kneser and Ney proposed an extension to absolute discounting which takes into account the number of unique words that precede a given n -gram [8]. As a motivating example, consider the bigram ‘*bottle cap*’. If ‘*bottle cap*’ has never been seen in the training data, then the absolute discounting model would backoff onto the unigram distribution for ‘*cap*’. Using the unigram distribution, ‘*Francisco*’ might be given a higher probability than ‘*cap*’ (assuming ‘*Francisco*’ occurs more frequently than ‘*cap*’). This would result in the bigram ‘*bottle Francisco*’ being given a higher probability than ‘*bottle cap*’. Clearly, this is undesirable, because ‘*Francisco*’ only ever follows ‘*San*’.

From this example, it seems intuitive to assign more probability to those n -grams that follow a larger number of unique words. Kneser and Ney encapsulate this intuition by replacing some of the absolute counts $c(w_i^j)$ with the number of unique words that precede the word sequence w_i^j , $N_{1+}(\bullet w_i^j)$:

$$N_{1+}(\bullet w_i^j) = |\{w \mid c(w w_i^j) \geq 1\}| \quad (2.8)$$

Kneser-Ney smoothing² is defined as follows:

$$\mathbb{P}(w_i | w_{i-n+1}^{i-1})_{\text{KN}} = \frac{\max\{\gamma(w_{i-n+1}^i) - D, 0\}}{\sum_w \gamma(w_{i-n+1}^{i-1}w)} \quad (2.9)$$

$$+ \frac{D}{\sum_w \gamma(w_{i-n+1}^{i-1}w)} N_{1+}(w_{i-n+1}^{i-1} \bullet) \mathbb{P}(w_i | w_{i-n+2}^{i-1})_{\text{KN}} \quad (2.10)$$

²This is actually the interpolated version of Kneser-Ney smoothing, which differs slightly in form to the equation presented in the original paper.

where

$$\gamma(w_{i-k+1}^i) = \begin{cases} c(w_{i-k+1}^i) & \text{for the outermost level of recursion (i.e. } k = n) \\ N_{1+}(\bullet w_{i-k+1}^i) & \text{otherwise} \end{cases} \quad (2.11)$$

and the unigram probability is given as:

$$\mathbb{P}(w_i)_{\text{KN}} = \frac{N_{1+}(\bullet w_i)}{\sum_w N_{1+}(\bullet w)} \quad (2.12)$$

Modified Kneser-Ney Smoothing

Chen and Goodman experimented with different discount values D in Kneser-Ney smoothing and noticed that the ideal average discount value for n -grams with one or two counts is substantially different from the ideal average discount for n -grams with higher counts. Upon this discovery, they introduced a modified version of Kneser-Ney smoothing [3]:

$$\begin{aligned} \mathbb{P}(w_i|w_{i-k+1}^{i-1})_{\text{MKN}} = & \frac{\max\{\gamma(w_{i-k+1}^i) - D(c(w_{i-k+1}^i), 0)\}}{\sum_w \gamma(w_{i-k+1}^{i-1}w)} \\ & + \lambda(w_{i-k+1}^{i-1})\mathbb{P}(w_i|w_{i-k+2}^{i-1})_{\text{MKN}} \end{aligned} \quad (2.13)$$

where γ is defined in equation 2.11, and λ is defined as:

$$\lambda(w_{i-k+1}^{i-1}) = \frac{D_1 N_1(w_{i-k+1}^{i-1} \bullet) + D_2 N_2(w_{i-k+1}^{i-1} \bullet) + D_{3+} N_{3+}(w_{i-k+1}^{i-1} \bullet)}{\sum_w \gamma(w_{i-k+1}^{i-1}w)} \quad (2.14)$$

and

$$D(c) = \begin{cases} 0 & \text{if } c = 0 \\ D_1 & \text{if } c = 1 \\ D_2 & \text{if } c = 2 \\ D_{3+} & \text{if } c \geq 3 \end{cases} \quad (2.15)$$

where N_1 , N_2 and N_{3+} are defined analogously to N_{1+} in equation 2.6. Chen and Goodman suggest setting the discounts to be:

$$D_1 = 1 - 2D \frac{n_2}{n_1} \quad D_2 = 2 - 3D \frac{n_3}{n_2} \quad D_{3+} = 3 - 4D \frac{n_4}{n_3} \quad (2.16)$$

where D is as defined in equation 2.7.

Katz Smoothing

Katz smoothing is a popular smoothing technique based on the Good-Turing estimate [9]. The Good-Turing estimate states an n -gram that occurs r times should be treated as occurring r^* times, where:

$$r^* = (r + 1) \frac{n_{r+1}}{n_r} \quad (2.17)$$

where n_r is the number of n -grams that occur r times. Converting this count into a probability simply involves normalising as follows:

$$\mathbb{P}(w_i|w_{i-n+1}^{i-1})_{\text{GT}} = \frac{c^*(w_{i-n+1}^i)}{\sum_{r=0}^{\infty} n_r r^*} \quad (2.18)$$

Katz smoothing [10] is then defined as:

$$\mathbb{P}(w_i|w_{i-n+1}^{i-1})_{\text{KATZ}} = \begin{cases} \mathbb{P}(w_i|w_{i-n+1}^{i-1})_{\text{GT}} & \text{if } c(w_{i-n+1}^i) > 0 \\ \alpha(w_{i-n+1}^{i-1})\mathbb{P}(w_i|w_{i-n+2}^{i-1})_{\text{KATZ}} & \text{otherwise} \end{cases} \quad (2.19)$$

where

$$\alpha(w_{i-n+1}^{i-1}) = \frac{1 - \sum_{\{w_i \mid c(w_{i-n+1}^i) > 0\}} \mathbb{P}(w_i|w_{i-n+1}^{i-1})_{\text{KATZ}}}{1 - \sum_{\{w_i \mid c(w_{i-n+1}^i) > 0\}} \mathbb{P}(w_i|w_{i-n+2}^{i-1})_{\text{KATZ}}} \quad (2.20)$$

In practice, the infinite sum in equation 2.18 cannot be computed. To get around this issue, Katz takes n -gram counts above some threshold k as reliable and only applies the Good-Turing estimate to those with a count less than or equal to k . Katz suggests $k = 5$. This modification requires a slightly different equation to 2.19 and is presented in Katz's original paper.

2.2 Recurrent Neural Network Models

In this section I give a brief introduction to neural networks, recurrent neural networks (RNNs) and how RNNs can be used in the context of language modelling. A thorough description of the more complex RNN architectures, gated recurrent units (GRUs) and long short-term memory (LSTM), is given in chapter 3.

2.2.1 An Overview of Neural Networks

The human brain is a furiously complicated organ, packed with a network of approximately 86 billion neurons³ that propagate electrochemical signals across connections called synapses. Artificial neural networks, or neural networks, were originally developed as a mathematical model of the brain [12], which despite being substantially oversimplified, now provides an effective tool for classification and regression in modern-day machine learning.

Neural networks consist of a series of nodes which are joined by directed and weighted connections. Inputs are supplied to some subset of the nodes and then propagated along the weighted connections until they reach the designated output nodes. In the context of the brain, the nodes represent neurons, the weighted connections represent synapses and the flow of information represents electrochemical signals.

An important distinction to be made is whether the neural network is cyclic or not. Acyclic neural networks are called feed-forward neural networks (FNNs), whereas cyclic neural networks are denoted recurrent neural networks (RNNs) which are covered in section 2.2.2. There are a variety of FNNs, but the most prominent is the multilayer perceptron (MLP) [13], which I will outline below.

³According to a study by Azevedo et al. [11].

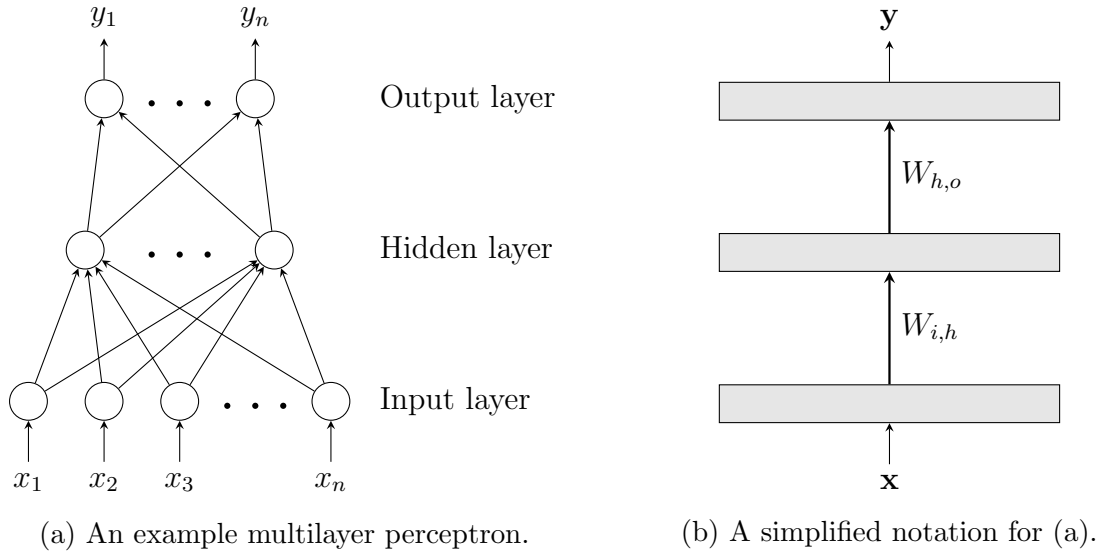


Figure 2.2: Two alternative notations for the same network: a multilayer perceptron with one hidden layer. In the simplified notation, the grey rectangle represents a vector of neurons, the thin arrow represents a vector and the thick arrow represents a matrix of weights.

The Multilayer Perceptron

The multilayer perceptron consists of layers of neurons, where each layer is fully connected to the next one. The first layer is the input layer, the last is the output layer and any layers in between are called *hidden layers*. If there is more than one hidden layer then it is called a *deep* neural network. The input neurons simply pass on the input values that they are given. The neurons in the subsequent layers of the network compute the weighted sum of their inputs, before passing that value through an *activation function* and outputting it:

$$\begin{array}{c}
 1 \\
 a_1^{l-1} \\
 \vdots \\
 a_n^{l-1}
 \end{array}
 \begin{array}{c}
 b_j^{l-1} \\
 w_{1,j}^{l-1} \\
 w_{n,j}^{l-1}
 \end{array}
 \rightarrow j \rightarrow \underbrace{\phi\left(b_j^{l-1} + \sum_i a_i^{l-1} w_{i,j}^{l-1}\right)}_{z_j^l - \text{the logit}}$$

a_j^l - the activation

Figure 2.3: The computation carried out by each neuron. The outputs or *activations* of the neurons from the previous layer, $(l-1)$, are denoted a_i^{l-1} , the weights on the input connections are $w_{i,j}^{l-1}$, the activation function is ϕ and b_j^{l-1} is the bias variable.

The unusual looking constant-valued input 1 is known as a *bias* input, which is weighted by a bias variable b_j^{l-1} . The purpose of the bias is to allow the neuron to shift the output of the activation function left and right by adjusting the value of b_j^{l-1} . Bias inputs are usually attached to every non-input neuron in the network.

It is also worth noting that activation functions should be chosen to be non-linear. Any

combination of linear operators is linear, which means that any linear MLP with multiple hidden layers is equivalent to an MLP with a single hidden layer. Non-linear neural networks, on the other hand, are more powerful, and can gain considerable performance by adding successive hidden layers to re-represent the input data at higher levels of abstraction [14] [15]. In fact, it has been shown that a non-linear MLP with a single hidden layer containing a sufficient number of neurons can approximate any continuous function on a compact input domain to arbitrary precision [16].

Frequently used activation functions include the sigmoid and the hyperbolic tangent functions. These are both non-linear functions squashed within the ranges $(0, 1)$ and $(-1, 1)$ respectively. Their steep slope at the origin is supposed to mimic an axon, which is a neuron that fires its output after the input reaches a certain potential. More importantly, these functions are differentiable, which allows for the network to be trained using gradient descent. This is discussed below.

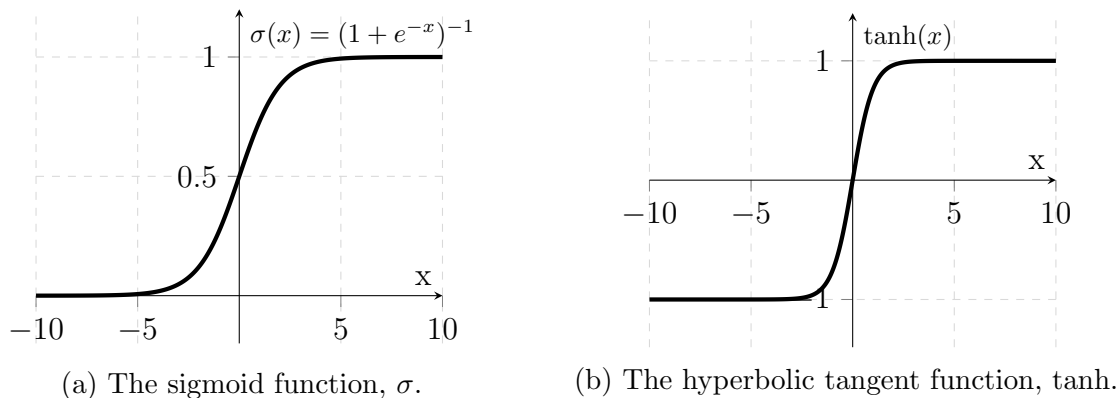


Figure 2.4: Two popular activation functions.

Backpropagation and Gradient Descent

FNNs compute a function, f , parameterised by the weights \mathbf{w} of the network, mapping an input vector $\mathbf{x} = (x_1, \dots, x_n)^T$ to an output vector $\mathbf{y} = (y_1, \dots, y_m)^T$.

$$\mathbf{y} = f_{\mathbf{w}}(\mathbf{x})$$

The entire point of training a neural network is to get it to learn a particular mapping from input vectors to output vectors. What the inputs and outputs represent depends on the problem at hand. For example, in the context of classifying pictures of animals, the input vector might represent the pixel values of an image and the output vector might represent a probability distribution over the set of animals in the classification task.

In order to train a neural network, you supply it with a training set, which is just a list of input-target pairs:

$$\mathbf{s} = ((\mathbf{x}_1, \mathbf{t}_1), \dots, (\mathbf{x}_N, \mathbf{t}_N))$$

Where \mathbf{t}_k is the vector that the network should output given the example input \mathbf{x}_k . A differentiable loss function, $\mathcal{L}(\mathbf{y}, \mathbf{t})$, is also defined, which essentially says how badly the

network output \mathbf{y} matches the target output \mathbf{t} . Then, a measure of how much error the network produces over the whole training set can be defined as follows:

$$\mathcal{L}_{\text{TOTAL}} = \sum_{k=1}^N \mathcal{L}(f_{\mathbf{w}}(\mathbf{x}_k), \mathbf{t}_k) \quad (2.21)$$

The end product of the *backpropagation* algorithm is the partial derivative:

$$\frac{\partial}{\partial w_{i,j}^l}(\mathcal{L}_{\text{TOTAL}})$$

for every weight $w_{i,j}^l$ in the neural network. It is a two stage algorithm that works as follows:

1. Calculate, $f_{\mathbf{w}}(\mathbf{x}_k)$ for each input example \mathbf{x}_k , and use those values to derive $\mathcal{L}_{\text{TOTAL}}$. The process of calculating the output of the neural network given the input is known as *forward propagation*, because the input is propagated through each layer in the network.
2. Given $\mathcal{L}_{\text{TOTAL}}$, calculate $\frac{\partial}{\partial w_{i,j}^l}(\mathcal{L}_{\text{TOTAL}})$ for each weight in each layer of the network by applying the chain rule backwards from the output layer to the input layer. This step is described in more detail below.

$\frac{\partial}{\partial w_{i,j}^l}(\mathcal{L}_{\text{TOTAL}})$ is calculated as follows:

$$\frac{\partial}{\partial w_{i,j}^l}(\mathcal{L}_{\text{TOTAL}}) = \sum_{k=1}^N \underbrace{\frac{\partial \mathcal{L}}{\partial f_{\mathbf{w}}(\mathbf{x}_k)} \frac{\partial f_{\mathbf{w}}(\mathbf{x}_k)}{\partial z_j^{l+1}}}_{\delta_j^{l+1}} \frac{\partial z_j^{l+1}}{\partial w_{i,j}^l} = \delta_j^{l+1} a_i^l \quad (2.22)$$

The term δ_j^{l+1} varies in form for each layer. For the output layer, it has the following form:

$$\delta_j^L = \frac{\partial \mathcal{L}}{\partial f_{\mathbf{w}}(\mathbf{x}_k)} \phi'(z_j^L) \quad (2.23)$$

For the preceding layers, δ_j^l can be defined recursively in terms of the delta values for the neurons in subsequent layers in the network. This is defined below, and a detailed justification of it is given in appendix A:

$$\delta_i^l = \phi'(z_i^l) \sum_j \delta_j^{l+1} w_{i,j}^l \quad (2.24)$$

Gradient descent is an optimisation technique that uses the derivatives produced by the backpropagation algorithm to adjust the weights such that the loss is minimised over the training set. The simplest form of gradient descent does this by changing each weight value in the direction of the negative gradient of the loss. This is equivalent to moving the network output downwards on the error surface defined by plotting the loss against each of the weights in the network:

$$w_{i,j}^l \leftarrow w_{i,j}^l - \eta \frac{\partial}{\partial w_{i,j}^l}(\mathcal{L}_{\text{TOTAL}}) \quad (2.25)$$

η is known as the *learning rate*, which determines how large the steps are in the negative direction of the gradient. Ways for setting η , and other useful techniques when training neural networks in practice are described in chapter 3.

2.2.2 Recurrent Neural Networks

In the context of language modelling, the input is a sequence of words and the output should be the same sequence of words shifted ahead in time by one word. For example, given ('the', 'cat', 'sat', ...) as input, the network should produce something similar to ('cat', 'sat', 'on', ...) as output. The problem with FNNs is that these sequences may have a varying number of words in them, yet FNNs have fixed input and output vector sizes. One way to get around this problem is to look at a finite window of words at any given time, and encode that window into a vector, but even for modest window sizes this can lead to a huge number of weights in the network, which becomes more difficult to train. A better solution is to use a recurrent neural network.

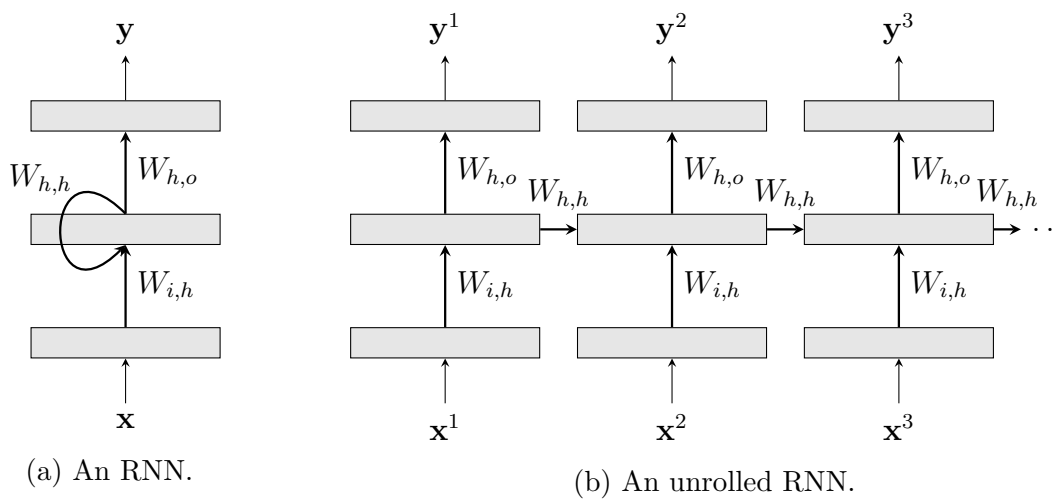


Figure 2.5: Two different representations of the same RNN.

A recurrent neural network can be constructed from an FNN by adding connections from each neuron in a hidden layer back to every neuron in that layer, including itself. This gives something like the network shown in figure 2.5a. An easier way to visualise an RNN is by *unrolling* it as shown in figure 2.5b. A nice property about RNNs is that they can be unrolled an arbitrary number of steps without the need for any additional parameters, because the same weight matrices are reused at each step. For this reason, RNNs are well suited for operating over sequences of input data. In the case of language modelling, this is ideal, because the inputs are sequences of words. The problem of representing textual words as numerical input vectors for an RNN is addressed in section 2.2.3.

As shown in figure 2.5b, RNNs can take a sequence of vectors as input. In most applications, the vectors in such sequences are typically ordered by time, so I use the notation \mathbf{x}^t and \mathbf{y}^t to denote the input and output vectors at time step t .

Backpropagation Through Time (BPTT)

Training an RNN is not so much different to training an FNN: the only difference is that gradients must also be calculated across time steps. As shown figure 2.5b, each weight is reused at each time step. For this reason, the weight update given in equation 2.22 is

summed across the time steps:

$$\frac{\partial}{\partial w_{i,j}^l}(\mathcal{L}_{\text{TOTAL}}) = \sum_t \delta_j^{l+1,t} a_i^{l,t}$$

The delta term for the output layer is defined analogously to equation 2.26:

$$\delta_j^{L,t} = \frac{\partial \mathcal{L}}{\partial f_{\mathbf{w}}(\mathbf{x}_k^t)} \phi'(z_j^{L,t}) \quad (2.26)$$

The delta term for the preceding layers can be shown to have the form below. This is justified in more detail in appendix A.

$$\delta_i^{l,t} = \phi'(z_i^{l,t}) \left(\sum_j \delta_j^{l+1,t} w_{i,j}^l + \sum_j \delta_j^{l,t+1} v_{i,j}^l \right) \quad (2.27)$$

where $v_{i,j}^l$ denotes the weight on the recurrent connection from neuron i in layer l back to neuron j in the same layer.

RNN Architectures

So far, the equations presented have assumed a *vanilla RNN*. That is, each neuron, or *cell*, simply takes the activations from the previous layer, $a_i^{l-1,t}$, and the activations from the same layer at the previous time step, $a_i^{l,t-1}$, and applies an activation function to their weighted sum to give $a_j^{l,t}$:

$$a_j^{l,t} = \phi \left(\sum_i a_i^{l-1,t} w_{i,j}^{l-1} + \sum_i a_i^{l,t-1} v_{i,j}^{l-1} + b_j^{l-1} \right) \quad (2.28)$$

This computation can be written for an entire layer of hidden neurons, $\mathbf{a}^{l,t} = (a_1^{l,t}, a_2^{l,t}, \dots, a_H^{l,t})$, using matrix notation as follows:

$$\mathbf{a}^{l,t} = \phi(\mathbf{W}^{l-1} \mathbf{a}^{l-1,t} + \mathbf{V}^{l-1} \mathbf{a}^{l,t-1} + \mathbf{b}^{l-1}) \quad (2.29)$$

Whilst this architecture provides a good baseline for RNNs, it struggles to maintain long-term information for reasons discussed in section 3.3.3. Various modifications to the way in which the cells are wired up have been proposed to improve upon this problem. In this project, I have implemented the two most prominent alternatives, which are Long Short-Term Memory (LSTM) [17] and the Gated Recurrent Unit (GRU) [18]. These architectures are outlined in figure 2.7 and discussed in much more detail in section 3.3.3.

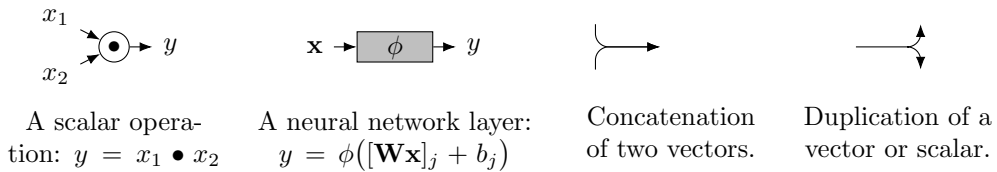


Figure 2.6: A notation for drawing RNN architectures.

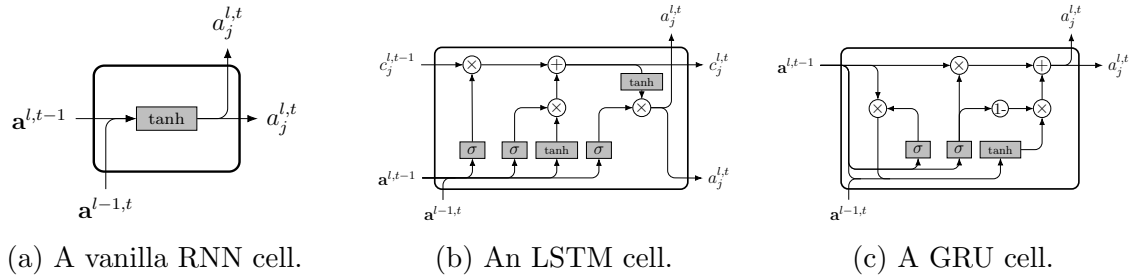


Figure 2.7: Various architectures of RNN neurons, using the notation from figure 2.6. Each diagram represents the computation of the j^{th} neuron or cell in layer l .

2.2.3 Word Embeddings

RNNs take a sequence of vectors as input, and produce a sequence of vectors as output. In the context of language modelling, these inputs and outputs should represent sequences of words. In order to do this, each word is represented by a single vector, which is known as a word *embedding*.

The aim behind word embeddings is that similar words should be represented by vectors that are closer together in the vector space than those for dissimilar words. One effective way to obtain such embeddings is to learn them through backpropagation, which is the method that was used in this dissertation. The idea is that each word is mapped to a large, randomly initialised vector. Then, when the RNN is trained on sequences of words, the embedding values are updated in a way that is analogous to how the weights are updated in the network.

2.3 Software Engineering

In this section details the requirements of the project and the early design decisions that were made to target those goals.

2.3.1 Requirements

The success requirements set out in the project proposal were:

R1	Language models (LMs) using the following techniques are implemented:
	• n -gram LMs with various smoothing techniques.
	• A vanilla RNN-based LM. • An LSTM-based LM.
R2	Comprehensible and reliable comparisons between the various LM implementations and their combinations are made regarding their accuracy, speed and resource consumption during both training and inference.
R3	A simple console application is developed to demonstrate the capability of the aforementioned language models in the context of next-word prediction.

The project proposal also suggested two possible extensions:

E1	Explore possible extensions to existing language models to improve their performance on error-prone text.
E2	Build a mobile keyboard that exhibits the functionality of the language models.

All requirements and extensions depend on **R1**, so I implemented that first. Afterwards, I built **R2** and **R3** in parallel before moving onto **E1** and **E2**.

2.3.2 Tools and Technologies Used

Below I describe and justify where necessary the tools and technologies that I selected to build my project.

Version Control and Build Tools

I hosted my project in a repository on GitHub, used Git for version control, and used Bazel for running builds and tests on my project.

Machine Learning

I chose to use TensorFlow, an open source machine learning library, to assist in the implementation of my recurrent neural networks. There are two key reasons for this:

1. TensorFlow supports *automatic differentiation*. That is, each tensor operation has a partial derivative associated with it, so when you want to apply the second stage of backpropagation, TensorFlow can automatically apply the chain rule through any sequence of tensor operations. Without this, I would have to hard-code new gradient calculations every time I change the structure of the network, which is tedious and error-prone.
2. TensorFlow supports GPU-accelerated matrix operations. Without this, training my recurrent neural networks would take significantly longer.

Of course, there are many other machine learning libraries that offer comparable features. The reason I picked TensorFlow in particular was because I was already familiar with its codebase from using and contributing to it during an internship at Google in 2016.

In order to train large scale models on NVIDIA GPUs, I used the University's High Performance Computing service.

Languages

At the time of writing my code, TensorFlow had both a C++ and a Python API. However, the C++ API only had support for loading and running inference on models, not for training them.

One of the possible extensions I set out in my project proposal was to implement a mobile keyboard that makes use of my language model implementations. Android and iOS both provide little or no support for applications written in Python, but they do support code written in C++ via the JNI and Objective-C++ respectively.

I used C++ to write the benchmarking framework and all n -gram language models, so that they could easily be used within a mobile application. As for the RNN-based language models, I wrote code to train and export them in Python, and then wrote some C++ classes for loading and running inference on them. This way, I could leverage the power of TensorFlow’s Python API whilst also supporting the use of such models in a mobile keyboard. In order to export and import trained language models across different languages, I used Google’s protocol buffers, a language neutral mechanism for serialising structured data.

Testing

I used Google Test for C++ unit tests and the `unittest` package for testing in Python.

2.3.3 Starting Point

My project codebase was written from scratch, with the assistance of the tools and libraries mentioned in section 2.3.2. Before the project, I had no understanding of language modelling. Apart from a basic knowledge of recurrent neural networks, I had to learn about language modelling, n -gram smoothing techniques, long short-term memory and gated recurrent units through personal reading.

Experience	Tools and Technologies
<i>Significant</i>	C++, Python, Git, GitHub, TensorFlow
<i>Some</i>	Bazel, iOS, Protocol Buffers
<i>None</i>	Google Test, Objective-C++, SLURM

Figure 2.8: A summary of my prior experience with the tools and technologies used in this project.

Chapter 3

Implementation

This chapter details how I implemented a variety n -gram and RNN-based language models, how I built the mobile keyboard on iOS and how I extended an existing language model to improve its performance on error-prone text.

3.1 System Overview

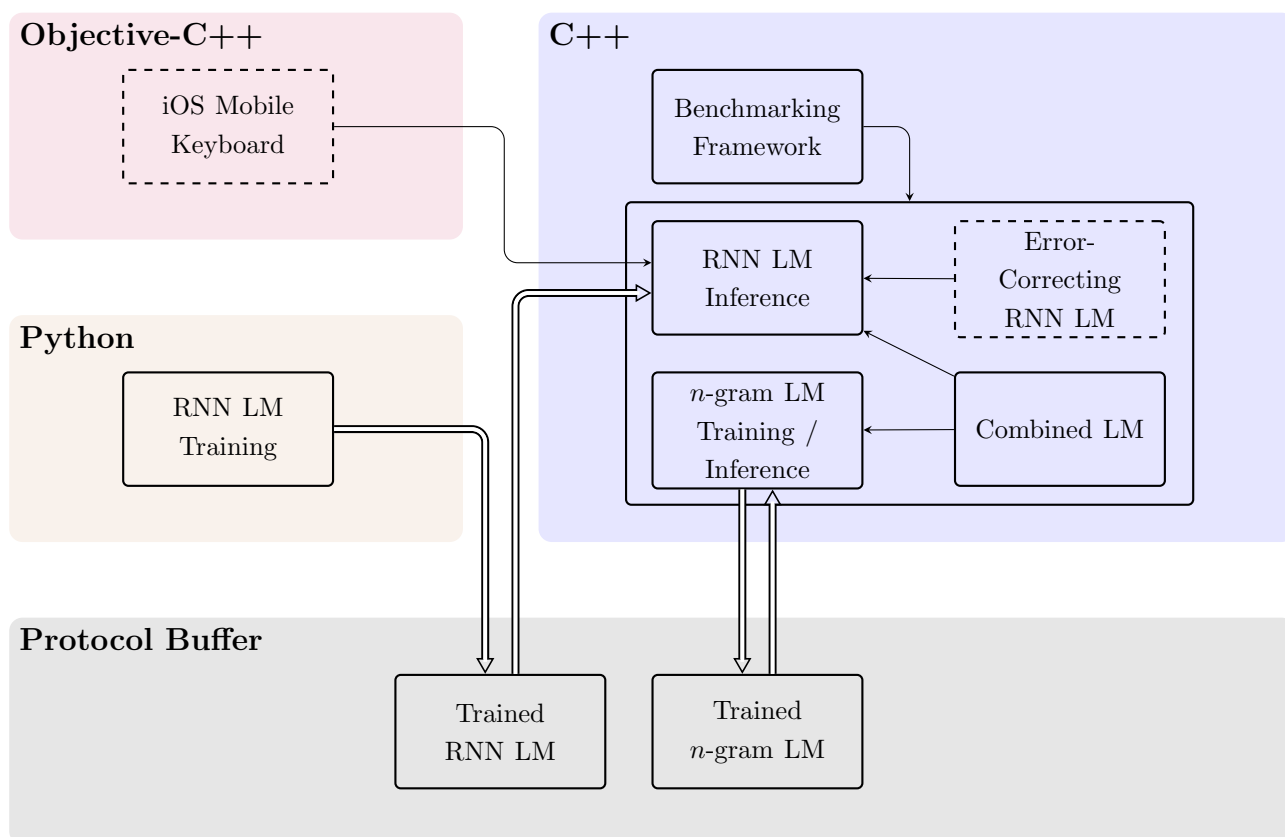


Figure 3.1: An overview of the project. Outlined rectangles denote groups of related files or code, of which the ones with a dashed outline represent project extensions. Thin arrows denote dependencies and thick arrows represent the movement of data.

The core components of my project are highlighted in figure 3.1. The choices of language for each of the system components are reasoned in section 2.3.2. I used Protocol Buffers

as a language neutral method of encapsulating a trained and exported language model that can be loaded into memory and used to make predictions.

As shown, all of the code for training and running inference on n -gram models was written in C++, and is described in section 3.2. The RNN-based language models were trained in Python using TensorFlow, before being exported into a small group of protocol buffers that can be loaded into C++ for running inference, all of which is detailed in section 3.3. The ‘Combined LM’ component represents a language model which interpolates the probabilities across one or more n -gram and/or RNN-based language models. This part is not described in depth, because it is mostly trivial. I outline in section 3.4 how I built a mobile keyboard in iOS that utilises one of the RNN-based language models. Finally, the ‘Error-Correcting RNN LM’ component denotes an extension to an existing RNN-based language model which I made in order to improve its performance on error-prone text, which is described in section 3.5.

The n -gram, RNN, error-correcting RNN and combined language models all inherit, either directly or indirectly, from a single language model superclass, which is all that the benchmarking framework depends on. This class is described briefly below. Details of the metrics used in the benchmarking framework are left for chapter 4.

3.1.1 The Language Model Hierarchy

All of my language model implementations in C++ were joined by the same inheritance hierarchy, at the top of which was a single language model superclass, LM:

```

1 class LM {
2 protected:
3     Vocab *vocab;
4     ...
5 public:
6     ...
7     // Given a sequence of words, 'seq', return the most probable next word, with its
8     // associated probability.
9     virtual void Predict(std::list<std::string> seq,
10                        std::pair<std::string, double> &pred);
11
12     // Given a sequence of words, 'seq', return the top k most probable next words, in
13     // order of decreasing probability.
14     virtual void PredictTopK(std::list<std::string> seq,
15                             std::list<std::pair<std::string, double>> &preds, int k);
16
17     // Return the probability of the word sequence, 'seq'.
18     virtual double Prob (std::list<std::string> seq) = 0;
19
20     // Given a sequence of words, 'seq', initialise 'probs' with the probability of each
21     // word in the vocabulary following that sequence.
22     virtual void ProbAllFollowing (std::list<std::string> seq,
23                                   std::list<std::pair<std::string, double>> &probs) = 0;
24
25     // The same as above, except the word-probabilities are stored in a character trie.
26     virtual void ProbAllFollowing (std::list<std::string> seq, CharTrie *probs) = 0;
27     ...
28 };

```

Every language model object stores a `Vocab` object, which contains the set of words in the language model vocabulary, and a mapping from those words to integer IDs. The integer IDs are used in place of the string representations of such words in various language model data structures that are discussed later in this chapter. This is simply a memory optimisation, because an integer requires less space than a string.

`Predict()` is a special case of `PredictTopK()`, which simply uses a priority queue to extract the top k predictions from result of calling `ProbAllFollowing()`. The implementations of `Prob()` and `ProbAllFollowing()` depend on the specific language model. One might think that `ProbAllFollowing()` can be implemented by calling `Prob()` for each word in the vocabulary, but this would be inefficient in the case of an RNN-based language model, which implicitly computes the probability distribution over every word in the vocabulary in each forward pass of the network, hence why the methods are kept separate.

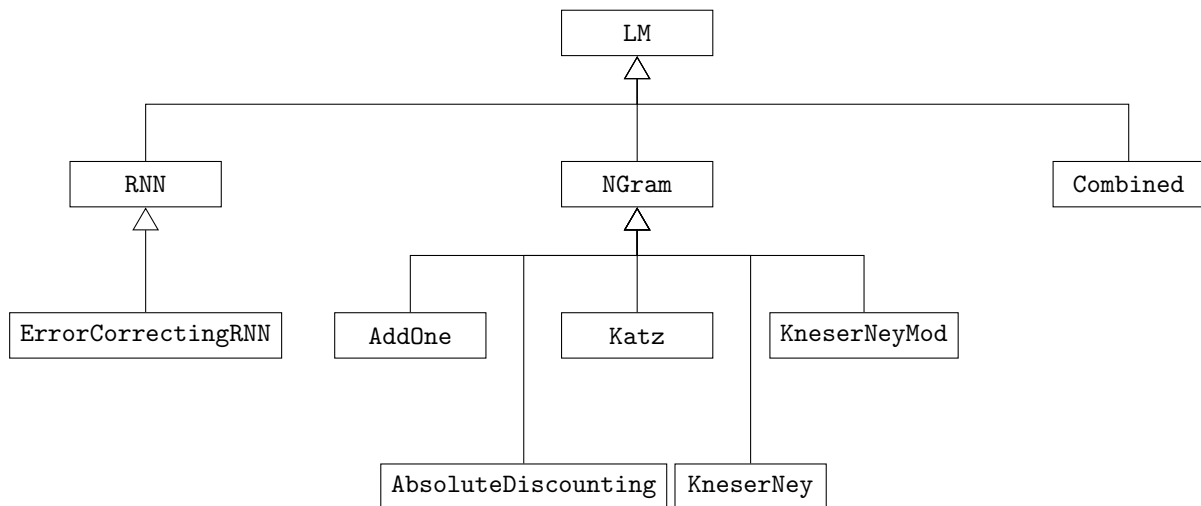


Figure 3.2: An outline of the C++ classes in the language model inheritance hierarchy. Methods and members are omitted for simplicity.

The reason that there is only a single RNN language model class is that the differences between the vanilla RNN, GRU and LSTM models arise in the Python code where the models are constructed and trained. The C++ code for the RNN-based models just loads and runs inference on the model graph described in the exported protocol buffer. See 3.3.1 for more details.

3.1.2 Testing

The correctness of my language model implementations in this project were always vital to its success. For this reason, every C++ and Python in the codebase that exhibits core language model functionality has an accompanying test class which contains one or more unit tests. In particular, the language model classes are paired with tests that query their `Prob()` and `ProbAllFollowing()` methods and compare the results with values that have been calculated by hand using the equations presented in chapter 2.

During the development of the project, all new code was tested before being committed to the codebase, and new test cases were immediately added upon discovery of any bugs.

3.2 n -gram Models

Below I detail my implementation of the n -gram models and their associated smoothing techniques.

One of the first things I discovered when implementing n -gram language models is that there is a tradeoff between how much computation is done at training time, and how much computation is done at inference time. For example, at one extreme you could compute only the number of times each n -gram occurs at training time, and leave the computation of the coefficients in the smoothing equations to be calculated on the fly. At the other extreme, you could compute and store the final value of the probability for every possible n -gram that can be formed from the words in the language model vocabulary, and then return the relevant probability immediately at inference time. The problem with the first approach is that the language model would be too slow at inference time, and the problem with the second approach is that storing absolute probabilities for every possible n -gram would occupy far too much space to be practical. In my implementation, I try to strike a balance between these two extremes by precomputing and efficiently storing all of the relevant coefficients for the smoothing equations, and then combining them at inference time to compute the required n -gram probability.

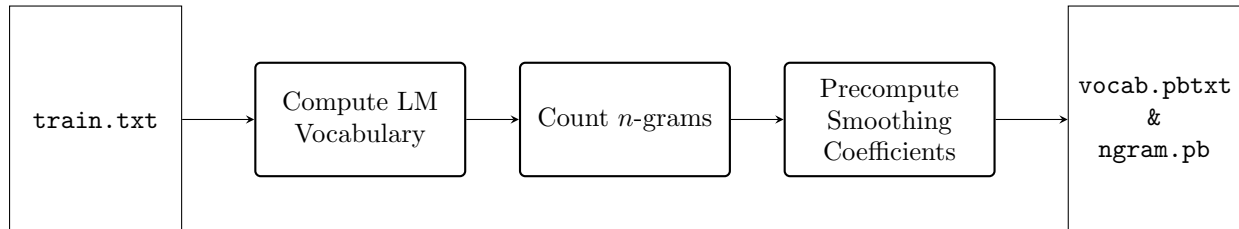


Figure 3.3: The n -gram training pipeline.

The process of training an n -gram model in my implementation is outlined in figure 3.3. The three stages of the pipeline are outlined in sections 3.2.1, 3.2.2 and 3.2.3 respectively.

3.2.1 Computing the Vocabulary

The first stage involves a single pass over `train.txt`, where words that occur at least k times are added to the vocabulary and mapped to an integer ID. A special marker `<unk>` is included in the vocabulary to represent any out-of-vocabulary words. That is, when the model encounters a word it has not seen before, it treats it as `<unk>`. The marker `<s>` is also included and is used to denote the end of a sentence.

3.2.2 Counting n -grams Efficiently

For n -gram models without any smoothing, the only information that needs to be computed at this stage is $c(w_{i-n+1}^i)$ and $\sum_w c(w_{i-n+1}^{i-1}w)$, which are the number of times each

n -gram w_{i-n+1}^i occurred in `train.txt`, and the sum of all such counts for all words that follow the $(n-1)$ -gram w_{i-n+1}^{i-1} respectively. Then, given these two precomputed numbers, when the n -gram model is queried for the probability of w_{i-n+1}^i , it can simply return the first number divided by the second, as required in equation 2.3.

This situation becomes more complex once smoothing techniques are included. For instance, absolute discounting requires the knowledge of $N_{1+}(w_{i-n+1}^{i-1} \bullet)$, the number of unique words that follow w_{i-n+1}^{i-1} , for every such $(n-1)$ -gram. Moreover, because it employs backoff, it actually requires this knowledge for every k -gram where $1 \leq k < n$. Kneser-Ney smoothing additionally requires $N_{1+}(\bullet w_{i-k+1}^i)$, the number of unique words that precede w_{i-k+1}^i , and $\sum_w N_{1+}(\bullet w_{i-k+1}^i w)$, the number of unique pairs of words that follow and precede w_{i-k+1}^i , for every k -gram where $1 \leq k < n$.

In my implementation, I precompute and store all of these parameters in a temporary word-level trie, as shown in figure 3.4. This way, each parameter is only ever computed once and is accessible in $O(n)$ time. Given that I only consider n -grams for $n = 2, 3, 4$ or 5, the access time practically constant.

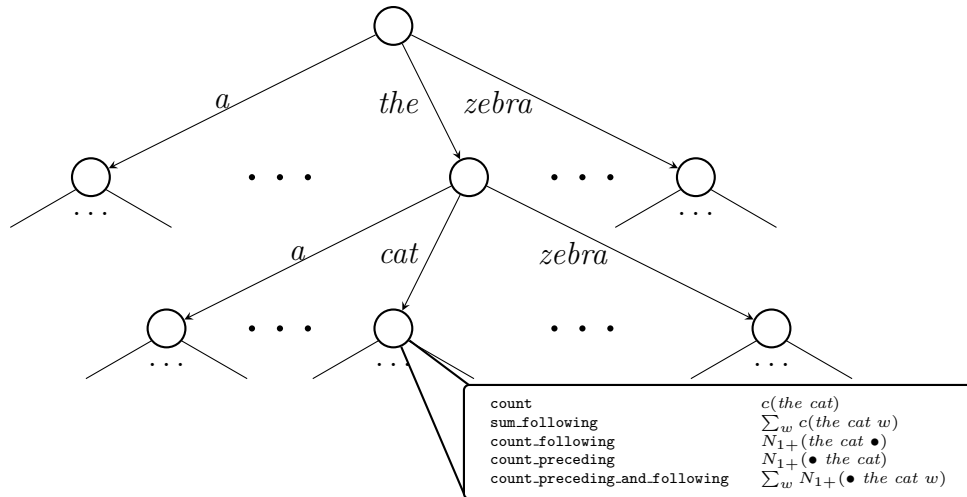


Figure 3.4: An example instance of the `CountTrie` class, which is the end result of the ‘Count n -grams’ stage in figure 3.3.

Populating a `CountTrie` instance works in two phases. In the first phase, all of the k -grams, where $1 \leq k \leq n$, from the training set are inserted into the trie. At every insertion, the relevant `count` and `count_preceding` values are updated. The second phase involves a full traversal of the trie in which the `sum_following`, `count_following` and `count_following_and_preceding` are computed. Each of these values can be derived for a particular node by looking at the `count` or `count_preceding` values of its children.

It is worth noting that this data structure is temporary. That is, it is only used to enable fast computation of the smoothing coefficients in the next stage.

3.2.3 Precomputing Smoothing Coefficients

With a few exceptions, most smoothing techniques can be generalised to the following form:

$$\mathbb{P}(w_i | w_{i-n+1}^{i-1})_{\text{SMOOTH}} = \alpha(w_{i-n+1}^i) + \beta(w_{i-n+1}^{i-1}) \mathbb{P}(w_i | w_{i-n+2}^{i-1})_{\text{SMOOTH}} \quad (3.1)$$

The probabilities of many different n -grams could share some of the same α and β coefficients. For example, $\mathbb{P}(\text{sat} \mid \text{the cat})_{\text{SMOOTH}}$ and $\mathbb{P}(\text{jumped} \mid \text{the cat})_{\text{SMOOTH}}$ would both require $\beta(\text{the cat})$. To avoid duplicated computation, the α and β values for every k -gram in the **CountTrie** instance are computed exactly once and stored in a new word-level trie. Each node in this trie stores two values, $\alpha(w_i^j)$ and $\beta(w_i^j)$, where w_i^j is the path of words that lead to the node. This way, the parameters can be fetched in constant time at n -gram model query time. The **CountTrie** instance from the previous stage indicate which k -grams the coefficients should be computed for, and also supplies the necessary parameters to derive them.

Once the word-level trie of smoothing coefficients has been populated, computing the probability of any n -gram can be done by traversing the data structure as shown in algorithm 1. Using this implementation, the only thing that each smoothing subclass has to override is the way in which the coefficients α and β are calculated when the data structure is initially populated. This made it much easier to add and test new smoothing techniques.

Algorithm 1 Computing $\mathbb{P}(w_i | w_{i-n+1}^{i-1})_{\text{SMOOTH}}$

```

1: procedure PROB( $w_i, w_{i-n+1}^{i-1}$ )
2:   if  $\text{length}(w_{i-n+1}^i) > 0$  then
3:      $\text{node} \leftarrow \text{GETNODE}(w_{i-n+1}^i)$ 
4:      $\alpha \leftarrow 0$ 
5:      $\beta \leftarrow 1$ 
6:     if  $\text{node}$  then
7:        $\alpha \leftarrow \text{node}.\alpha$ 
8:       if  $\text{length}(w_{i-n+1}^{i-1}) > 0$  then
9:          $\text{b\_node} \leftarrow \text{GETNODE}(w_{i-n+1}^{i-1})$ 
10:        if  $\text{b\_node}$  then
11:           $\beta \leftarrow \text{b\_node}.\beta$ 
12:        return  $\alpha + (\beta \cdot \text{PROB}(w_i, w_{i-n+2}^{i-1}))$ 
13:   return 0

```

Katz smoothing is an exception in that it does not quite match the form presented in equation 3.1. To get around this issue, the Katz subclass simply overrode algorithm 1 to traverse the same data structure in a slightly different way.

3.3 Recurrent Neural Network Models

Below, I explain how I implemented my RNN-based language models. To do this, I first give a brief overview of TensorFlow in section 3.3.1. Then, in section 3.3.2, I explain the structure of the RNNs I use and how they apply to language modelling. In section 3.3.3

I give the details of the GRU and LSTM architectures. Section 3.3.4 details the word embeddings used in the models. Finally, section 3.3.5 describes how I tuned the RNN hyperparameters and applied techniques to encourage convergence and avoid overfitting.

3.3.1 TensorFlow

In this section I aim to give a brief overview of TensorFlow [19], the open-source machine learning library that I used to implement my RNN-based language models.

Computation in TensorFlow is defined in terms of a directed graph, where each node represents an operation and each directed edge represents the flow of a tensor, hence the name. TensorFlow programs consist of two key sections: the first is constructing the graph, and the second is executing one or more paths through the graph in what is called a ‘session’. In order to execute a path through a graph, you specify one or more nodes to evaluate, and then TensorFlow automatically executes all operations that are required to evaluate the tensor at the specified node(s).

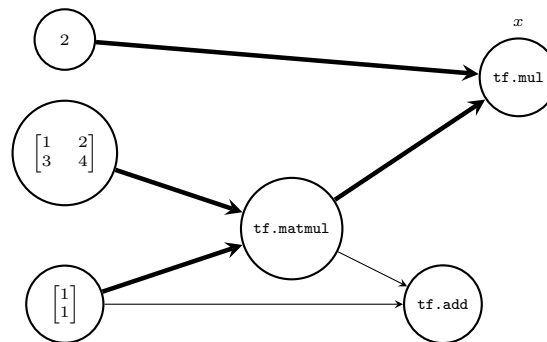


Figure 3.5: The toy TensorFlow computation graph. On evaluation of node x , only the operations on the highlighted paths are computed, and the scalar (rank-0 tensor) 20 is returned.

Each operation takes zero or more tensors as input, and produces a tensor as output. There are some notable operations that take zero inputs: a *constant*, which outputs a tensor whose value cannot be changed; a *variable*, which is the same except the tensor that it outputs can be changed; and a *placeholder*, which is a promise to provide a tensor at execution time.

The computation of a (recurrent) neural network can be described as a directed graph of tensor operations. For example: a layer of neurons can be represented as a vector and the weights in between each layer can be represented as a matrix. Then, computing the activations of the next layer is simply a case of matrix multiplying the weights with the activations from the previous layer, and applying an activation function element-wise to the result. In the context of TensorFlow, a placeholder would be used to provide the input training data, and variables would be used to store all of the trainable parameters of the network, namely the weights, bias variables and the word embeddings.

Once a neural network has been constructed, forward propagation can be computed by evaluating the output nodes of the network in a session. TensorFlow has a useful method

called `tf.gradients(ys, xs)`, which returns the partial derivatives of the sum of the tensor(s) in `ys` with respect to the variables in `xs`, and is known as *automatic differentiation*. By setting `xs` to the trainable parameters of the network and `ys` to the loss function, I avoided having to hardcode the backpropagation computation. Using these partial derivatives, any optimisation technique such as gradient descent can be applied to update the network.

After having trained the RNN in Python, I would *freeze* the computation graph and export it into a format that can be loaded into C++. This involved converting all of the variables in the graph to constants and saving the network, along with some metadata, into a protocol buffer. Using the C++ API, I could then load the network and run inference on it by evaluating the relevant tensors in a session, without having to know anything about the internal structure of the RNN.

3.3.2 Network Structure

The general structure of all of the RNN-based language models I trained is described below.

Before the RNN is constructed, an initial pass is made over the training set to construct a vocabulary V . This is done in exactly the same way as for the n -gram models, as described in section 3.2.1, so I will not repeat the detail.

At the very start of the network is a word embedding lookup table, that maps words to vectors that represent points in a high dimensional space. In my implementation, I set the number of dimensions of this space (that is, the number of vector elements in each word embedding) to be the same as the number of hidden neurons in each hidden layer of the RNN. This is a configurable parameter, which was typically set to 256.

After the embedding lookup is the input layer of the RNN into which the word embeddings are fed. There are then two hidden layers which are recurrently connected to themselves, followed by the output layer.

There are $|V|$ neurons in the output layer, each of which represent a word in the language model vocabulary. Rather than feeding the logits of the output neurons through an activation function, they are passed through the softmax function. That is, if the vector of logits of the output neurons is denoted \mathbf{z} , then the output of neuron j after the softmax function is defined as $\text{SOFTMAX}(\mathbf{z})_j$, where:

$$\text{SOFTMAX}(\mathbf{z})_j = \frac{e^{z_j}}{\sum_i e^{z_i}} \quad (3.2)$$

The outputs of the softmax function sum to 1 and are used to represent a probability distribution over all of the possible next words. That is, $\mathbb{P}(\text{next word} = j) = \text{SOFTMAX}(\mathbf{z})_j$, where j is the ID of a particular word in the vocabulary. Each ID lies in the range $[0, |V|)$.

3.3.3 Long Short-Term Memory and Gated Recurrent Units

As mentioned in section 2.2.2, there are various architectures that can be used for the cells/neurons of RNNs, three of which I implemented in this project. The vanilla RNN

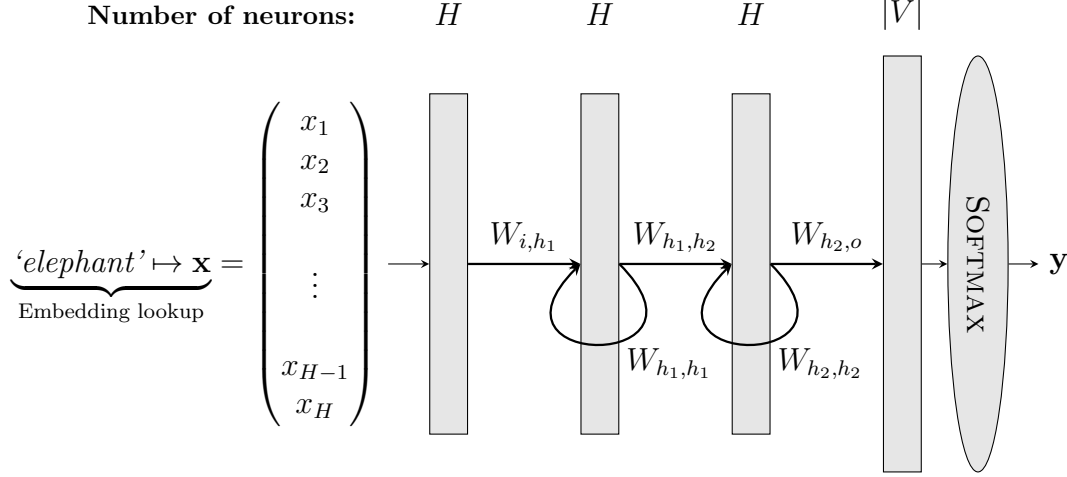


Figure 3.6: The structure of my RNN-based language models.

architecture simply takes the activations from the previous layer, $\mathbf{a}^{l-1,t}$, and the activations from the same layer at the previous time step, $\mathbf{a}^{l,t-1}$, and applies an activation function to their weighted sum. This is shown in equation 2.29, which is repeated below. The fact that the vanilla RNN keeps hold of the activations of the hidden layers from the previous time step means that it implicitly maintains some state:

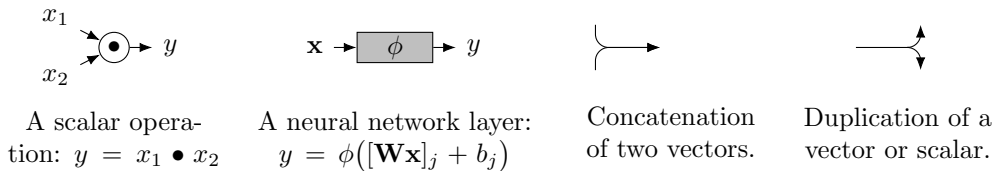
$$\mathbf{a}^{l,t} = \phi(\mathbf{W}^{l-1} \mathbf{a}^{l-1,t} + \mathbf{V}^{l-1} \mathbf{a}^{l,t-1} + \mathbf{b}^{l-1})$$

Unfortunately, it has been observed that it is difficult to train vanilla RNNs to capture long-term dependencies, because the influence of the inputs as they cycle through the recurrent connections in the network often either decay or blow up exponentially [20] [21]. These phenomena result in weights earlier on in the network either vanishing or exploding disproportionately, and are known as the *vanishing gradient problem* and the *exploding gradient problem* respectively. Various improvements have been made to improve upon these problems, two of which I implemented in my project and are outlined below.

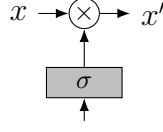
Long Short-Term Memory

Hochreiter and Schmidhuber [17] originally proposed an architecture that directly tackles the vanishing gradient problem, called Long Short-Term Memory (LSTM). The state of each hidden layer l in their architecture at time step t consists of a tuple $(\mathbf{a}^{l,t}, \mathbf{c}^{l,t})$. $\mathbf{a}^{l,t}$, just like in the vanilla RNN cells, allows the LSTM cells to make decisions over a short period of time. The second state vector, $\mathbf{c}^{l,t}$, serves to retain longer-term information. At time step $t + 1$, each element $c_j^{l,t}$ of $\mathbf{c}^{l,t}$ passes through a single LSTM cell in layer l , and can be modified by a series of three *gates*: the *forget gate*, the *input gate* and the *output gate*. These gates can remove, add and output information respectively from $c_j^{l,t}$.

In what follows, I use the notation from figure 2.6, which is repeated below for convenience:



Each LSTM gate has the following structure:



It is named a gate because it controls how much information from x is let through the scalar operation. The sigmoid layer outputs a scalar value in the range $[0, 1]$, where 0 prevents any information from x passing through and 1 allows all of it to pass through.

The LSTM forget, input and output gates in cell j of layer l at time step t are denoted $f_j^{l,t}$, $i_j^{l,t}$ and $o_j^{l,t}$ respectively. If $[\mathbf{x}_1, \mathbf{x}_2]$ denotes the concatenation of vectors \mathbf{x}_1 and \mathbf{x}_2 , and $[\mathbf{x}]_j$ denotes the j^{th} element of vector \mathbf{x} , then the gate values are computed as follows:

$$f_j^{l,t} = \sigma([\mathbf{W}_f^l[\mathbf{a}^{l,t-1}, \mathbf{a}^{l-1,t}]]_j + [\mathbf{b}_f]_j) \quad (3.3)$$

$$i_j^{l,t} = \sigma([\mathbf{W}_i^l[\mathbf{a}^{l,t-1}, \mathbf{a}^{l-1,t}]]_j + [\mathbf{b}_i]_j) \quad (3.4)$$

$$o_j^{l,t} = \sigma([\mathbf{W}_o^l[\mathbf{a}^{l,t-1}, \mathbf{a}^{l-1,t}]]_j + [\mathbf{b}_o]_j) \quad (3.5)$$

Firstly, the forget gate determines how much of the long-term state value, $c_j^{l,t-1}$, should be ‘forgotten’. Next, there is the opportunity for information derived from the short-term state, $\mathbf{a}^{l,t-1}$, and the previous layer activation $\mathbf{a}^{l-1,t}$ to be added to the long-term state. The input gate decides how much of this information should be added. The application of the forget and input gates gives the new value of the long-term state value, $c_j^{l,t}$:

$$c_j^{l,t} = f_j^{l,t} \cdot c_j^{l,t-1} + i_j^{l,t} \cdot \tanh([\mathbf{W}_c^l[\mathbf{a}^{l,t-1}, \mathbf{a}^{l-1,t}]]_j + [\mathbf{b}_c]_j) \quad (3.6)$$

Finally, the output gate is used to control how much of the long-term state is used for the cell output, $a_j^{l,t}$:

$$a_j^{l,t} = o_j^{l,t} \cdot \tanh(c_j^{l,t}) \quad (3.7)$$

This structure is summarised in figure 3.7, and is the version of LSTM that I implemented in my project. There are many different variants of LSTM, but the reason I chose this one in particular is because it seemed to be the best representation of an LSTM baseline. That is, most of the other versions are some form of extension of the one presented above. A notable variant from Gers and Schmidhuber [22] makes use of *peephole connections*, whereby $\mathbf{c}^{l,t-1}$ is added to the input of the forget and input gates, and $\mathbf{c}^{l,t}$ is added to the input of the output gate.

Gated Recurrent Unit

Whilst the LSTM architecture provides a significant improvement for retaining long-term information, it is computationally more expensive to run. Cho et al. propose the Gated Recurrent Unit [18], which is a simplified version of LSTM that is cheaper to compute. The most significant changes are that it uses two gates rather than three, and it merges the long and short-term states into one.

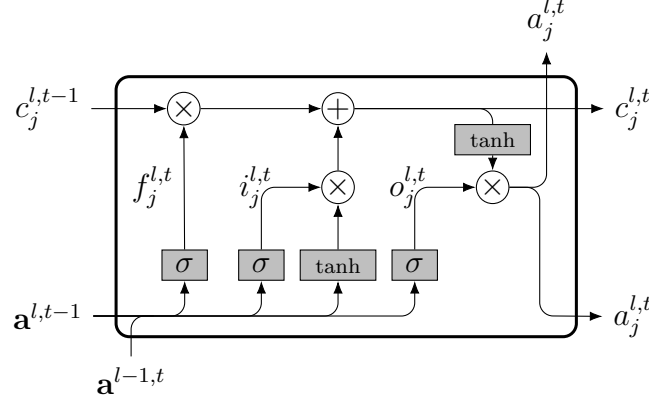


Figure 3.7: A summary of the structure of an LSTM cell. Specifically, this diagram represents the computation of the j^{th} cell in layer l at time step t .

The GRU consists of a *reset gate*, $r_j^{l,t}$, and an *update gate*, $u_j^{l,t}$:

$$r_j^{l,t} = \sigma([\mathbf{W}_r^l[\mathbf{a}_j^{l,t-1}, \mathbf{a}^{l-1,t}]]_j) \quad (3.8)$$

$$u_j^{l,t} = \sigma([\mathbf{W}_u^l[\mathbf{a}_j^{l,t-1}, \mathbf{a}^{l-1,t}]]_j) \quad (3.9)$$

The cell output, $a_j^{l,t}$, is then computed as follows:

$$a_j^{l,t} = u_j^{l,t} \cdot a_j^{l,t-1} + (1 - u_j^{l,t}) \cdot \tilde{a}_j^{l,t} \quad (3.10)$$

where

$$\tilde{a}_j^{l,t} = \tanh([\mathbf{W}_a^l[\mathbf{a}_j^{l,t-1} \odot \mathbf{r}^{l,t}, \mathbf{a}^{l-1,t}]]_j) \quad (3.11)$$

and \odot represents element-wise multiplication.

The reset gate acts similarly to the LSTM forget gate, in that it allows the cell to drop information that it deems irrelevant. The update gate controls how much information from the previous state carries over to the new one. With the short-term and long-term states combined into one, it follows that cells either learn to capture short or long-term information. Those that capture short-term information tend to have more active reset gates, whereas those that capture long-term information tend to have more active update gates.

Abstraction Over the RNN Cells

The vanilla RNN, GRU and LSTM architectures only differ in how their cells are connected together internally. I exploited this in my implementation by abstracting the implementation of the RNN cell from the overall network structure that was presented in section 3.3.2, which made it very easy to add and test new RNN cell structures. A minor exception to this observation is that the LSTM cell maintains a tuple of vectors for its state, whereas the vanilla RNN and GRU cells only have a single state vector. To cater for this difference, I simply register which RNN cell is being used within the generic RNN instance and use that information to determine whether to expect a one or two state vectors at each time step.

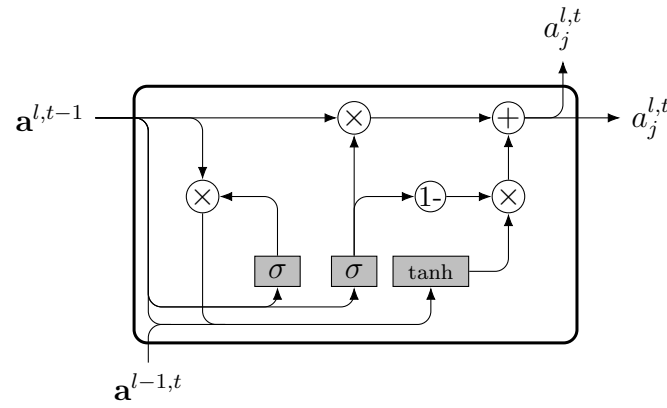


Figure 3.8: A summary of the structure of a GRU cell. Specifically, this diagram represents the computation of the j^{th} cell in layer l at time step t .

3.3.4 Word Embeddings

As mentioned in section 3.3.2, each input word is first mapped to a vector representation known as a *word embedding* which is learned as the network is trained. The goal behind representing words as points in a high-dimensional space is that similar words will cluster close together. This way, the RNN-based language model can more easily generalise to combinations of words that it has not seen before, by understanding which words can be substituted for one another.

I wanted to be able to visualise the word embeddings after each RNN-based language model had been trained. To do this, I used principle component analysis (PCA) to reduce the word embeddings to two-dimensional points that can be plotted on a scatter diagram. PCA is a dimensionality reduction technique that transforms the input data to a new coordinate system such that the first coordinate represents some projection that achieves the greatest variance, the second coordinate the second greatest variance, and so on.

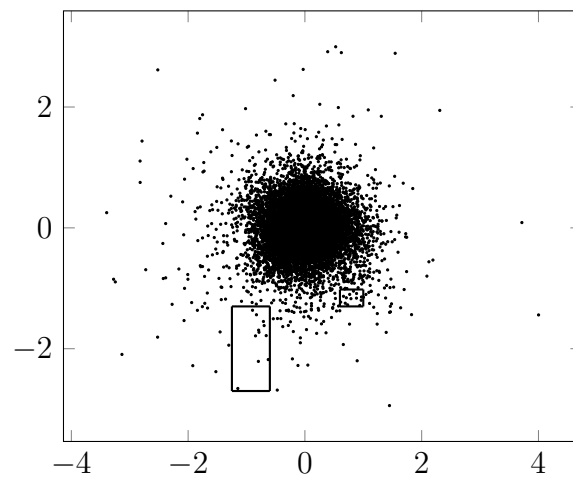


Figure 3.9: LSTM-based language model word embeddings projected onto two dimensions using PCA. The highlighted regions and are shown in figure 3.10.

The words that have been assigned a more distinguishable representation are given by

the points which are further from the centre in figure 3.9. Two regions of the plot have been highlighted, and are shown in figure 3.10.

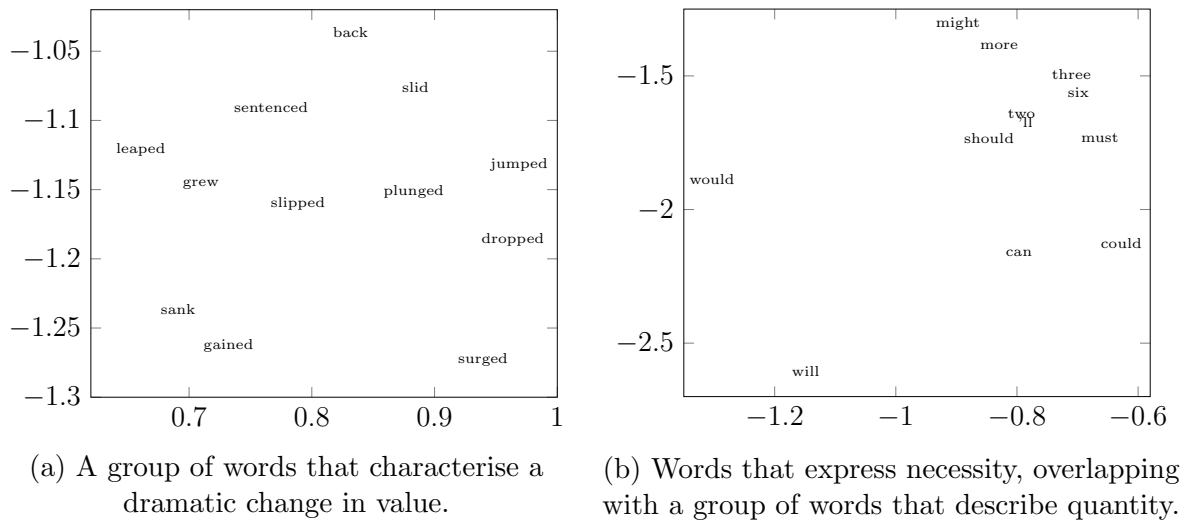


Figure 3.10: Two zoomed in regions of figure 3.9.

3.3.5 Putting Theory into Practice

In this section, I address a number of problems that are faced when training RNNs in practice and describe how I went about mitigating them.

Underfitting and Overfitting

In any supervised machine learning problem, there are one or more underlying patterns in the training data that need to be learned. *Underfitting* occurs when the model fails to capture these patterns, and is typically the case when the model is too simple for the training data, or it is not trained for long enough. This problem was easy to mitigate by adding more parameters to the model, such as the number of hidden neurons, and training it for a sufficient amount of time.

As well as the underlying patterns in the training data, there also often exists some noise which you do not want the model to learn. *Overfitting* occurs when the model learns the noise in the training data, which makes it difficult for it to generalise to new data that it has not seen before. When machine learning models are trained, a portion of the data called the test set is reserved to test the performance of the model. The model is not allowed to see any of the data from the test set during training, so that its ability to generalise to unseen data is fairly evaluated. In order to avoid overfitting my RNN-based language models, I also set aside a portion of data called the validation set. In my implementation, I evaluate the perplexity¹ of my language model on the validation set at the end of each training iteration.

In my implementation, the learning rate starts at some initial value, η_0 . Once the improvement of the perplexity on the validation set between two iterations drops below some

¹Perplexity is the standard performance measure for language models, explained in depth in section 4.1.1.

threshold, Δ , the learning rate is then halved on every new iteration. Finally, to prevent overfitting, when the improvement of the perplexity on the validation set drops below 0.1%, the training process is halted. I tried various values for η_0 and Δ , and chose the best ones for each RNN architecture: $\eta_0 \in \{1, 0.5, 0.1, 0.01, 0.001\}$, $\Delta \in \{10\%, 5\%, 1\%, 0.5\%\}$.

RNNs with a large number of parameters can be particularly susceptible to overfitting. Srivastava proposes a technique called *dropout* for reducing this problem [23], whereby some of the neurons and their connections are randomly dropped from the network during training. Zaremba et al. found that dropout works best in RNNs if it is only applied to the non-recurrent connections [24], which is the method I used in my RNN implementation.

Convergence

Whilst the LSTM and GRU architectures directly combat the aforementioned vanishing gradient problem, there is still the possibility that the gradients can explode, which can make the network unstable during training. In my implementation, I used *gradient clipping* to mitigate this issue, whereby the gradients are clipped once their norm exceeds a certain threshold.

The convergence of a neural network is also dependent on the size of the learning rate. If the learning rate is too large, then the weight updates may overshoot the minimum which they are directed towards. On the other hand, learning rates that are too small result in very slow convergence or the network becoming stuck in local minima. As described above, I tackled this tradeoff by trying various initial values for the learning rate and employing decay after performance improvement drops below some threshold.

3.4 Mobile Keyboard

As a project extension, I decided to apply one of my RNN-based language models in a practical context by building a mobile keyboard. Below I describe how I went about doing this, as well as the lessons I learnt in the process.

iOS applications are written in either Objective-C or Swift. An advantage of using Objective-C is that you can make use of Objective-C++ to link C++ code into your application, which is exactly what I needed to do. My C++ code, however, depended on TensorFlow, so in order to make my language models accessible within the mobile keyboard, I compiled both my RNN-language model implementation and a subset of TensorFlow's C++ API into a single static library by modifying one of the makefiles already present in the TensorFlow codebase.

To build a mobile keyboard in iOS, you have to implement what is called an 'App Extension'. App Extensions differ from normal applications in iOS in that their functionality is system-wide. Apple insist that such ubiquitous components should be extremely responsive and light-weight, which means that they have much stricter memory and CPU limitations than normal applications. Although they do not specify exactly what these limitations are, many people in the developer community estimate that the memory limit is approximately 30 MB.

In order to get an RNN-based language model to run under such memory and CPU limitations, I had to trade-off some of its prediction accuracy. In particular, I reduced the number of hidden neurons and the size of the vocabulary to reduce the space occupied by the model in memory. I also removed the softmax layer to save reduce the CPU overhead, which technically is not needed in text prediction where the sum of the output neurons does not need to sum to 1.

3.4.1 Updating Language Model Predictions On the Fly

A mobile keyboard should present the user with next-word predictions that are updated as they type. Querying the language model is a relatively expensive operation, so in my implementation I try to do this as little as possible. For example, when the user presses the spacebar, I query the language model to retrieve the most probable words that could follow the space. These words are returned in the form of a character-level trie, where their associated probabilities are stored in the nodes of the trie. This way, when the user types subsequent characters, the new top next-word predictions can be retrieved by traversing downwards from the root of the trie and taking the words with the highest probability on the current branch.

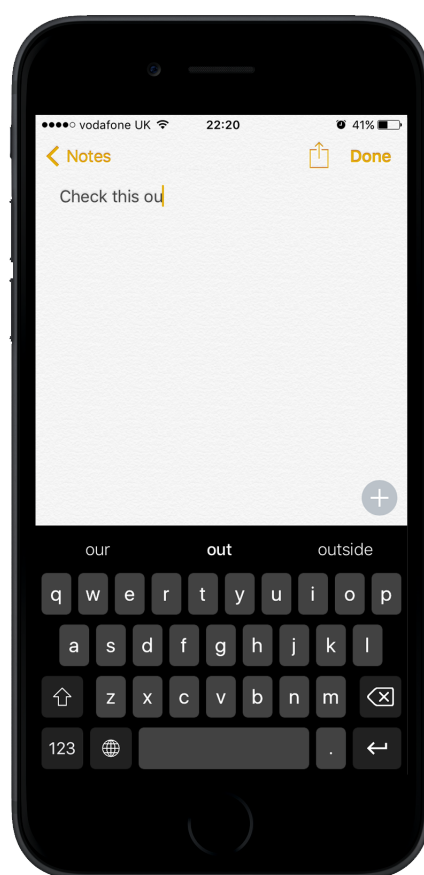


Figure 3.11: The mobile keyboard on iOS. The top three predictions are displayed at the top of the keyboard, and are updated with every character.

3.5 Extending Models to Tackle Error-Prone Text

An issue with existing language models is that their predictions suffer when presented with error-prone text. For example, consider the two word sequences ‘*the weather is*’ and ‘*the whether is*’. Whilst in the first case a language model might successfully predict words such as ‘*sunny*’ or ‘*nice*’, in the second case it will likely struggle to suggest reasonable word candidates. Evidence to support this claim is given in figure 4.8 in section 4.2.3.

In the context of text prediction, it is unfortunately very rare that a human user will not make any mistakes. It is therefore an important issue to address and in this section I propose an extension to a language model which aims to improve its performance on text containing errors.

3.5.1 The Approach

My approach is to attack the source of the problem directly by adding a layer of error correction between the input text and the language model. There are lots of different types of error that can occur in the English language, but my implementation only focuses on correcting spelling mistakes.

Error-correction can be broken up into two parts: detecting the errors and making suggestions to fix them. In the case of spelling mistakes, detection can be done with a single dictionary lookup. Candidate corrections for a spelling mistake should be similar to the original word and probable in the surrounding context. In my implementation, I use the Levenshtein distance as a measure of how similar two words are, and I use the language model itself to gauge how probable candidate corrections are in the surrounding context. Specifically, when searching for candidate replacements for a word w that is preceded by the words w_i^j : the language model is used to produce a probability distribution over the words in its vocabulary for how likely they are to follow w_i^j , these words are then traversed in decreasing order of probability and the first one, if any, that has an edit distance within some threshold δ of the original spelling is selected as a replacement, as described in algorithm 2. I implemented this approach as an extension to my RNN-based language models.

Algorithm 2 Computing $\mathbb{P}(w_j | w_i^{j-1})_{\text{ERROR_CORRECTING}}$

```

1: procedure PROB( $w_j, w_i^{j-1}$ )
2:    $outputs \leftarrow \emptyset$ 
3:    $state = 0$ 
4:   for  $w$  in  $w_i^{j-1}$  do
5:      $w_{correct} \leftarrow w$ 
6:     if  $w \notin \text{dictionary}$  and  $outputs \neq \emptyset$  then
7:       SORTBYDECREASINGPROBABILITY( $outputs$ )
8:       for  $w'$  in  $outputs$  do
9:          $d \leftarrow \text{LEVENSHTEINDISTANCE}(w, w')$ 
10:        if  $d \leq \delta$  then
11:           $w_{correct} \leftarrow w'$ 
12:          break
13:     $outputs, state \leftarrow \text{FEEDFORWARD}(w_{correct}, state)$ 
14:  return  $outputs[w_i]$ 

```

Chapter 4

Evaluation

In this chapter, I will first describe the benchmarking framework that I built to evaluate the language models, before proceeding onto the results. The results section is threefold: firstly I present the performance of the existing language models that I implemented, secondly I focus on the tradeoffs faced when employing those models on a mobile device, and finally I display my findings in language modelling on error-prone text.

4.1 Evaluation Methodology

4.1.1 Metrics

In the context of text prediction, there are essentially two questions one might want to answer when evaluating a language model:

1. How accurately does the language model predict text?
2. How much resource, such as CPU or memory, does the language model consume?

In order to answer these questions, I implemented a generic benchmarking framework that can return a series of metrics that fall into one of the two aforementioned categories when given a language model. These metrics include perplexity, average-keys-saved, memory usage and average inference time. The first two are concerned with the accuracy of language models and the latter two relate to the resource usage. I also recorded how long it took to train each model.

Perplexity

Perplexity is the most widely-used metric for language models, and is therefore an essential one to include so that my results can be compared with those of other authors. Given a sequence of words $w_1^N = w_1 w_2 \dots w_N$ as test data, the perplexity PP of a language model L is defined as:

$$\text{PP}_L(w_1^N) = \sqrt[N]{\frac{1}{\mathbb{P}_L(w_1^N)}} = \sqrt[N]{\prod_{i=1}^N \frac{1}{\mathbb{P}_L(w_i | w_1^{i-1})}} \quad (4.1)$$

where $\mathbb{P}_L(w_i | w_1^{i-1})$ is the probability computed by the language model L of the word w_i following the words w_1^{i-1} . The key point is that **lower values of perplexity indicate**

better prediction accuracy for language models trained on a particular training set.

This somewhat arbitrary-looking formulation can be better understood from a touch of information theory. In information theory, the cross-entropy $H(p, q)$ between a true probability distribution p and an estimate of that distribution q is defined as:¹

$$H(p, q) = - \sum_x p(x) \log_2 q(x)$$

It can be shown that $H(p, q) = H(p) + D_{KL}(p||q)$ where $D_{KL}(p||q) \geq 0$ is the Kullback-Leibler distance between p and q . Generally speaking, the better an estimate q is of p , the lower $H(p, q)$ will be, with a lower bound of $H(p)$, the entropy of p .

The perplexity PP of a model q , with respect to the true distribution p it is attempting to estimate, is defined as:

$$\text{PP} = 2^{H(p, q)} \quad (4.2)$$

Language models assign probability distributions over sequences of words, and so it seems reasonable to use perplexity as a motivation for a measure of their performance. In the context of language modelling, however, we do not know what the underlying distribution of p is, so it is approximated with Monte Carlo estimation by taking samples of p (i.e. sequences of words from the test data) as follows:

$$\text{PP}_L(w_1^N) = 2^{-\frac{1}{N} \sum_i \log_2 \mathbb{P}_L(w_i | w_1^{i-1})}$$

With a little algebra, this can be rearranged to give equation 4.1.

One issue with perplexity is that it is undefined if $\mathbb{P}_L(w_i | w_1^{i-1})$ is 0 at any point. To get around this in my implementation, I replaced probability values of 0 with the small constant `1e-9`. Results that use this approximation are marked.

Average-Keys-Saved

It is typical for the top three next-word predictions to be displayed and updated as the user types in a mobile keyboard, as described in section 3.4. Clearly, it is in the interest of the mobile keyboard developer to minimise the amount of typing a user has to do before the correct prediction is displayed. Average-keys-saved is based on this incentive, and is defined as the number of keys that the user would be saved from typing as a result of the correct next word appearing in the top three predictions, averaged over the number of characters in the test data.

As an example, if the user is typing `science` and the word `science` appears in the top three predictions after they have typed `sc`, then that would count as 5 characters being saved, averaging at $\frac{5}{7}$ keys saved per character. Averaging over the number of characters in the test data ensures that the results are not biased by the data containing particularly long words, which are easier to save characters on.

¹Note that $H(p, q)$ is often also used to denote the joint entropy of p and q , which is a different concept.

Memory Usage

This is measured as the amount of physical memory in megabytes occupied by the process in which the language model under test is instantiated.

Training Time

The amount of time it took to train the language model.

Average Inference Time

This is measured as the amount of time in milliseconds that the language model takes to assign a probability to all of the words in its vocabulary given a sequence of words, averaged over a large number of sequences.

4.1.2 Datasets

I used three datasets throughout the evaluation of my project:

Penn Tree Bank (PTB) Dataset

The Penn Tree Bank is a popular dataset for measuring the quality of language models, created from text from the Wall Street Journal. It has already been preprocessed such that numbers are replaced with `N`, rare words are replaced with `<unk>` and the text has been split up into one sentence per line. It has been split up into a training, validation and test set. The training set has 10,000 unique words and 887,521 words overall.

Given that the Penn Tree Bank is so widely adopted, I used it for all tests in which the size of the training data is fixed, and will refer to it as PTB.

One Billion Word (1BW) Benchmark

This is a much larger dataset produced by Google of approximately 1 billion words [1]. I used this dataset for tests in which the size of the training data is a variable under investigation, and will refer to it as 1BW.

Cambridge Learner Corpus (CLC)

The Cambridge Learner Corpus is a dataset of 1,244 exam scripts written by candidates sitting the Cambridge ESOL First Certificate in English (FCE) examination in 2000 and 2001 [25]. The original dataset contains the scripts annotated with corrections to all of the mistakes by the candidates. In this project I make use of a preprocessed version of the dataset, in which there is one file containing the error-free version of the exam scripts and there is another file containing the original exam scripts with their errors. These two files are aligned line by line. I used this dataset when exploring the performance of language models on error-prone text, and will refer to it as CLC.

4.2 Results

4.2.1 Existing Models

Smoothing techniques and the value of n in n -gram models

The first set of language models that I built were n -gram models, along with a series of smoothing techniques for improving their predictions on less frequent n -grams.

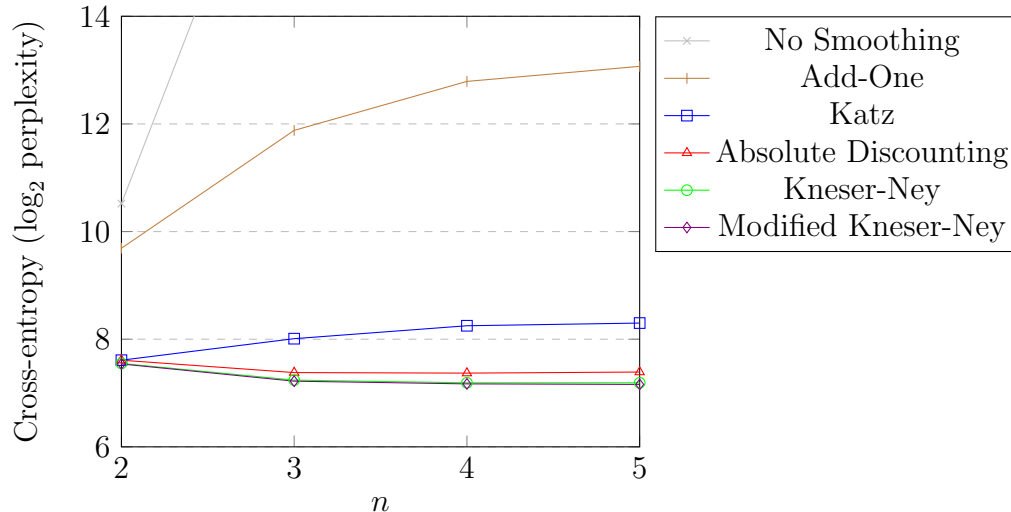


Figure 4.1: Cross-entropy of n -gram models trained on the PTB dataset.

Recall from equation 4.2 that cross-entropy is just the binary logarithm of perplexity, and that lower perplexity scores indicate better prediction accuracy. With this in mind, it is clear that modified Kneser-Ney smoothing offers the best prediction performance amongst the n -gram models.

The change in performance with the value of n is interesting. Intuitively, one might expect that increasing n will always yield better results, because this corresponds to increasing the number of words you use to make a prediction. However, for n -gram models with no smoothing, add-one smoothing or Katz smoothing, this is not the case. For n -gram models with add-one or no smoothing, this is because they do not employ backoff. At higher values of n , n -grams are much more sparse, so without any backoff higher n -gram models can only rely on sparse counts, resulting in lower probabilities being assigned to plausible sequences of words. Katz smoothing does employ backoff, and achieves much better performance, but it still distributes too much probability to rare n -grams.

A comparison of RNN-based models with n -gram models

As described in chapter 3, I also implemented three RNN-based language models which differ in the RNN cell architecture: vanilla RNN, Gated Recurrent Unit and Long Short-Term Memory.

At a qualitative level it could be seen that the RNN-based language models, with the LSTM architecture in particular, were better at making predictions with long-term word dependencies than the n -gram smoothed models. This is demonstrated in figure 4.2, and

confirmed with the quantitative readings presented in figure 4.3.

money is the root of a new york city.	money is the root of the nation's largest public bank.
the meaning of life is a unit of the company's shares outstanding.	the meaning of life is the best way to get the way to the public.
science and was up N N from \$ N million or N cents a share.	science's most recent issue of the nation 's largest economic indicators.
(a) 5gram + modified Kneser-Ney smoothing	(b) 2-layer LSTM with 256 hidden neurons

Figure 4.2: Sentences generated from two different language models trained on the PTB dataset. The words highlighted in bold indicate the words given to the model, and the non-bold words indicate the words that were generated by repeatedly taking the next most probable word until a full stop is predicted.

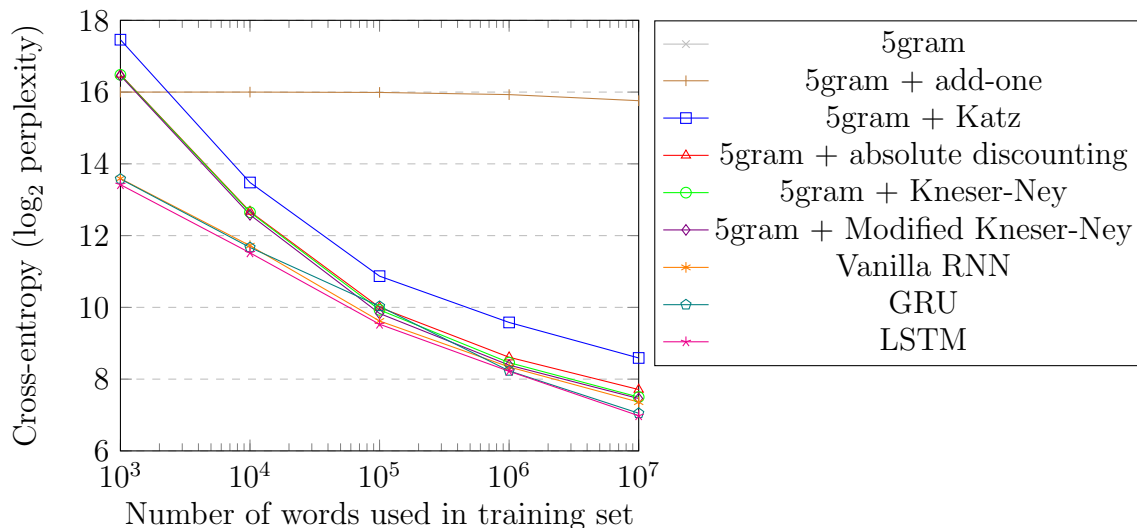


Figure 4.3: Cross-entropy of various language models with respect to the training set size, using the 1BW dataset.

As can be seen in figure 4.3, the RNN-based language models consistently outperform the n -gram based models in terms of prediction accuracy. This difference is most significant for small training sets of less than 10^5 words, and starts to become more pronounced again once the training set grows large beyond 10^6 words. The performance gain from adding more training data seems to start flattening out towards 10^7 words for the smoothed n -gram models, whereas the RNN-based ones seem to continue improving.

Some datasets are more predictable than others, which means that metrics like perplexity and average-keys-saved depend on which dataset is used. In order to produce metrics that are comparable with the work of other authors, I have tested my language models on the Penn Tree Bank dataset, as shown in figure 4.4. The perplexity was calculated over the whole of the test set, whereas average-keys-saved, a more expensive metric to compute, was taken over the first 1000 words of the test set. Indeed, my results for 3-gram modified Kneser-Ney smoothing, 5-gram modified Kneser-Ney smoothing and the vanilla RNN are

comparable to those from Mikolov [26], who also benchmarked these three models.

Language Model	Perplexity	Average-Keys-Saved	Memory Usage (MB)	Training Time (min,secs) [†]	Average Inference Time (ms)
3-gram	$4.54 \times 10^{5*}$	0.35014	266.91	11s	62
3-gram + add-one	3764.96	0.53063	266.94	11s	41
3-gram + Katz	256.95	0.68482	266.71	14s	88
3-gram + absolute disc.	166.03	0.72178	266.78	13s	63
3-gram + KN	150.73	0.72466	266.88	14s	54
3-gram + modified KN	149.54	0.72355	266.97	14s	54
5-gram	$1.96 \times 10^{8*}$	0.07167	737.36	26s	130
5-gram + add-one	8610.45	0.33886	737.30	26s	63
5-gram + Katz	314.49	0.67154	737.43	41s	156
5-gram + absolute disc.	167.38	0.72333	737.43	40s	126
5-gram + KN	146.35	0.72598	737.37	44s	114
5-gram + modified KN	142.68	0.72554	737.53	50s	116
Vanilla RNN	131.03	0.72776	253.67	15m 10s	39
Gated Recurrent Units	114.52	0.73993	271.39	28m 35s	37
Long Short-Term Memory	112.47	0.73617	287.13	18m 59s	38
LSTM, 5-gram + MKN (av)	96.07	0.75719	929.39	19m 49s	189
LSTM, 5-gram + MKN (int)	94.70	0.75830	927.20	19m 49s	190

**These perplexity scores use the approximation mentioned in section 4.1.1 where 0-valued probabilities are replaced with a small constant to avoid division by 0. [†]The n-gram models were trained on my laptop, whereas the neural models were trained on a GPU cluster on the High Performance Computing Service. Nevertheless, the n-gram models were still much faster to train.*

Figure 4.4: A benchmark of various language models on the PTB dataset.

There is a lot of information that can be drawn from figure 4.4, so I will highlight only the most important and interesting points:

- In general, the RNN-based language models outperform the n -gram based models in terms of prediction accuracy. In my implementation, they run slightly faster at inference time and also consume less memory, although this is partly due to the fact that they are entirely unpruned. The RNN-based language models take significantly longer to train and obtain the correct hyperparameters for.
- Ignoring insignificant differences in memory and time, the discounting-based smoothing techniques (absolute discounting, Kneser-Ney smoothing and modified Kneser-Ney smoothing) equal or outperform add-one and Katz smoothing in every metric.
- LSTM yields marginally better prediction accuracy than the GRU architecture, but requires more memory. Apart from the memory overhead, GRUs significantly outperform the vanilla RNN architecture. The differences in training time for the RNN-based models should be taken lightly, because each model was trained for a

different number of iterations, depending on how long it took to converge to a steady perplexity score on the validation set.

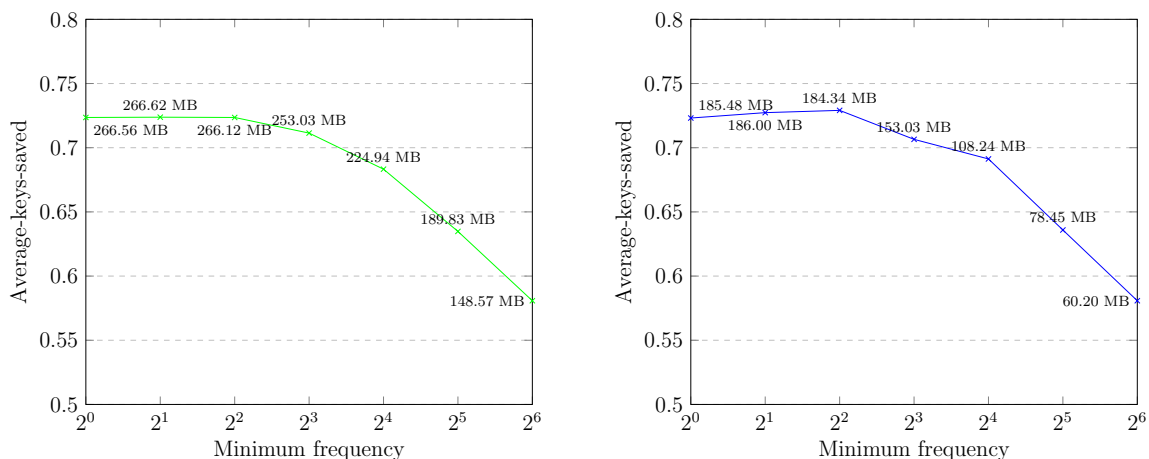
- Interestingly, a substantial improvement in perplexity and average-keys-saved can be obtained by simply averaging the probabilities produced by an n -gram and an RNN-based language model. This seems to imply that the two classes of language model complement one another, in the sense that RNN-based models make strong predictions when some of the n -gram predictions are weak, and vice versa.
- The combined language models can be improved even more by interpolating their probabilities rather than just averaging them. The result shown for the interpolation of 5-gram modified Kneser-Ney with LSTM used $\lambda = 0.38$, where the probability was calculated as $\lambda(n\text{-gram probability}) + (1 - \lambda)(\text{LSTM probability})$.

4.2.2 On a Mobile Device

Incorporating an RNN-based language model implementation into a mobile keyboard, as described in section 3.4, presented a different series of challenges. The strict memory and CPU limitations of App Extensions in iOS forced me to explore the tradeoffs between resource consumption and prediction performance.

As shown in figure 4.4, the vanilla RNN architecture presents the smallest memory overhead, and so this is what I used in the following experiments, along with the 3-gram modified Kneser-Ney smoothing model.

Two obvious ways of decreasing the memory footprint of a language model are decreasing the vocabulary size and, for RNN-based models, decreasing the number of hidden neurons. What is less obvious, is how quickly the average-keys-saved drops with these two parameters.



(a) 3gram + modified Kneser-Ney smoothing (b) Vanilla RNN with 256 hidden neurons

Figure 4.5: The effect of *minimum frequency* and therefore vocabulary size on average-keys-saved and memory usage. Both models were trained using the PTB dataset.

The vocabulary size was altered by changing a parameter called the *minimum frequency* in my model implementations. This parameter defines the minimum number of times a word

must appear in the training set to be considered part of the language model vocabulary. Any words that were not included in the vocabulary were mapped to the unknown-word marker, `<unk>`.

For both the n -gram and vanilla RNN models, the minimum frequency can be increased to approximately 2^4 before the average-keys-saved starts dropping rapidly. In the case of the RNN model, a saving of 77.24 MB is made by increasing the minimum frequency to 2^4 , with a drop of only 0.03187 in average-keys-saved. The savings made on the n -gram model are less substantial.

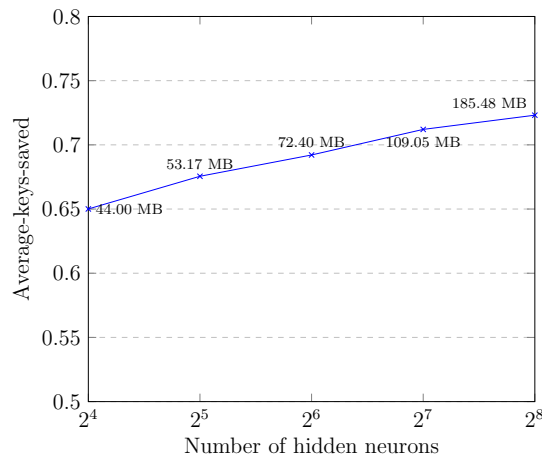


Figure 4.6: The effect of the number of hidden neurons on average-keys-saved and memory usage in a 2-layer vanilla RNN.

From figure 4.6, it can be seen that much greater reductions in memory usage can be achieved for a given loss in average-keys-saved by decreasing the number of hidden neurons. Dropping from 256 to 32 hidden neurons in each layer gives a memory saving of 132.31 MB, with an average-keys-saved loss of 0.04758. It is intriguing that a 2-layer vanilla RNN with only 32 hidden neurons in each layer can achieve a better average-keys-saved score than the 5-gram model with Katz smoothing from figure 4.4.

Of course, there are several other techniques that can be used to optimise the memory and CPU overhead of language models for mobile platforms. These include, but are not limited to: using 16-bit floats, rather than 32-bit or 64-bit, to store the model parameters; removing the softmax layer from the RNN, since the outputs do not need to sum to 1 in text prediction; and memory-mapping the model to reduce memory pressure when it is initially being loaded.

4.2.3 On Error-Prone Text

In this section I present my findings for the performance of language models on error-prone text. In the ideal case, when presented with error-prone text, a RNN-based language model should be able to predict the next words that would follow the error-free version of the text, as shown in figure 4.7.

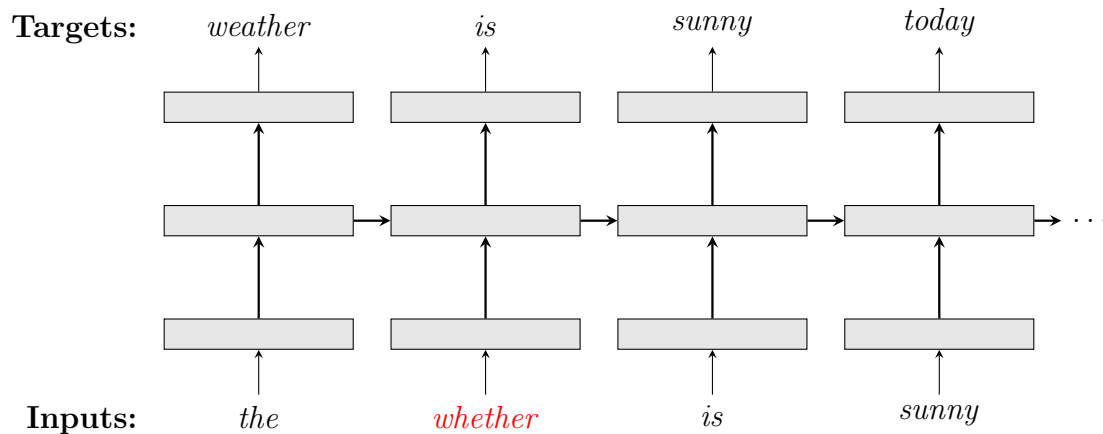


Figure 4.7: The ideal language model that can make predictions regardless of errors in the input text.

Normally, language models are evaluated by feeding them a sequence of input words (the *inputs*) and seeing how well they predict that same sequence of input words shifted forward in time by one word (the *targets*). When evaluating my error-correcting language models, I did exactly the same thing, except that the inputs were allowed to contain error in them and the targets were not.

In order to obtain sequences of words with and without errors in them I used the CLC dataset. I split this dataset up into six files:

```
train.correct.txt  train.incorrect.txt
valid.correct.txt  valid.incorrect.txt
test.correct.txt   test.incorrect.txt
```

The training, validation and test sets all consisted of one file containing uncorrected text and another file containing the corresponding corrected text. I trained an 2-layer LSTM-based language model with 256 hidden neurons on `train.correct.txt`, and used `valid.correct.txt` to guide the learning rate decay. In order to focus the evaluation on the error-correcting ability of my language models, I removed any pairs of lines from the test set files that were identical. I also removed any pairs of lines that contained corrections that involve the insertion or deletion of words, because my model was not designed to handle these types of errors, as explained in section 3.5.

Before exploring ways to improve predictions on error-prone text, I first established an approximate lower and upper bound on the performance of my model. The ideal error-correcting language model would predict words as if the input were not error-prone, therefore, to obtain the upper bound in performance, I evaluated the LSTM-based language model using `test.correct.txt` for the input words and `test.correct.txt` for the target words. In order to establish a lower bound, I evaluated the performance of the same model using `test.incorrect.txt` for the input words and `test.correct.txt` for the target words. This model did not contain any error-correcting capabilities.

As explained in section 3.5, my error-correcting language model attempted to replace words in the context that were not contained in a dictionary of English words, by going through the language model vocabulary in order of decreasing probability and picking the

first candidate word within an edit distance δ of the original word.

Language model	Description	Perplexity	Average-keys-saved
<i>Upper bound</i>	The best performance achievable assuming perfect error-correction.	77.83	0.65483
<i>Error-correcting LSTM</i>	$\delta = 1$	87.99	0.64453
<i>Error-correcting LSTM</i>	$\delta = 2$	88.47	0.64376
<i>Error-correcting LSTM</i>	$\delta = 3$	88.57	0.64273
<i>Lower bound</i>	The performance that would be achieved with no error-correction.	89.70	0.63912

Figure 4.8: The performance of my error-correcting language model with respect to the upper and lower bound, using the CLC dataset.

It can be seen from figure 4.8 that this approach improves the performance of language model predictions on text containing errors that do not require the insertion or deletion of words. The best result arises from an edit distance threshold of 1, which seems attributed to the fact that a larger threshold allows for more false corrections, as shown in figure 4.10.

Input	i am more than happy to give you the nessessary information.
Target	i am more than happy to give you the necessary information.
Prediction	i am more than happy to give you the necessary information.

Figure 4.9: An example of an accurate correction made by my model, with $\delta = 2$.

Input	the bus will pick you up right at your hotel entery .
Target	the bus will pick you up right at your hotel entrance .
Prediction	the bus will pick you up right at your hotel every .

Figure 4.10: An example of a false correction made by my model, with $\delta = 2$.

There is still a gap to be closed, which is demonstrated in figure 4.11 by the fact that my model does not correct non-spelling mistakes.

Input	appart of that , there is no recommendation as to what to wear.
Target	apart from that, there is no recommendation as to what to wear.
Prediction	apart of that, there is no recommendation as to what to wear.

Figure 4.11: An example of a correction that my model failed to make, with $\delta = 2$.

Chapter 5

Conclusions

The project was a success: I implemented a variety of n -gram and RNN-based language models, and benchmarked them in terms of both their prediction performance and their applicability to running on a mobile platform, as set out in my project proposal. Additionally, I built a mobile keyboard for iOS that uses my language models as a library, and proposed a novel extension to an existing language model that improves its performance on error-prone text.

My implementation of existing language models, which consists of an n -gram model with 5 different smoothing techniques and an RNN-based model of 3 different architectures, saw a level of performance that was comparable with measures made by existing authors. For example, my 5-gram modified Kneser-Ney smoothing model, vanilla RNN model and LSTM model all matched the results from Mikolov on the Penn Tree Bank dataset [26] [27].

As well as confirming that modified Kneser-Ney is the best smoothing technique for n -gram models, and that LSTM-based models offer the strongest predictions, I also demonstrated that such models complement one another, in that the interpolation of their probability distributions surpasses the performance of all stand-alone language models in this dissertation. This suggests that the best option for language modelling on a mobile device could be to interpolate several extremely light-weight n -gram and RNN-based language models.

I explicitly demonstrated that language model performance is deterred by error-prone text. Whilst I made some modifications to an existing model that saw an increase in performance on such text, there is still room for improvement. If I were given more time to work on the project, I would investigate ways to cater for non-spelling corrections. From the CoNLL shared tasks of 2013 [4] and 2014 [5], it is clear that there is a plethora of methods that already exist for grammatical error-correction. It would be interesting to see which of these could be applied to improve the performance of language models on error-prone text.

Throughout this dissertation, I have learnt a great deal about language modelling, but what have my language models learnt from me? Here is a one line summary of my work, in the words of an LSTM-based language model, trained on the text of this L^AT_EXfile:

Blah blah blah.

Bibliography

- [1] C. Chelba, T. Mikolov, M. Schuster, Q. Ge, T. Brants, P. Koehn, and T. Robinson, “One billion word benchmark for measuring progress in statistical language modeling,” *arXiv preprint arXiv:1312.3005*, 2013.
- [2] S. F. Chen and J. Goodman, “An empirical study of smoothing techniques for language modeling,” in *Proceedings of the 34th annual meeting on Association for Computational Linguistics*, pp. 310–318, Association for Computational Linguistics, 1996.
- [3] S. F. Chen and J. Goodman, “An empirical study of smoothing techniques for language modeling,” *Comput. Speech Lang.*, vol. 13, pp. 359–394, Oct. 1999.
- [4] H. T. Ng, S. M. Wu, Y. Wu, C. Hadiwinoto, and J. Tetreault, “The conll-2013 shared task on grammatical error correction,” 2013.
- [5] H. T. Ng, S. M. Wu, T. Briscoe, C. Hadiwinoto, R. H. Susanto, and C. Bryant, “The conll-2014 shared task on grammatical error correction,” in *CoNLL Shared Task*, pp. 1–14, 2014.
- [6] W. E. Johnson, “Probability: The deductive and inductive problems,” *Mind*, vol. 41, no. 164, pp. 409–423, 1932.
- [7] H. Ney, U. Essen, and R. Kneser, “On structuring probabilistic dependences in stochastic language modelling,” *Computer Speech & Language*, vol. 8, no. 1, pp. 1–38, 1994.
- [8] R. Kneser and H. Ney, “Improved backing-off for m-gram language modeling,” in *Acoustics, Speech, and Signal Processing, 1995. ICASSP-95., 1995 International Conference on*, vol. 1, pp. 181–184, IEEE, 1995.
- [9] I. J. Good, “The population frequencies of species and the estimation of population parameters,” *Biometrika*, pp. 237–264, 1953.
- [10] S. Katz, “Estimation of probabilities from sparse data for the language model component of a speech recognizer,” *IEEE transactions on acoustics, speech, and signal processing*, vol. 35, no. 3, pp. 400–401, 1987.
- [11] F. A. Azevedo, L. R. Carvalho, L. T. Grinberg, J. M. Farfel, R. E. Ferretti, R. E. Leite, R. Lent, S. Herculano-Houzel, *et al.*, “Equal numbers of neuronal and nonneuronal cells make the human brain an isometrically scaled-up primate brain,” *Journal of Comparative Neurology*, vol. 513, no. 5, pp. 532–541, 2009.
- [12] W. S. McCulloch and W. Pitts, “A logical calculus of the ideas immanent in nervous activity,” *The bulletin of mathematical biophysics*, vol. 5, no. 4, pp. 115–133, 1943.

- [13] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, “Learning internal representations by error propagation,” tech. rep., DTIC Document, 1985.
- [14] G. E. Hinton, S. Osindero, and Y.-W. Teh, “A fast learning algorithm for deep belief nets,” *Neural computation*, vol. 18, no. 7, pp. 1527–1554, 2006.
- [15] Y. Bengio, Y. LeCun, *et al.*, “Scaling learning algorithms towards ai,” *Large-scale kernel machines*, vol. 34, no. 5, pp. 1–41, 2007.
- [16] K. Hornik, M. Stinchcombe, and H. White, “Multilayer feedforward networks are universal approximators,” *Neural networks*, vol. 2, no. 5, pp. 359–366, 1989.
- [17] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [18] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, “Learning phrase representations using rnn encoder-decoder for statistical machine translation,” *arXiv preprint arXiv:1406.1078*, 2014.
- [19] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, *et al.*, “Tensorflow: A system for large-scale machine learning,” in *Proceedings of the 12th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*. Savannah, Georgia, USA, 2016.
- [20] S. Hochreiter, *Untersuchungen zu dynamischen neuronalen Netzen*. PhD thesis, diploma thesis, institut für informatik, lehrstuhl prof. brauer, technische universität münchen, 1991.
- [21] Y. Bengio, P. Simard, and P. Frasconi, “Learning long-term dependencies with gradient descent is difficult,” *IEEE transactions on neural networks*, vol. 5, no. 2, pp. 157–166, 1994.
- [22] F. A. Gers and J. Schmidhuber, “Recurrent nets that time and count,” in *Neural Networks, 2000. IJCNN 2000, Proceedings of the IEEE-INNS-ENNS International Joint Conference on*, vol. 3, pp. 189–194, IEEE, 2000.
- [23] N. Srivastava, *Improving neural networks with dropout*. PhD thesis, University of Toronto, 2013.
- [24] W. Zaremba, I. Sutskever, and O. Vinyals, “Recurrent neural network regularization,” *arXiv preprint arXiv:1409.2329*, 2014.
- [25] H. Yannakoudakis, T. Briscoe, and B. Medlock, “A new dataset and method for automatically grading esol texts,” in *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pp. 180–189, Association for Computational Linguistics, 2011.
- [26] T. Mikolov, “Statistical language models based on neural networks,” *Presentation at Google, Mountain View, 2nd April*, 2012.
- [27] T. Mikolov, A. Joulin, S. Chopra, M. Mathieu, and M. Ranzato, “Learning longer memory in recurrent neural networks,” *arXiv preprint arXiv:1412.7753*, 2014.

Appendix A

Backpropagation Recurrence Relation

To see how the recurrence relation for δ_i^l is derived, consider the neuron i in the first hidden layer of the neural network shown in figure A.1. Without loss of generality, suppose this layer is layer l .

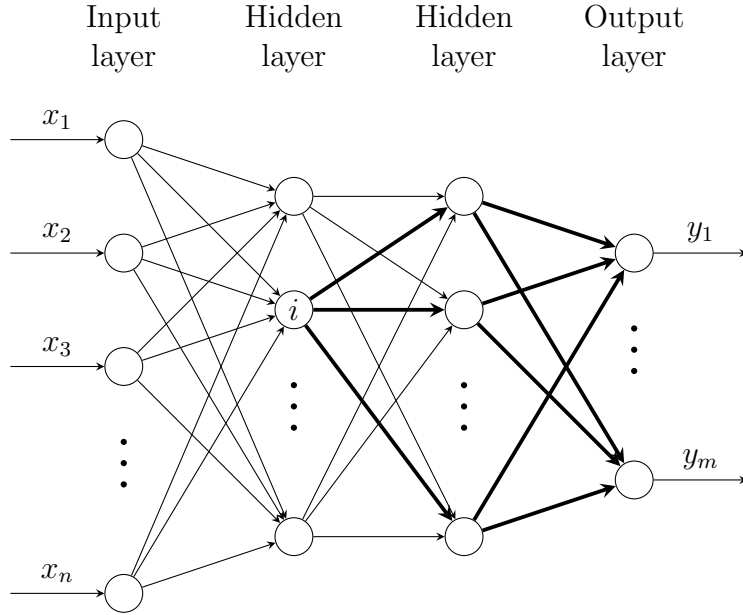


Figure A.1: An illustration of the connections along which the partial derivatives are chained for backpropagating to a particular node i .

From the diagram, it is clear that the derivative $\frac{\partial \mathcal{L}}{\partial z_i^l}$ depends on the derivatives of all of the neurons in layer $l + 1$, and so from the chain rule it follows that:

$$\frac{\partial \mathcal{L}}{\partial z_i^l} = \sum_j \frac{\partial \mathcal{L}}{\partial z_j^{l+1}} \frac{\partial z_j^{l+1}}{\partial z_i^l}$$

Plugging in the definition of δ_i^l , and simplifying gives:

$$\delta_i^l = \sum_j \delta_j^{l+1} \frac{\partial z_j^{l+1}}{\partial z_i^l} = \sum_j \delta_j^{l+1} \frac{\partial z_j^{l+1}}{\partial y_i^l} \frac{\partial y_i^l}{\partial z_i^l} = \phi'(z_i^l) \sum_j \delta_j^{l+1} w_{i,j}^l$$

as required.

These ideas can be extended to a recurrent neural network, whereby the derivative $\frac{\partial \mathcal{L}}{\partial z_i^l}$ depends not only on the derivatives of the neurons in layer $l + 1$, but also on the neurons in layer l at time step $t + 1$:

$$\frac{\partial \mathcal{L}}{\partial z_i^{l,t}} = \sum_j \frac{\partial \mathcal{L}}{\partial z_j^{l+1,t}} \frac{\partial z_j^{l+1,t}}{\partial z_i^{l,t}} + \sum_j \frac{\partial \mathcal{L}}{\partial z_j^{l,t+1}} \frac{\partial z_j^{l,t+1}}{\partial z_i^{l,t}}$$

Again, plugging in the definition of $\delta_i^{l,t}$ and simplifying gives:

$$\delta_i^{l,t} = \sum_j \delta_j^{l+1,t} \frac{\partial z_j^{l+1,t}}{\partial z_i^{l,t}} + \sum_j \delta_j^{l,t+1} \frac{\partial z_j^{l,t+1}}{\partial z_i^{l,t}} = \phi'(z_i^{l,t}) \left(\sum_j \delta_j^{l+1,t} w_{i,j}^l + \sum_j \delta_j^{l,t+1} v_{i,j}^l \right)$$

as required.

Appendix B

Project Proposal