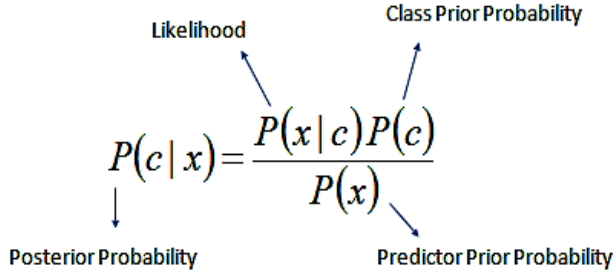


[illegible]



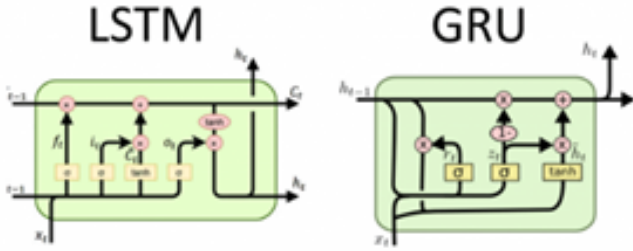
$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

Fig. 1. Snownlp is based on native bayes method to classify the sentence into sentiments score between 0-1.

B. Technical Index

We gathered data of TSMC from 2018 to 2020, which contains news and technical analysis factors. There are several technical analysis factors that we use, such as moving average(MA), open, close, high, low, KD, RSI, MACD, standard deviation of stocks. The data is normalized for training.

We propose to apply an LSTM model because it is powerful with sequence prediction problems, and the LSTM model could store past information. There will be total 14 features in our modeling input. Furthermore, we determine to apply GRU, a LSTM with a forget gate, but has fewer parameters than LSTM. GRUs have been shown to exhibit better performance on certain smaller and less frequent datasets.



III. EVALUATION

Result

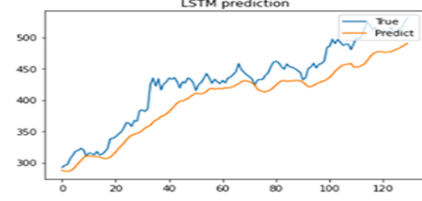


Fig1 : Use technical index and sentiments to predict close via LSTM

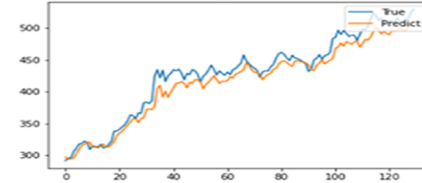


Fig2 : Use technical index and sentiments to predict close via GRU

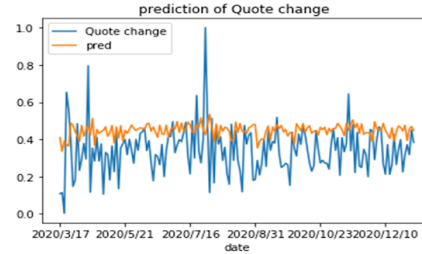


Fig3 : Use snownlp to predict quote change via GRU

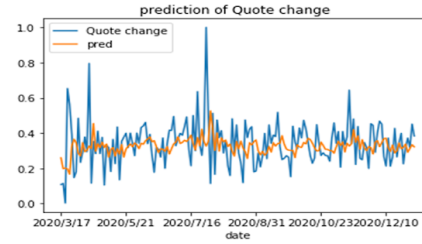


Fig4 : Use bert to predict quote change via GRU

In our experiments result, we find that the input features including 'BERT' and 'SnowNLP' can make the predict closing price closer to the true closing price. And GRU model perform better than LSTM in this case, but both of them cannot predict the true closing price perfectly. As the result, we hope to obtain the correlation between stock quote change ratio and sentiments score. Figure 3 and Figure 4 shows that BERT has higher correlation with closing price than SnowNLP, we consider that the main reason is because the default pre-training dataset in SnowNLP is about buying/selling comments, there must lost some financial words and terms in the dataset. However, BERT is pretrained by the social posts on Twitter, which contains lots of categories of word, there may involve financial terms in the pre-train dataset.

IV. DISCUSSION

2020 is a tough year for every investor and stock market, if we hope to get more precise performance, there are some improvement we can discuss in the future.

First, we only obtain 465 days PPT information in dataset, and split 300 days for training set, the others for testing set. The small dataset may cause training difficult, we shall crawl more data to revise the result. Second, 14 features are not enough for predicting stock price obviously, a better way for natural language processing is to use word-vector method to do analysis, using more vectors and features to make prediction may have different result. The last improvement is to train SnowNLP dataset by our stock news, in order to make the sentiments score has more correlation with the true stock price effect.