

INTERNATIONAL UNIVERSITY

# MACHINE LEARNING AND ITS APPLICATION IN FINANCE

by

Do Viet Ho Tam Thuc

A thesis submitted in partial fulfillment for the  
degree of Bachelor of Science

in the  
International University  
Department of Mathematics

June 2017

*If we knew what it was we were doing, it would not be called research, would it?*

— Albert Einstein

INTERNATIONAL UNIVERSITY

## *Abstract*

International University  
Department of Mathematics

Bachelor of Science

by Do Viet Ho Tam Thuc

Recent development in the field Artificial Intelligence have brought financial engineers new advance tools for modeling the stock market. Nowadays, supercomputers are integrated into most stock exchanges to capture important signals from market and try to predict the movements of the stock. In this thesis, we propose a deep learning model to capture and understand the events in newspaper, then predict the movement of stock price. The proposed model is simple and does not need to use a supercomputer, which is suitable for a individual trader, it achieves more than 65% accurate performance and highly comparable with the state of the art methods.

# *Acknowledgements*

I would first like to thank my thesis advisor PhD. Nguyen Tien Dung of the applied mathematics department at Polytechnic University - HCM city. The door to Prof. Dung office was always open whenever I ran into a trouble spot or had a question about my thesis or writing. He consistently allowed this thesis to be my own work, but steered me in the right the direction whenever he thought I needed it.

I would also like to acknowledge Dang Gia Huy, Nguyen Ngoc Son An and Mr.Phan Thai Duy Tan as the second reader of this thesis, and I am gratefully indebted to them for their very valuable comments on this thesis.

Finally, I must express my very profound gratitude to my parents for providing me with unfailing support and continuous encouragement throughout my years of study and through the process of researching and writing this thesis. This accomplishment would not have been possible without them. Thank you.

Author

Tam Thuc

# Contents

<b>Abstract</b>	<b>ii</b>
<b>Acknowledgements</b>	<b>iii</b>
<b>Abbreviations</b>	<b>vi</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Previous works . . . . .	2
1.2 Goal of thesis . . . . .	2
<b>2 Deep learning Models</b>	<b>4</b>
2.1 What is Deep Learning ? . . . . .	4
2.1.1 Deep Feedforward Networks . . . . .	6
2.1.2 Convolutional Neural Networks . . . . .	7
2.1.3 Recurrent Neural Networks . . . . .	10
2.2 Optimization for training deep models . . . . .	12
2.2.1 Learning process and traditional optimization . . . . .	12
2.2.2 Basic learning algorithm . . . . .	13
<b>3 Methods</b>	<b>16</b>
3.1 Data Preprocessing . . . . .	16
3.2 Information Extraction . . . . .	16
3.3 Embedding Vector of Information . . . . .	17
3.3.1 Model Architecturre . . . . .	17
3.3.2 Training Process for Events Embedding . . . . .	18
3.4 Stock Price Movement Prediction Model . . . . .	19
<b>4 Experimental Setup</b>	<b>20</b>
4.1 Data Collection . . . . .	20
4.2 Baseline stock prediction models and proposed models . . . . .	21
4.3 Apple Stock Prediction . . . . .	21
<b>5 Results &amp; Discussion</b>	<b>23</b>
5.1 Prediction Accuracy and Baseline Comparing . . . . .	23

---

<b>A Detail Optimization Process</b>	<b>25</b>
<b>Bibliography</b>	<b>26</b>

# Abbreviations

<b>NLP</b>	<b>N</b> atural <b>L</b> anguage <b>P</b> rocessing
<b>MLPs</b>	<b>M</b> ulti <b>L</b> ayer <b>P</b> perceptrons
<b>LSTM</b>	<b>L</b> ong- <b>S</b> hort <b>T</b> erm <b>M</b> emory network
<b>SGD</b>	<b>S</b> tochastic <b>G</b> radient <b>D</b> escent

*Dedicated to me*



# Chapter 1

## Introduction

Nowadays, it is the beginning of New-Age Predictive Analytics, the world of data scientists has turn their attention to Deep Learning and solved complex problem using data-driven approach. This is due to the fact that there are enormous amount data available to explore and huge advancements in computing power such that deep learning architecture can take place.

As we focus only on financial market, our works will fully base on the fact that financial market is informationally efficient, which was shown in Fama et al, 1956 [1]. We also realize that Natural Language Processing (NLP) technologies had huge advancement in recent years, such as *OLLIE*, an Open Information Extraction Software which can break the sentences into simple and processible structures, another technique introduced in [2] which can represent words in mathematical ways that computers can understand. However, the events in financial market are extremely spare and full of noise, false news and temporal effect come into the market and cause incorrect prediction. Therefore, instead using all events extracted from newspaper as proposed in previous works, we conduct our experiment only on a fraction of events which is much more dense and informative.

In our experiment, we train event embeddings model using the architecture proposed in Ding et al 2015 [3] but with some modification described in later chapter, this will learn mathematical ways to represent events in newspaper. For the movement price prediction, we use many difference deep learning models to compare and find the most favorable one.

## 1.1 Previous works

Recently, the deep learning approach has been successfully applied to tasks such as classification in [4], speech recognition in [5], and some famous deep machine learning approaches such as convolutional neural networks introduced above, deep belief networks in [6] and RNN - LSTM. However, its application in financial forecasting remain a relatively unexplored area. The challenge is that whether some deep machine model can be modified and learning to decrease the overall risk for trading strategies.

An example of deep learning in financial was introduce by Yoshihara et al. 2014 [7]. In their report, they used recurrent deep neural network to predict market trend by modeling temporal effects of past events. The same problem was approached in Ding et al, 2015 [3] where a combination of neural tensor network and a deep convolutional network was used to model short-term and long- term influences of events on stock price movements. This thesis is conducted base on their approach about events modeling, however, we use another architecture to model the price movement. This will shorten the training process and reduce the computation resource which make our model become easy to applied on small machine but still gain high accuracy on stock movement prediction.

## 1.2 Goal of thesis

To our knowledge, there are many deep learning model that were study and applied to many field. As the goal of our thesis, we will conducted some experiments on "the best reported results in the literature" as shown in Ding et al, 2015 [3] but with some modification. The changes we make during the experiment will only help us gain insight about the model, how the events were capture by the neural network and presented as driven-factor of stock movements.

Therefore, before going further into detail, we will first summarize some famous deep learning model in chapter 2 such as neural tensor network, deep convolutional network and recurrent deep neural network. These model largely contributed to the present success of the field artificial intelligence.

In addition, we also design other model to compare to the proposed model, this will help us realized some key factor that lead to the success of the model. The baseline models and the main model will be illustrated in chapter 4, Experimental Setup along with the detail architecture of the proposed model in chapter 3. As the final step of the thesis,

---

we discuss some further approach on the problem also some drawback of the model, how to improve it.

## Chapter 2

# Deep learning Models

### 2.1 What is Deep Learning ?

Deep learning can be considered a special type of machine learning. In order to deeply understand it, one must need a strong background in machine learning. Therefore, as the beginning of our thesis, we will provide the most important general principles about machine learning as warm up knowledge before going into detail about our work.

According to the book Machine Learning, Tom Mitchell (1997), a machine learning algorithm is an algorithm that learns from data as the definition "A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ ". We can imagine that how general machine learning algorithm is with this definition about task  $T$ , experience  $E$  and performance measure  $P$ . Hence, we will provide some examples of the tasks and their performance measures.

The task in most machine learning algorithms can be described as how a machine learning system should process an example also called as a sample in data. A sample mostly presents as a set of features which is measured from an object, an event that a machine can process. For instance, a picture is a set of vectors, each vector is a sequence of numbers measuring a value proportional to the light intensity at that particular location, mostly known as pixel. For simplicity, in this thesis, let us consider  $\mathbf{x} \in \mathbb{R}^k$  and each  $x_i$  is a feature. Below is some task that can be solved by machine learning:

- **Classification:** This task is the most fundamental task in machine learning problems. In this type of task, the machine must learn to classify which of  **$k$  categories** that the **input sample** belongs to. To be able to solve the task, the algorithm

must return a function  $f : \mathbb{R}^k \rightarrow \{1, \dots, k\}$ , the model  $y = f(x)$  can assign the class of the object encode as an index  $y \in 1, \dots, k$ . An example of this type of task is stock price trend forecasting in [8], input is an vectors of feature such as PE ratio, PX volume, PX EBITDA etc, and output is the movement prediction in next n-days encode as 0, 1.

- **Regression:** In this type of task, the machine is asked to predict a numerical value given some input. The machine learning algorithm solve this task by learn the function  $f : \mathbb{R}^k \rightarrow \mathbb{R}$ . An example of this task in finance is that predicting the expected claim amount that an insured person will make or predict the future prices of securities.
- **Feature learning:** In this task, feature learning or representation learning, the program is asked to transform raw data input to a presentation that other machine learning task can exploit effectively. This task is motivated by the fact that most machine learning task such as **Classification** often require the input to be computationally convenient to process. But some real world data such as words, literature, or speech are too complex and unique. Therefore, it is always recommended for any model in NLP to begin with a feature embedding process. An example of this task is the famous articles learn to transform words into dense vectors [2].

The second important part of a machine learning algorithm is the performance measure  $P$ . We must design a quantitative measure to evaluate the goodness of the algorithm, for instance, accuracy is often used as a measure for classification task. Usually, we use a **test set** to measure how well the model is on the unseen data and keep the training data separately.

The next part, which is no less important is the experience  $E$ . We often define  $E$  as the dataset. It is a collection of many sample, or more formally call as data points. By the characteristic of the dataset, we divide learning algorithm into two well known type:

- **Unsupervised learning algorithms:** The dataset in this type containing many features, and the jobs of the learning algorithm is to learn useful characteristic, and structure of the data. The classic machine learning task that fit on this type of algorithm is **Feature learning** task. Some other learning algorithm perform other tasks, like clustering, which is divide the origin dataset into clusters of similar data points.

- **Supervised learning algorithms:** experience a dataset containing features, but each data points is also associated with labels  $\mathbf{y}$ . A supervised learning algorithm can learn from data to classify objects into some categories.

Roughly speaking, unsupervised learning involves in learn from data points of a random vector  $x$  and attempting output the probability distribution  $p(x)$ , or some interesting properties of that distribution, while supervised learning involves learning to predict  $y$  from  $x$ , often by estimating  $p(y | x)$  where vector  $y$  is the associated value of random vector  $x$ .

### 2.1.1 Deep Feedforward Networks

**Deep feedforward networks**, also often called **feedforward neural networks**, or **multilayer perceptrons** (MLPs) is the most fundamental deep learning model. General goal of the model is to approximate some function  $f^*$ . For instance, a classification function  $y = f^*(x)$  map input vectors  $x$  to the label  $y$ . The model formally define as  $y = f(x; \theta)$  and the learning algorithm will learn the parameters  $\theta$  that output the best function approximation. The word *feedforward* originated from the fact that the information flows thought the model begin from vector  $x$  as input layer, then to the intermediate computational nodes used to define  $f$  and finally to the output  $y$ . There are no **feedback nodes** in which the computation nodes receive again its output from previous flows, these feedback nodes often known as **recurrent neural networks**, presented in subsection 1.1.4

Feedforward networks are the most important architecture to machine learning practitioners. They are the backbond of many follow architecture present later in this thesis. For instance, the convolutional networks used for recognizing object from images, presented in subsection 1.1.3.

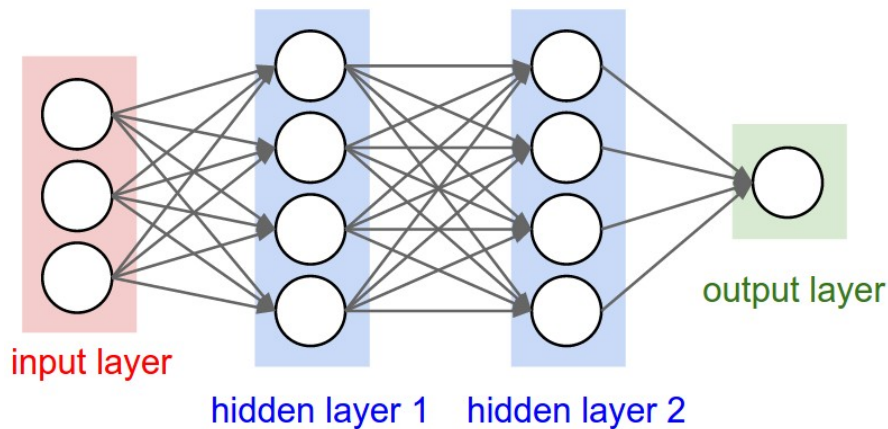


FIGURE 2.1: Feedforward networks model architecture

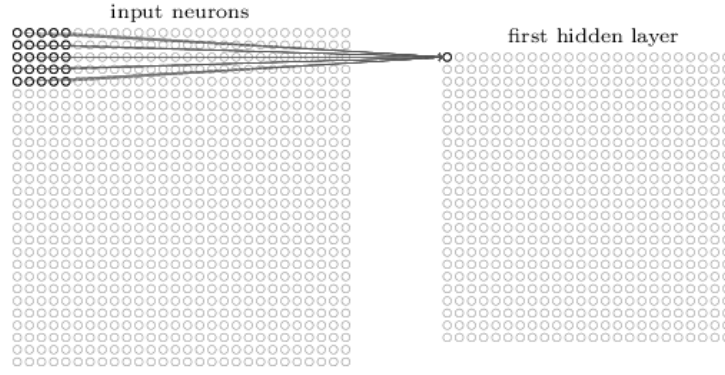
The word **network** describes the model as a composition of many different functions, they connect together and become a directed acyclic graph. A simple and most commonly graph can be written as a chain structure  $f(x) = f^3(f^2(f^1(x)))$ , where  $f^3, f^2, f^1$  are separately called as **output layer**, **hidden layer 2**, **hidden layer 1** and  $x$  is **input layer**. The overall number of layer gives the **depth** of the model. It is from this terminology that the name *deep learning* arises. Then, during the training process, the learning algorithm drive the  $f(x)$  to approximate  $f^*(x)$ . Last output  $f^*(x)$  of the training process will best classify the object  $x$  in to categories  $y$  base on the observation data points.

### 2.1.2 Convolutional Neural Networks

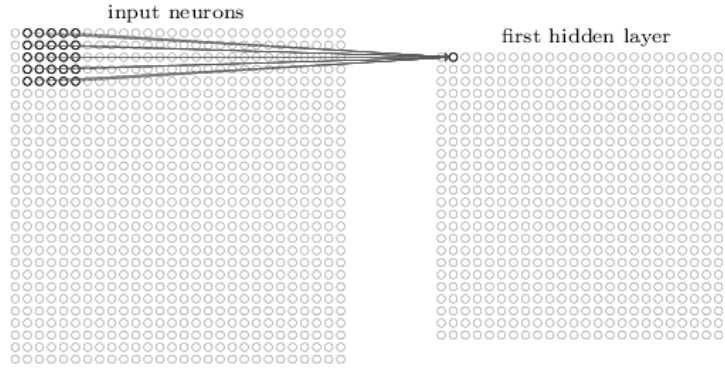
To describes convolutional networks architecture, we will first introduce a simple deep learning model, in which each layer is fully connected to the next layers as described in figure 2.1, with the function  $f^i(x) = \sigma(W_i^T x + b_i)$ ,  $i \in \{1, 2, 3\}$ ,  $W_i \in \mathbb{R}^{k \times h}$  where  $k, h$  is number of nodes in input, output layers separately and  $\sigma$  is an activation such as tanh function. This model work pretty well on simple input such as a feature vector  $x$  with no spatial structure. However, in real world problem, data points usually have spatial structure such as an images, price of many stocks in days or a sequence of events presented as vectors  $x_i \in \mathbb{R}^k$ ,  $i \in \{1, 2, \dots, N\}$ . Therefore, the convolutional networks (CNNs) was designed to take advantage of the spatial structure in data points.

CNNs was base on three basic ideas: local receptive fields, shared weights, and pooling. Each of these ideas are presented as follow:

- **Local receptive fields:** In figure 2.1, the inputs were depicted as a vertical line of neurons. In CNNs, input is usually presented as a matrix  $N \times N$  square of neurons (nodes), as the logic of feedforward networks model, we will connect the input neuron to a layer of hidden neurons, however, instead of connect all  $N \times N$  neurons to the next layers, many small, localized regions of the input matrix of neurons are connected with next neurons. To be more specific, each neuron in hidden layers next to the input will be connected to a small region of the input neurons, for example, a  $5 \times 5$  region of  $28 \times 28$  matrix of neurons. To illustrate this concretely, next figure contains a local receptive field in the top-left corner that connect to a neuron in next hidden layer:



Then we move the local receptive field to the right by one nodes (i.e., by one neuron), to create a second hidden neuron as below figure:



And so on, the first hidden layer was built next to the input matrix of neurons. In fact, different moving step which mostly known as stride length is used in other model.

- **Shared weights and biases:** Each hidden neuron has a bias and  $5 \times 5$  weights connected to its local receptive field in input matrix neurons in the example above, and the same weights and bias is used for all hidden neuron. In other words, for the  $j, k$ th hidden neuron  $h_{j,k}$ , it is calculated by:

$$\sigma \left( b + \sum_{l=0}^4 \sum_{m=0}^4 w_{l,m} a_{j+l, k+m} \right) \quad (2.1)$$

where  $\sigma$  is the neural activation function,  $b$  and  $w_{l,m}$  are the shared parameter across all local receptive field and  $a_{x,y}$  is the input neuron at position  $x, y$ . This show that at different location in the input matrix, the exactly same feature was detected and presents in the first hidden layer. For instance, if the input is a images, and  $b, w_{l,m}$  can recognize a circle in the local receptive field, this ability must be useful at other location, so all location should be applied with  $b$  and  $w_{l,m}$ . To make it more abstract, CNNs are well suitable to the location variation of a input matrix, move the matrix of neurons away, CNNs still detect the local feature in the matrix.



For this reason, the mapping from the input matrix to the hidden matrix is often called as a *convolutional kernels* with  $b$  and  $w_{l,m}$  are said to define a kernel or filter. Therefore, a complete convolutional layer contains many different convolutional kernels associated with different mapping as describe below, note that some deep architecture can use many more kernels than 3 kernels:

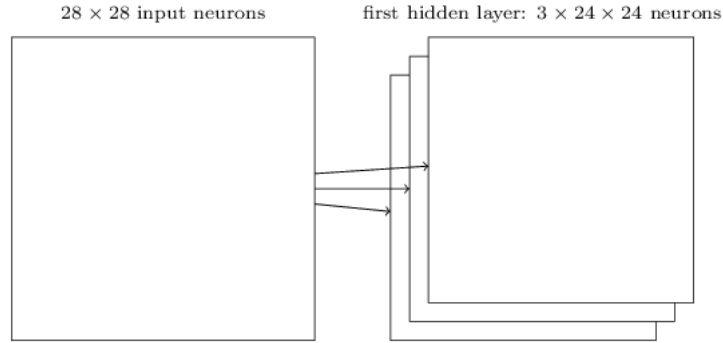


FIGURE 2.2: A complete convolutional layer

- **Pooling layers:** In addition to the convolutional layer, convolutional neural networks also consist of pooling layers. The jobs of pooling layers is to summarize the output from convolutional layer. In detail, each neurons in it summarize a region of  $n \times n$  neurons in the output from previous convolutional layer. A common procedure is max-pooling layer, whose each unit simplify outputs in the  $2 \times 2$  input region by the max function as illustrated in the following diagram:

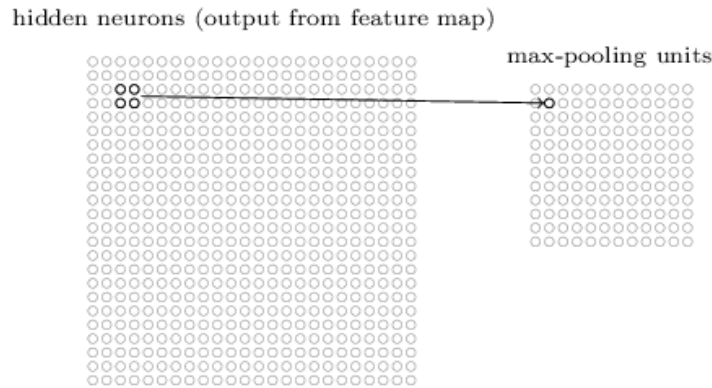


FIGURE 2.3: A pooling layer

We can interpret the pooling layer as a way to collect information from larger local field using information from smaller local receptive field from input layer. This method will output a condensed feature map. Another benefit from this layer is that it reduce computational power needed in later optimization by reducing number of parameters but still capture needed information.

We now can have a full understanding about a complete convolutional neural network. The architecture of the network can be visualized by a simple image:

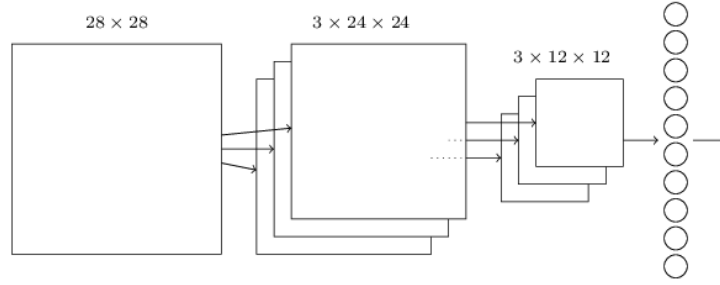
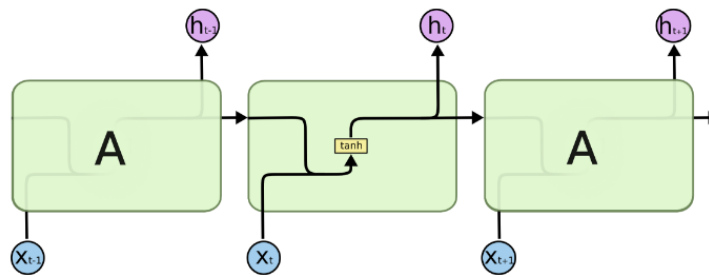


FIGURE 2.4: Network architecture

### 2.1.3 Recurrent Neural Networks

Human always need context to fully understand a word, for instance the word "bodies" in the sentence "Some celestial bodies, such as the planets and stars, can be seen with the naked eye". Words combine together to have other meaning, as well as events in newspapers in different days or sound when we talk. So, traditional feedforward network seem to fall when trying to address this issue, but Recurrent Neural Networks do not. They are networks that built from smaller and repeated network in them, allowing information to be learned repeatedly. In the last couple years, there have been many dramatic innovations in applying RNNs to AI problems such as: speech recognition, language modeling, translation, image captioning. Essential to these innovations is **LSTM** architecture, .



The repeating module in a standard RNN contains a single layer.

FIGURE 2.5: Standard RNN architecture

Long Short Term Memory networks mostly known as LSTM are a kind of RNN, capable of learning long-term relationship in sequence of samples. This network was introduced by Hochreiter and Schmidhuber in [9] and designed to avoid the long-term dependency problem in RNN. This problems arises when an information, in present time  $t$ , depend on another one, which is long time ago at time  $t - L$ , so that it slowly vanish

when the network processes other irrelevant information between time  $t$  and  $t - L$ . The figure 2.5 above will illustrate standard RNNs, where the repeating cell have a simple tanh layer structure. Each time the loop process the information  $x_t$  while considering all previous information embedded in  $y_{t-1}$ , then produce a output  $h_t$  and new state of information  $y_t$  which presented as the new summary of all processed information. The detail mathematical calculation is given below:

$$y_t = \sigma(W_y x_t + U_y y_{t-1} + b_y)$$

$$h_t = \sigma(W_h y_t + b_h)$$

Hence, we can realized that the information is embedded into vector  $y_t$  through each time step  $t$ , but the longer the time of the information the less likely that  $y_t$  can capture it, since the same parameters are applied at each step and the later information will overwrite the old one. Therefore, LSTM network has a different structure to address the issue, the cell of it is described in figure 1.6.

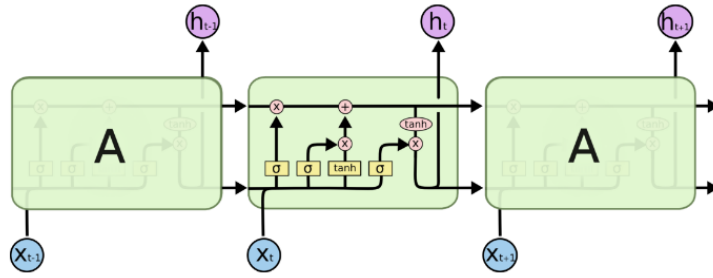


FIGURE 2.6: LSTM architecture

At first step, the module decide the information can be eliminated from the previous cell state  $C_{t-1}$  presented in the horizontal line running into every cell in the figure 1.6:

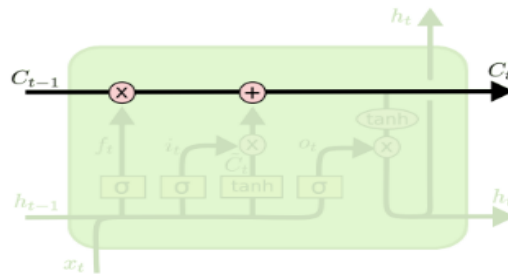


FIGURE 2.7: Cell state line

This is done by calculate the forgetting gate  $f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$  where  $\sigma$  is sigmoid activation function,  $h_{t-1}$  is previous output of the LSTM cell and  $x_t$  is the

present information. Then, the cell make a pointwise multiplication operation  $f_t * C_{t-1}$ , where  $f_t \in [0, 1]$ , a 0 represents "forget all previous information", and 1 represents "keep all previous information".

The next part of the cell is decide what new information will be stored in cell state. This is done by calculating two value, input gate  $i_t$  and the candidate state values  $\tilde{C}_t$ .

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (2.2)$$

$$\tilde{C}_t = \sigma(W_{\tilde{C}} \cdot [h_{t-1}, x_t] + b_{\tilde{C}}) \quad (2.3)$$

Then the cell will calculate the new cell state  $C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$ . Lastly, the cell continue to decide what is the output of this state though calculation of  $h_i$  given by:

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \quad (2.4)$$

$$h_i = o_t * \tanh(C_t) \quad (2.5)$$

What we illustrated so far is a standard LSTM. But not all LSTMs are the same as the above. A dramatic different LSTM is the Gated Recurrent Unit, or GRU, introduced by Cho, et al. (2014) [10]. It replace the forget and input gates by a single update gate. It also combines the cell state and hidden state in to a simple state, and makes some other changes in ouput gate. The GRU model is simpler than LSTM models, and has been became popular lately due to its fast convergence ability compares with LSTM in the training process. However, in this context, we use LSMT as ours main model because of the complexity of ours problems.

## 2.2 Optimization for training deep models

### 2.2.1 Learning process and traditional optimization

Optimization process used in machine learning for deep models have several difference with traditional optimization. In most machine learning case, the objective of the model is minimizing the performance measure  $P$  that use on the test set, but during the training process, only training data is used. Therefore, the learning algorithm use a indirect optimization process for  $P$ , it trained on  $\mathbb{J}(\theta)$ , a different cost function, and this will improve  $P$ . This is contrast to traditional approach, in which minimize performance measure  $P$  is the only goal.

Generally, the  $\mathbb{J}(\theta)$  function can be presented as an average value of all training samples, such as

$$\mathbb{J}(\theta) = \mathbb{E}_{(\mathbf{x}, y) \sim \hat{p}_{data}} \mathbb{L}(f(x; \theta), y) \quad (2.6)$$

This  $\mathbb{J}(\theta)$  defines the objective of the model with respect to the *data*, where  $\hat{y} = f(x; \theta)$  is the predicted value when  $\mathbf{x}$  is the input, the target output is  $y$  in the case supervised learning. The equation 1.2 can be re-written as:

$$\mathbb{J}(\theta) = \frac{1}{m} \sum_{i=1}^m \mathbb{L}(f(x^{(i)}; \theta), y^{(i)}) \quad (2.7)$$

The optimizing process relied on minimizing this objective is often called as **empirical risk minimization**. However, empirical risk minimization is prone to overfitting, the situation in which deep learning model can memorize training set and poorly perform on test set. In addition, some useful loss function, such as 0-1 loss, have the first derivatives is zero or undefined everywhere else, which is not useful. However, some effective modern optimization algorithms are based on gradient descent, so that in the context of deep learning, another slightly approach is taken.

We use another loss function called as **surrogate loss function**. For instance, the negative log-likelihood (NLL) of the correct class is used as an alternative for the 0-1 loss. NLL allows the model to calculate the conditional probability of target  $y$ , given the input  $x$ , then choose the class with higher probability, and also yield least error. In some cases, the NLL help the model to learn more about data, this can be seen as the 0-1 loss of test set continue to decrease even when 0-1 loss of training data has reached zero. This can be done when the learning algorithm continue to push the probability between the classes further apart and learn deeper feature in training data.

### 2.2.2 Basic learning algorithm

We have briefly described the objective of the learning process in previous section, but largely account for model correctness is the algorithm used to minimize the objective

function. Therefore, below presents the most used optimization algorithms for deep learning.

**Input:** Learning rate  $\epsilon_k$

**Input:** Initial parameter  $\theta$

**while** *Objective function value larger stopping criterion* **do**

Sample a minibatch of  $m$  examples from the training set  $\{x_1, \dots, x_N\}$  with corresponding targets  $y_i$ ;

Compute gradient w.r.t all parameters:  $\hat{g} \leftarrow +\frac{1}{m} \cdot \frac{\delta \sum_i \mathbb{L}(f(x_i; \theta), y(i))}{\delta \theta}$  .;

Apply update at step  $k$ :  $\theta \leftarrow \theta - \epsilon_k \cdot \hat{g}$

**end**

**Algorithm 1:** Stochastic gradient descent (SGD) update at training iteration  $k$

In SGD algorithm, the learning rate is considered as the most crucial hyper-parameter, this is due to the source of noise (the random sampling of  $m$  training examples) that keep the variance of data remain after the objective function arrives at minimum. Therefore, in real world problem, it is recommended to decrease the learning rate over iteration, such as a linearly decay of the learning rate until iteration  $\tau$ :  $\epsilon_k = (1 - \alpha)\epsilon_0 + \alpha\epsilon_0$  where  $\alpha = \frac{k}{\tau}$ .

While SGD remain as a statue of learning algorithm, in some complex deep learning model, its convergence took very long and huge computational resource. Therefore, the method of momentum is designed to overcome this issue in [11]. The momentum algorithm accumulates past gradients and continues to move in their direction. The update rule is given by:

$$v \leftarrow \alpha v - \epsilon_k \frac{\delta \sum_i \mathbb{L}(f(x_i; \theta), y(i))}{\delta \theta} \quad (2.8)$$

$$\theta \leftarrow \theta + v \quad (2.9)$$

The velocity  $v$  accumulates the gradient elements exponentially, the larger  $\alpha$  relative to learning rate at iteration  $k$  the more gradient affect the direction of parameters. As we see, momentum helps accelerate SGD in the relevant direction and dampens oscillations by include past gradients in the update. Now, the size of the step at each update depend on how consistent between gradients at each time step, this is how the name momentum derived, it is similar as velocity and force, in which cumulative gradients is considered as force that generate acceleration push the particles (the objective function) down the hill (learning curve) faster.

There are many variants of the momentum strategy, such as Nesterov accelerated gradient (1983) introduces in [12], in which the authors designed a strategy to address the issue of momentum on blindly updating parameters base on previous directions.

---

Another considerable strategy is Adaptive Moment Estimation (ADAM), it computes adaptive learning rates for each parameter which illustrated in [13] but also keeps an exponentially decaying average of past gradients. In this thesis, we use ADAM as our optimization algorithm for all the model.

## Chapter 3

# Methods

My work uses the idea in Ding et al. [3, 14], which was introduced as an empirical analysis of events embedding method. In this method, each event is presented as a feature vector in a compacted form to solve the problem about sparsity of event vectors. Then, machine training process will "learn" from these vector to predict the movement of the stock. The detail of my method will be explained deeply in following sections.

### 3.1 Data Preprocessing

One of the problems in NLP is noisy and huge data, to overcome this, unlike previous work of Ding et al., 2015 [3], we only chooses the most relevant news about the financial market and a specific company to conduct our experiment. This will make the model focus more on the relevant information.

However, with the constrains above, we realize that using the headline of news like previous works [3] will limits the training data. This brings us to the idea of using the body of the news - which used by newspaper as the summary of the relevant news. This step will make the raw data becomes more informative than using the headline of the news only.

### 3.2 Information Extraction

We extract the structured events from the raw text data using one of the most reliable Open Information Extraction technology - *OLLIE*, introduced in *Faderetal.*, 2012[15]. Given a sentences from the raw text data, *OLLIE* will extract events from the sentence and output the result as list of events  $E_i = (O_1, P, O_2), i \in \{1, 2, 3, \dots\}$  where,  $O_1$  and  $O_2$



are the objects that  $P$ , presented as the action in the sentence, is performed on. Then, I adds the published time  $T$  to the events for the training step.

The process can be explained in the following example. The sentence "Yahoo fell 0.6 percent to \$18.27 in late trading" in  $T = "2010 - 04 - 20"$  is modeled as: ("Yahoo", "fell 0.6 percent in", "late trading", "2010-04-20")

### 3.3 Embedding Vector of Information

As the result from above step, the events is mostly unique and represent a fraction of information. Therefore, it is necessary to embed the events into a dense vector but still contain that information. Follow the idea in [3], we use a neuron network model with the same architecture as they described in their paper to present events by vectors of feature extracted by the model.

#### 3.3.1 Model Architecture

The goal of their model is using machine learning to learn distributed representations of the events. The model's architecture is described in the Figure 2.1 below.

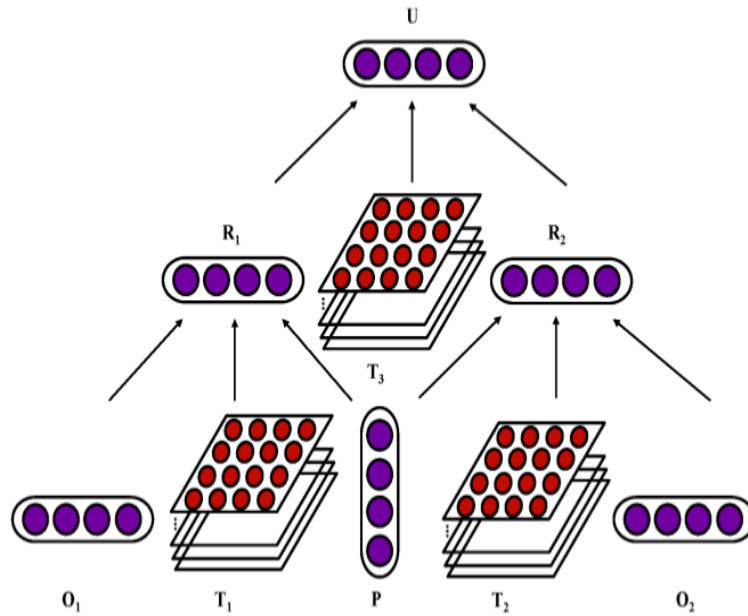


FIGURE 3.1: Model architecture

The input of model is word vectors and the output is an event vector. Since, all the events argument, objects and actions are lists of words, so they represent an object or

an action as the average of its words embedding,  $O_i \in \mathbb{R}^d$  and  $P \in \mathbb{R}^d$ . As illustrated in Figure 2.1, the roles of  $O_1$  and  $O_2$  are learned by two tensors,  $T_1$  and  $T_2$  independently,  $O_1T_1P$  and  $PT_2O_2$  will present as the role-dependent embeddings  $R_1$  and  $R_2$ , then, the third layer  $T_3$  is used to semantic compositionality over  $R_1$  and  $R_2$  and generate a complete structured U of the event. From the figure 3.1,  $R_1 \in \mathbb{R}^d$  is calculated by:

$$R_1 = f(O_1^T T_1^{[1:k]} P + W \begin{bmatrix} O_1 \\ P \end{bmatrix} + b) \quad (3.1)$$

where  $T_1^{[1:k]} \in \mathbb{R}^{d \times d \times k}$  and the bilinear product  $O_1^T T_1^{[1:k]} P$  is a vector  $r \in \mathbb{R}^k$ . Other parameter  $W$ , and  $b$  is a simple feed-forward network, in which  $W \in \mathbb{R}^{k \times 2d}$  and  $b \in \mathbb{R}^k$  and  $f$  is tanh function

### 3.3.2 Training Process for Events Embedding

The model using more than 5 thousands events from Bloomberg financial news only about Apple company from the year 2005 to 2009 as the training samples for modeling the events. The training data is smaller than previous work in [3] but much more specific about Apple company, and this time frame is chosen because it is when Apple release its first iPhone, which strongly impact the company in later years. The purpose is to model the events in this time frame and see how it affects Apple stock.

---

#### Algorithm 1: Event Embedding Training Process

---

**Input:**  $\mathcal{E} = (E_1, E_2, \dots, E_n)$  a set of event tuples; the model  $EELM$   
**Output:** updated model  $EELM'$

- 1 random replace the event argument and got the corrupted event tuple
- 2  $\mathcal{E}^r \leftarrow (E_1^r, E_2^r, \dots, E_n^r)$
- 3 **while**  $\mathcal{E} \neq []$  **do**
- 4      $loss \leftarrow \max(0, 1 - f(E_i) + f(E_i^r) + \lambda \|\Phi\|_2^2)$
- 5     **if**  $loss > 0$  **then**
- 6          $Update(\Phi)$
- 7     **else**
- 8          $\mathcal{E} \leftarrow \mathcal{E} / \{E_i\}$
- 9 **return**  $EELM$

---

The training algorithm is described as above, they assume that to modeling the events, the model have to be able to classify the true events  $E$  with corrupted one  $E^r$ . The corrupted event is created by replacing one of the objects in true events with a no

meaning objects randomly selected from the vocabulary of training data. Then, the objective of the model is to minimize the loss function below:

$$\text{loss}(E, E^r) = \max(0, 1 - f(E) + f(E^r)) + \lambda \times \|\theta\|_2^2 \quad (3.2)$$

where  $\theta = (T_1, T_2, T_3, W, b)$  is the set of parameters that need to update in every iterations to achieve the minimum of the objective function. The standard  $L_2$  regularization on parameters  $\theta$  except  $b$  is added by hyper-parameters  $\lambda$ . The algorithm will repeatedly looping over the training data, each loop, the model will calculate the loss function between the sample (one event tuple) and the corrupted one, if loss larger 0 then the model will update the parameters to minimize the loss, using the ADAM optimization algorithm [13] which base on the standard back-propagation algorithm. If loss equal 0, then skip the sample and continue to the next sample.

### 3.4 Stock Price Movement Prediction Model

After the learning process to represent the events by vectors, we model the effect of the events over the past week on stock prices based on the recurrent neural network (RNN).

The input sample is a sequence of daily event vectors in five days, which the daily vectors is the average of all events in that day. The output of model is the binary class represents the movement of stock price where 1 is increase and 0 is decrease.

The natural architecture model that fit extremely well on sequence data is RNN, so we use it as our second framework for predict stock movement. In previous paper, the Convolutional Neural Network (CNN) was used for these training data, but we recognize that our data is small and using such large architecture for different time spans (monthly, weekly, daily) to model specifically about Apple is costly. Therefore, a simple Long-Short Term Memory network (LSTM) is used on weekly time span, with smaller amount of parameters. The result of our model in test data is very competitive with the previous work, 64.44% accuracy using LSTM.

Formally, given a sequence of daily event vectors  $U = (U_1, U_2, U_3, U_4, U_5)$ . For each time step  $i$  with input  $U_i$  and previous hidden state  $h_{i-1}$  we compute the updated hidden state  $h_i, C_i = \text{LSTM}(x_i, h_{i-1}, C_{i-1})$ . The last hidden state is of the LSTM layer is used in the next feed-forward layers. The output of feed-forward layer is the prediction of the model which is a vector  $y \in \mathbb{R}^2$  represent the score for 2 class increase/decrease of the Apple stock price.

## Chapter 4

# Experimental Setup

### 4.1 Data Collection

We collect news from Bloomberg news API over the periods from January 2005 to January 2010, and instead using all news titles, we searched news with the key word "Apple". The result is a list of news titles with the paragraph contains the relevant information. We collected the paragraphs and break it into sentences, then we extracted events using only those sentences. The events output from *OLLIE* were then filtered by its score from *OLLIE*, events with score less than 0.7 will be filtered out. Detail about our training data is shown in table 4.1 below.

	Number of sample
New paragraphs with key word "Apple"	34514
Events extracted from paragraphs	57502
Filtered Events	6857

TABLE 4.1: Table to summarize training Events

The experiments in our paper focus on predicting the Apple stock prices obtaining from Yahoo Finance. We collect the stock price data and divide into three groups for training, validating and testing in the stock trend prediction model. The summary of data set is illustrated in table 4.2.

	Number of sample (days)
Training Data	727
Validating Data	80
Testing Data	90

TABLE 4.2: Table to summarize training stocks model

## 4.2 Baseline stock prediction models and proposed models

The baseline models are some simple model using standard feedforward neural network on the flatten event embeddings vectors in 5 days and CNN model on sequence of event vectors in 5 days. These model were described in detail as baselines model for comparing:

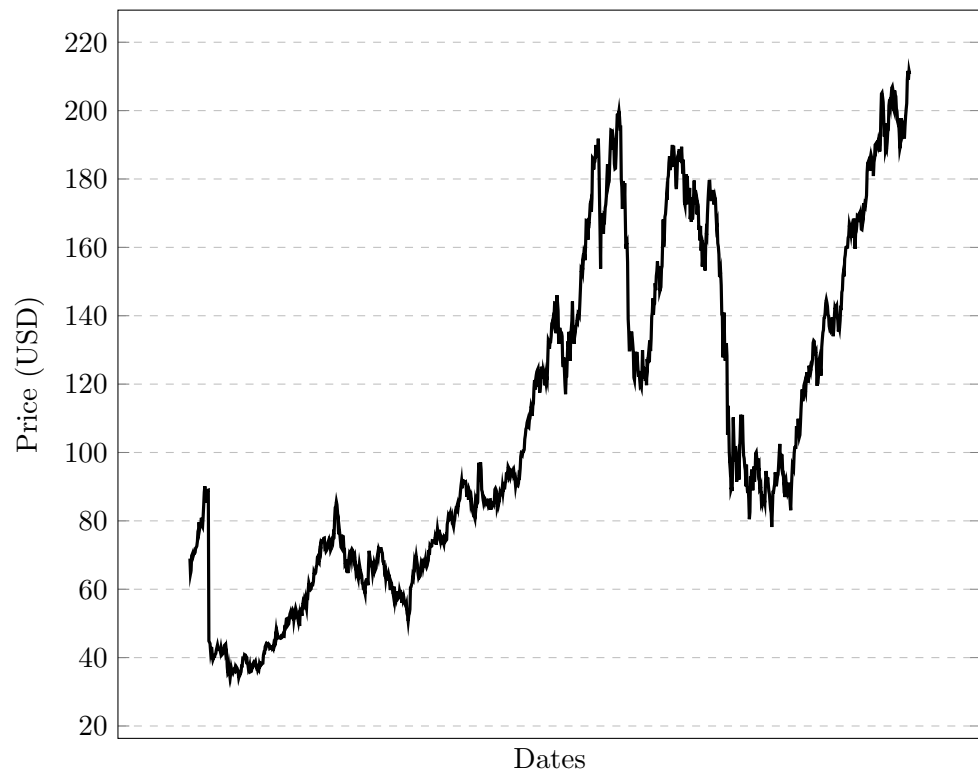
- Model 1: In this model, we use a simple fully-connected neural network, each  $t$  days samples were flatten and fully connects to one node in next hidden layer. This model contains 2 hidden layers with 512 and 1024 nodes separately, using sigmoid activation function, the output is 2-d vectors presented as the prediction. During training process, the dropout layers were added for regulation purpose.
- Model 2: The second base model was design base on CNN architecture, it contains 2 convolutional layer, each layers has 128 kernels and one max-pooling layer. The training process is the same as the first base model.

In contrast to the baselines, we use the a neural tensor network to learn event vectors as in [3] but then use a deep Long-Short term memory network to build the price movement prediction model. This is ours proposed model, note that all model were trained using ADAM optimization process. The detail result of those baselines and proposed model are shown by the table in the next chapter.

## 4.3 Apple Stock Prediction

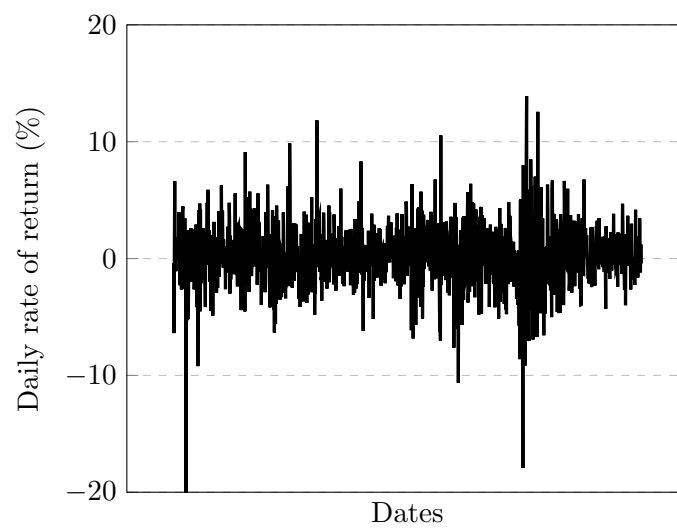
Apple Inc is an American multinational technology company that designs, develops, and sells consumer electronics, computer software, and online services. The success of the company were build on their reputation, their enormous user, and amazing marketing plan. These reason are important factors that make us decided to choose AAPL stock price, because this show that AAPL is strongly impacted by news, important events and sensitive to the behaviors of their customers.

Apple Stock Price from 2005 to 2010



[h!]

Apple Stock rate of return from 2005 to 2010



## Chapter 5

# Results & Discussion

### 5.1 Prediction Accuracy and Baseline Comparing

Model name	Accuracy(%)	Recall(%)	Precision(%)	F-measure(%)
Fully-connected neural network 5-days	56.67	60.7	62.0	60.1
Convolutional neuron network 5-days	60.0	0.5717	0.5352	0.5296
Proposed Model: LSTM 5-days	64.4	76.3	66.04	70.2

TABLE 5.1: Table to summarize training stocks model

As shown in table 5.1, the proposed model using LSTM achieved highest performance compare to other model. This shown that the LSTM is able to learn the information embedded into events vectors and takes advantaged of the time structure in data to predict the stock price movement.

In spite of small training sample, the CNN model achieved very good performance as pointed out by previous work in [3], the model capture the events vectors and predict with 60% accuracy. However, we have to mention that this model only use 5-days events time frame as different from their previous works, which use both 5-days and 30-days time frames as input for the model.

Therefore, we conclude that although the long-term time frame events might contribute to the performance of the model, short-term time frame events still be considered as key-driven factor for the model predictions.

Good performance of the model on real data shown that it can be applied in the real market as auto trading robot. The output of the model can be converted into a vector of probabilities presents the movement of stock price. Then, from the most simple trading strategy to the most complex can be applied to make use of those events vectors. Applying trading strategies is out of the scope of the thesis, but the model shown good performance, hence, it is promising to make further experiments of the model as the auto trading robot.



## Appendix A

# Detail Optimization Process

# Bibliography

- [1] Eugene F Fama. The behavior of stock-market prices. *The journal of Business*, 38 (1):34–105, 1965.
- [2] Jeffrey Pennington, Richard Socher, and Christopher D Manning. Glove: Global vectors for word representation. In *EMNLP*, volume 14, pages 1532–1543, 2014.
- [3] Xiao Ding, Yue Zhang, Ting Liu, and Junwen Duan. Deep learning for event-driven stock prediction. In *IJCAI*, pages 2327–2333, 2015.
- [4] Honglak Lee, Roger Grosse, Rajesh Ranganath, and Andrew Y Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *Proceedings of the 26th annual international conference on machine learning*, pages 609–616. ACM, 2009.
- [5] Xiao-Lei Zhang and Ji Wu. Deep belief networks based voice activity detection. *IEEE Transactions on Audio, Speech, and Language Processing*, 21(4):697–710, 2013.
- [6] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. A fast learning algorithm for deep belief nets. *Neural computation*, 18(7):1527–1554, 2006.
- [7] Akira Yoshihara, Kazuki Fujikawa, Kazuhiro Seki, and Kuniaki Uehara. Predicting stock market trends by recurrent deep neural networks. In *Pacific Rim International Conference on Artificial Intelligence*, pages 759–769. Springer, 2014.
- [8] Yuqing Dai and Yuning Zhang. Machine learning in stock price trend forecasting.
- [9] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [10] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.

- [11] Ning Qian. On the momentum term in gradient descent learning algorithms. *Neural networks*, 12(1):145–151, 1999.
- [12] Yurii Nesterov. A method for unconstrained convex minimization problem with the rate of convergence  $O(1/k^2)$ . In *Doklady an SSSR*, volume 269, pages 543–547, 1983.
- [13] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [14] Xiao Ding, Yue Zhang, Ting Liu, and Junwen Duan. Using structured events to predict stock price movement: An empirical investigation. In *EMNLP*, pages 1415–1425, 2014.
- [15] Michael Schmitz, Robert Bart, Stephen Soderland, Oren Etzioni, et al. Open language learning for information extraction. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 523–534. Association for Computational Linguistics, 2012.