# Big Data Application Development - Team Project Proposal

## Part 1. General Information

**Team Members:**

**Lizi Chen, Chun-Yi Yang**

**Project Title:**

**Stock Price Prediction by Financial News Sentiment Analysis**

**Project Description:**

*In this project, we are going to build a stock price prediction model that takes in real-time financial news and produces investment suggestion such as buy or sell.*

**Who is a typical user of your application:**

1. Individual Investor

2. Portfolio Manger

3. Stock Market Learner and Researcher

4. Business Strategist

**What insight will you derive from the data?**

1. By verify the ground truth and doing the sentiment analysis on financial news, we can understand sometimes a good news may not be an immediate boost for the stock price, and sometimes bad news does not add enough fluctuation to the market either.

2. By measure the time lag between a published news and how the stock price may change, we can tell the category of news that have a stronger assertion to provide more fluctuation in stock price or less.

3. By generally train the model to make investment suggestions according to news, we can see how blast of news and its language can affect the market.

**What actuation decisions can be made based on the actionable insight?**

1. buy a stock

2. sell a stock

# Big Data Application Development - Team Project Proposal

## Part 2. Data Source Information

### Name of Data Source 1: Historical Intra-day Stock Price

### Data Source Description:
We are going to retrieve historical intra-day stock price dataset according to the retrieved news. Each business entity will have its own CSV data file.
https://stooq.com/db/
https://www.google.com/finance/getprices?i=60&p=2d&q=IBM
https://www.quantshare.com/sa-426-6-ways-to-download-free-intraday-and-tick-data-for-the-us-stock-market

| Data Collection Frequency | Data Size | Data Frequency |
|---|---|---|
| • Are you collecting the data in realtime? Or are you collecting it periodically?<br>• Are you collecting static data? (e.g. historic data that you load once) | • Estimate size of the data you will store, e.g. MB, GB, TB, PB | • If realtime data, what is the frequency and volume of data (how often and how much data will you collect at a time)?<br>• If batch data, how often will you collect it? |
| ☐ **Realtime (ongoing near-realtime collection)**<br><br>☐ Batch (multiple non near-realtimecollections)<br><br>☐ Static (one time collection) | ☐ **MB**<br><br>☐ 1-10 GB<br><br>☐ 10-100 GB<br><br>☐ 100-300 GB<br><br>☐ 300-500 GB<br><br>☐ > 500 GB | **If realtime data:**<br>• How often will you collect data?<br>  ☐ Every second, or every few seconds<br>  ☐ **Every minute, or every few minutes**<br><br>• What is the size of data you will collect at each interval?<br><br>_____**Very minimal KB size data**_____<br><br>**If _not_ realtime data:**<br>• Will you collect a batch of data periodically or just once (static)?<br>  ☐ Just once<br>  ☐ Every hour, or every few hours<br>  ☐ Every day, or every few days, or every week, or every month<br><br>• How much data that will be collected at each interval?<br><br>_____ |

# Big Data Application Development - Team Project Proposal

## Part 2. Data Source Information

**Name of Data Source 2:** Reuters News Archive

**Data Source Description: We will retrieve Reuters News from Reuter's New Archive for the past several months. Also, we will keep tracking the upcoming new data for final product and for updating the new model.**

| Data Collection Frequency | Data Size | Data Frequency |
|---|---|---|
| • Are you collecting the data in realtime? Or are you collecting it periodically?<br>• Are you collecting static data? (e.g. historic data that you load once) | • Estimate size of the data you will store, e.g. MB, GB, TB, PB | • If realtime data, what is the frequency and volume of data (how often and how much data will you collect at a time)?<br>• If batch data, how often will you collect it? |
| ☐ **Realtime (ongoing near-realtime collection)**<br><br>☐ Batch (multiple non near-realtimecollections)<br><br>☐ Static (one time collection) | ☐ MB<br><br>☐ **1-10 GB**<br><br>☐ 10-100 GB<br><br>☐ 100-300 GB<br><br>☐ 300-500 GB<br><br>☐ > 500 GB | **If realtime data:**<br>• How often will you collect data?<br>  ☐ Every second, or every few seconds<br>  ☐ **Every minute, or every few minutes**<br><br>• What is the size of data you will collect at each interval?<br><br>  _____**Around 20 MB for raw HTML files**_____<br><br>**If _not_ realtime data:**<br>• Will you collect a batch of data periodically or just once (static)?<br>  ☐ Just once<br>  ☐ Every hour, or every few hours<br>  ☐ Every day, or every few days, or every week, or every month<br><br>• How much data that will be collected at each interval?<br><br>  _____ |

# Big Data Application Development - Team Project Proposal

## Part 2. Data Source Information

**Name of Data Source 3:** WSJ News Archive

**Data Source Description:**
Similar to the Reuter News Archive, we will retrieve WSJ news for the past several months for model construction as well as keep tracking fresh news for further model refinement.

| Data Collection Frequency | Data Size | Data Frequency |
|---|---|---|
| • Are you collecting the data in realtime? Or are you collecting it periodically?<br>• Are you collecting static data? (e.g. historic data that you load once) | • Estimate size of the data you will store, e.g. MB, GB, TB, PB | • If realtime data, what is the frequency and volume of data (how often and how much data will you collect at a time)?<br>• If batch data, how often will you collect it? |
| ☐ **Realtime (ongoing near-realtime collection)**<br><br>☐ Batch (multiple non near-realtimecollections)<br><br>☐ Static (one time collection) | ☐ MB<br><br>☐ **1-10 GB**<br><br>☐ 10-100 GB<br><br>☐ 100-300 GB<br><br>☐ 300-500 GB<br><br>☐ > 500 GB | **If realtime data:**<br>• How often will you collect data?<br>  ☐ Every second, or every few seconds<br>  ☐ **Every minute, or every few minutes**<br><br>• What is the size of data you will collect at each interval?<br><br>_____**Similar to Reuter, around 20 MB for Raw HTML files.**_____<br><br>**If _not_ realtime data:**<br>• Will you collect a batch of data periodically or just once (static)?<br>  ☐ Just once<br>  ☐ Every hour, or every few hours<br>  ☐ Every day, or every few days, or every week, or every month<br><br>• How much data that will be collected at each interval?<br><br>_____ |

| Big Data Application Development - Team Project Proposal |
| --- |
| *Pending Data Source (To Be Evaluated)* |
| *1. GDELT World News Real-time Archive* |
| *2. CityFalcon Rated News Archive* |
| |
| |