

# Improving the Stock Market Prediction with Representation Learning



Name: Ding Fan

BUPT number: 2016213516

QM number: 161188243

Class: 2016215119

## **Contents**

- 1. Introduction**
- 2. Background**
- & Design and Implementation**
- 3. Results and Discussion**
- 4. Conclusion and Further work**

## 1. Introduction: Purpose

Learning representations of the data: make it easier to extract useful information when building classifiers or other predictors. Such as: word2vec (embedding words to vectors), node2vec (embedding graph nodes to vectors)

previous studies: as a time series problem-based on the historical price data (solution: RNN)

NEWS: "Microsoft successfully completed the acquisition of company xxx"

Which stocks' price will be affected?

- Microsoft ↑
- The supply companies ↑  
The customer companies ↑

- quantitative trading data : trading price
- qualitative descriptions(un-quantitative) : news and reports.
- relationship information (un-quantitative) : stock relationship Graph

Learn representations  
that encode these stock-related information



# 1. Introduction: Target

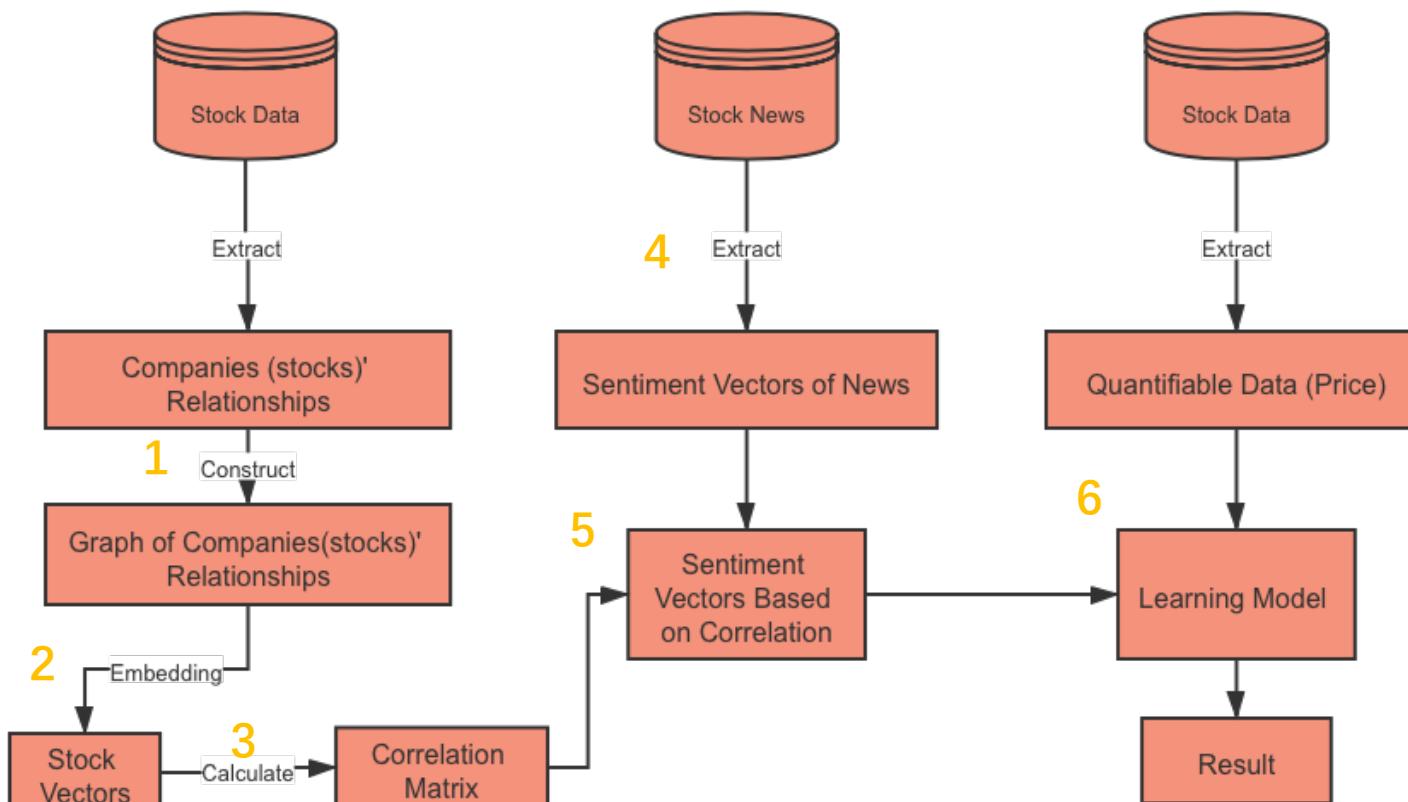
## Tasks:

1. Data preprocessing and feature extraction (**Finished**)
2. Representation learning model building (**Finished**)
3. Learning method design and implementation (**Finished**)
4. Software with interface that report the results and visualize the representations (**Finished**)

## Outcomes:

- Representations that encode the stock-related information
- The prediction algorithm and the prediction results
- Working software that can visualize the representations

## 2. Design and Implementation Overall Framework



1. Construct the graph of companies' relationships. **word2vec principle**
2. Embedding the graph of companies' relationships into vectors to generate the stock vectors. **node2vec**
3. Calculate the similarity between every2 stocks and construct the correlation matrix of 50 stocks.
4. Get the sentiment score of each stock news
5. Use the correlation matrix to update the sentiment score of the target stock
6. Input the quantifiable data and the updated sentiment score of target stock to learning model **LSTM**

## 2. Background-1 word2vec

# 1-Stock Relationship Graph Construction Model

### Word2vec: learning low-dimension representation of words

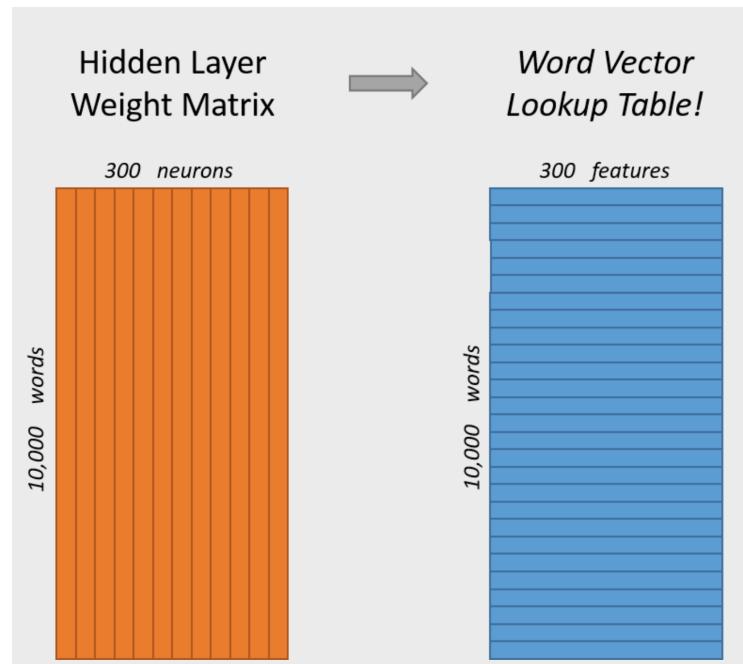
- Represent each word with a vector
- Has better performance when need to calculate the similarity between different words or find relevant words.

- Input: a set of training sentences (**n** different words)
- Step 1: generate training samples from each sentence
- Step 2: train 2-layer network (hidden layer+ output layer)
- Output: **n** vectors represents **n** words (**Hidden layer**)

Source Text	Training Samples
The quick brown fox jumps over the lazy dog. →	(the, quick) (the, brown)
The quick brown fox jumps over the lazy dog. →	(quick, the) (quick, brown) (quick, fox)

Trained

word	One-hot	Embedding
King	0001	[0.99,0.99,0.55,0.7]
Queen	0010	[0.99,0.05,0.93,0.6]
woman	0100	[0.02,0.01,0.99,0.5]
princess	1000	[0.98,0.02,0.94,0.1]



## 2. Design and implementation

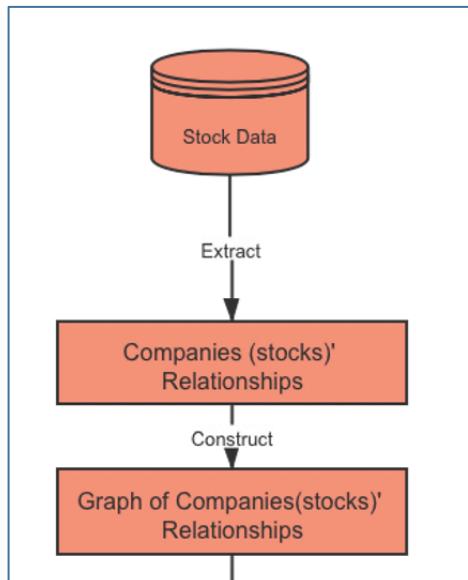
### 1-Stock Relationship Graph Construction Model

similarity of Stock price fluctuation  $\xrightarrow{\text{word2vec}}$  stock relationship graph

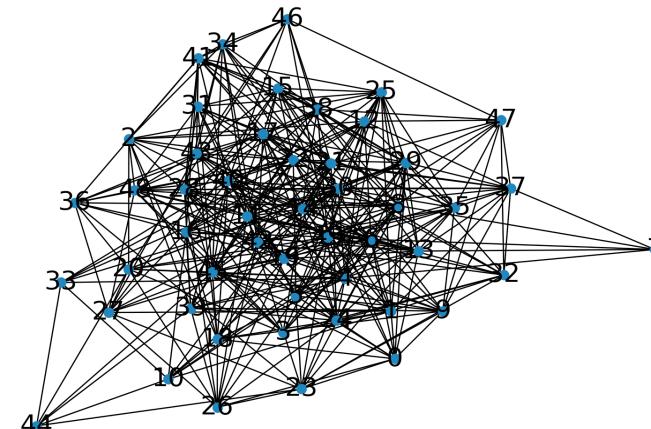
For every 2 stocks: if price fluctuate **more similarly**, the 2 stocks appear together in stock sentence **more frequently**

- **Input: 50 stock prices** (244 trading days in 2015)
- **Step1:** Each day, put stocks name in order of price change from high to low to generate 244 sentences  

- **Step2:** based on 244 sentences, use word2vec to get the **50 vectors (representation of 50 stocks)**
- **Step3:** calculate the cosine similarity **between every 2 stocks** to construct the **stock relationship graph**.



if **similarity exceeds the threshold value**, set an edge between the 2 stock nodes.



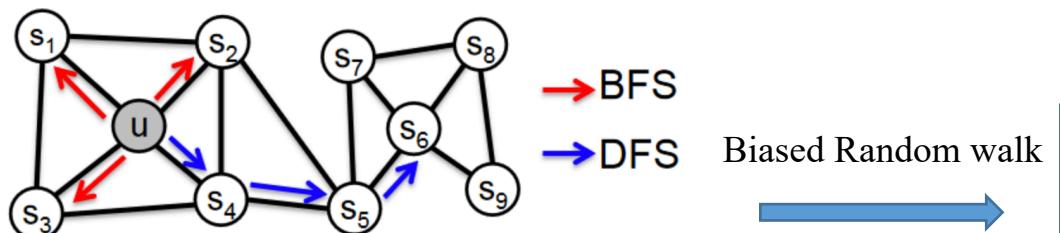
## 2. Graph Embedding Model

### 2. Background-2 node2vec

- **Input:** a Graph with **n** nodes
- **Step1: Biased Random walk** on input graph to generate a set of **node sequence sentences for training**
- **Step2: word2vec**
- **Output:** **n** vectors represents **n** nodes in the Graph

#### Graph Embedding: node2vec (derived from word2vec)

- Each vector represents a node in graph
- contains the characteristics of the node itself
- Also, characteristics of the relationship with other nodes



Sentence\_1: u S<sub>4</sub> S<sub>5</sub> S<sub>6</sub>  
Sentence\_2: u S<sub>2</sub> S<sub>1</sub> S<sub>3</sub>

like  
=

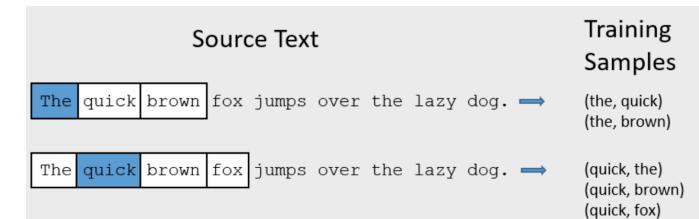
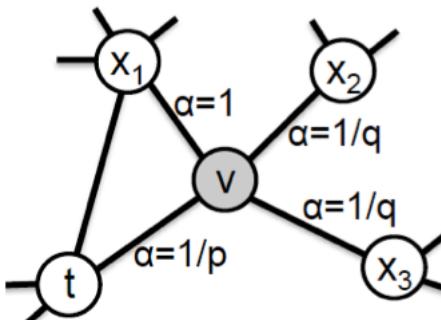


Figure 1: BFS and DFS search strategies from node  $u$  ( $k = 3$ ).



#### Biased Random walk

$$\pi_{vx} = \alpha_{pq(t,x)} \cdot w_{vx}$$

$$\alpha_{pq}(t, x) = \begin{cases} \frac{1}{p} & \text{if } d_{tx} = 0 \\ 1 & \text{if } d_{tx} = 1 \\ \frac{1}{q} & \text{if } d_{tx} = 2 \end{cases}$$

- $p$ : controls the probability of revisiting a node in the walk
- $q$ : controls the probability of exploring “outward” nodes



## 2. Design and implementation

- Sentiment score of stock x on i-th day is denoted as  $\text{sentiment}_i^x$
- Quantifiable features is denoted as  $Q_i^x = \{q_i^{x_1}, q_i^{x_2}, \dots, q_i^{x_n}\}$
- Similarity between stock x, y is  $S_{x,y}$
- Use the correlation matrix to update the sentiment score of stock x by:

$$\widetilde{\text{sentiment}}_i^x = \sum_{y=0}^n \text{sentiment}_i^x * S_{x,y}$$

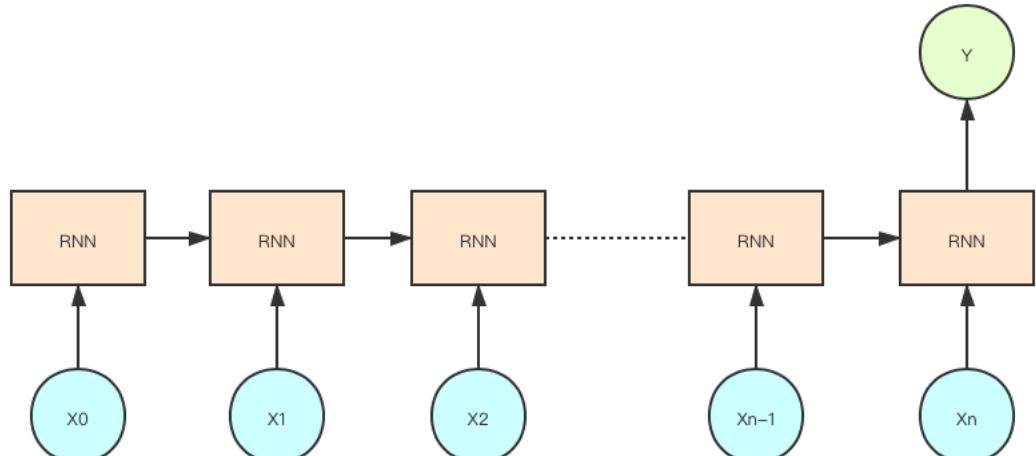
- $LSTM_{input} = \{(\widetilde{\text{sentiment}}_1^x, Q_1^x), (\widetilde{\text{sentiment}}_2^x, Q_2^x), \dots, (\widetilde{\text{sentiment}}_n^x, Q_n^x)\}$

## 3. Learning Model Based on LSTM

a kind of RNN used to process sequence information data.

Quantifiable data  
Correlation matrix  
News data

**Sentiment score of news:** use the NLP Python-SDK provided by Baidu-ai, the reason is that it has a well pre-trained model.



- n refer to the length of the time-series
- Prediction of stock x on n+1-th day, based on n previous days
- $X_n = (\widetilde{\text{sentiment}}_n^x, Q_n^x)$
- Y is the prediction of the stock x on n+1-th day(1: ; 0: )

## 2. Design and implementation

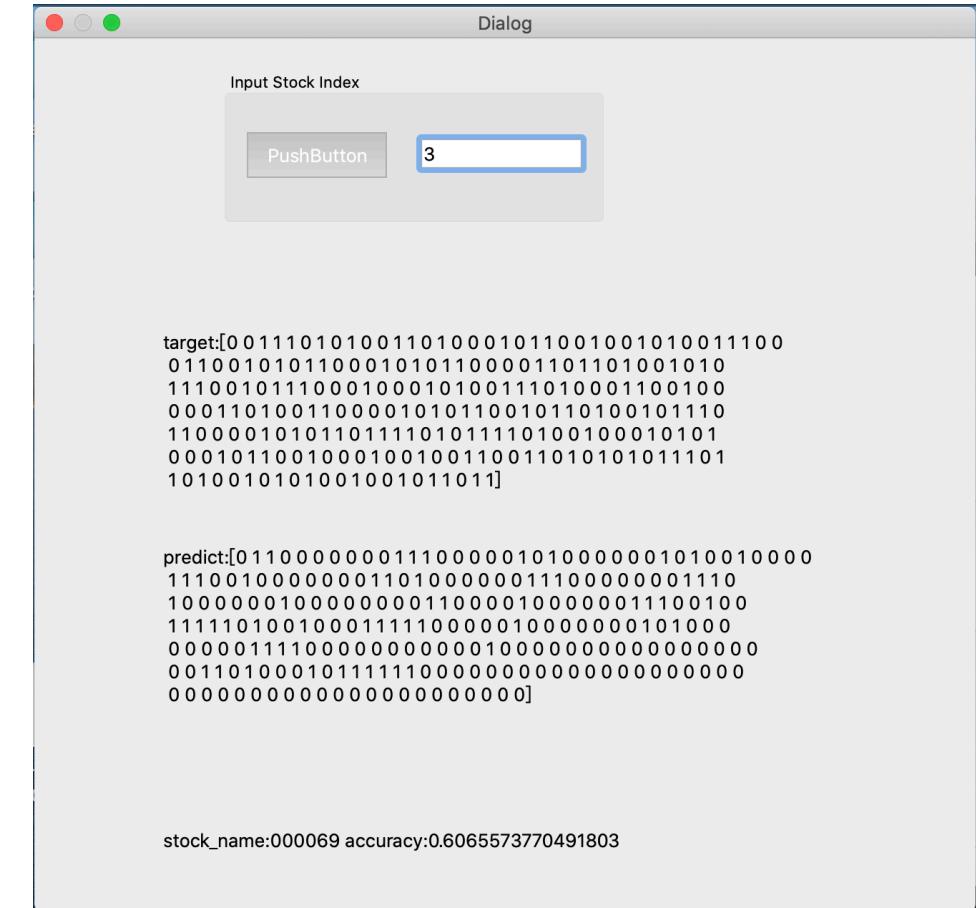
## 4. Working Software

Dialog

Input: 'all' or num in [0,50)

A network diagram illustrating connections between numbers from 0 to 45. The central node is connected to 15 other nodes. Nodes 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, and 15 are shown as blue circles with their values in red. Other nodes are represented by small black dots.

## Window to display the Representation of Stock Graph



## Window to show the prediction result

### 3. Experiments Results

#### 2. Setup

Group	Description
Group1 (Quantifiable)	give prediction only depends on the quantifiable data like stock prices.
Group2 (Quantifiable + News)	give prediction depends on the quantifiable data like stock prices and the news which reflect the sentiment
Group3 (Quantifiable + News +Stock Relationships)	give prediction depending on the quantifiable data like stock prices, the news which reflects the sentiment and the stock relationships.

- n: Prediction of each stock on n+1-th day, based on n previous days  
• Accuracy (ACC) and F-1 score of following table is the arithmetic mean of 50 stock models

	n=2		n=3		n=4		n=5		n=6		n=7	
	ACC	F-1										
Group 1	0.5209	0.4807	0.5175	0.4891	0.5213	0.4958	0.5203	0.4980	0.5209	0.4883	0.5180	0.4842
Group 2	0.5297	0.5180	0.5313	0.5115	0.5331	0.5123	0.5351	0.5194	0.5402	0.5217	0.5331	0.5170
Group 3	<b>0.5473</b>	<b>0.5438</b>	<b>0.5514</b>	<b>0.5441</b>	<b>0.5547</b>	<b>0.5445</b>	<b>0.5532</b>	<b>0.5432</b>	<b>0.5443</b>	<b>0.5406</b>	<b>0.5386</b>	<b>0.5366</b>

## 4. Conclusion and Further work

C:

- Using Representation Learning, helps improve the stock prediction performance.
- However, there also are some problems and weaknesses which need to be solved and improved in the future or further work.
- Constructing the stock relationship graph is based on the price fluctuation similarity, which maybe cannot accurately reflect the relationships between different stocks.

F. Try exploring other methods to construct the stocks' relationship or finding a well-ready make graph to improve the stock prediction performance in our model.

**Thanks for Watching**