



Predicting Future Financial Volatility from Financial Reports with SEC 10-K Report Benchmark

**Team : Zhong Qiaoyang, Zhang Cheng
Wu Bingxun, Wu Yifan, Zhang Zhiyin**

Our Project Overview



01 Main goal of our project

Taking corporate annual financial reports as the object of analysis, this study focuses on exploring the predictive ability of textual features on future stock price volatility. Based on empirical tests of machine learning regression models, we systematically evaluate the differences in the predictive effectiveness of two types of textual representations: sparse features (BOW Model) and dense features (BERT, Word2Vec like).

02 Expected Result

In the representation extracted from textual datasets, a novel representation that combines TF-IDF with LongFormer, can effectively captures significant information from specific financial terminology to improve TF-IDF remains context-agnostic. Because LongFormer provides a context-based representation. This combined may approach demonstrates enhanced performance.

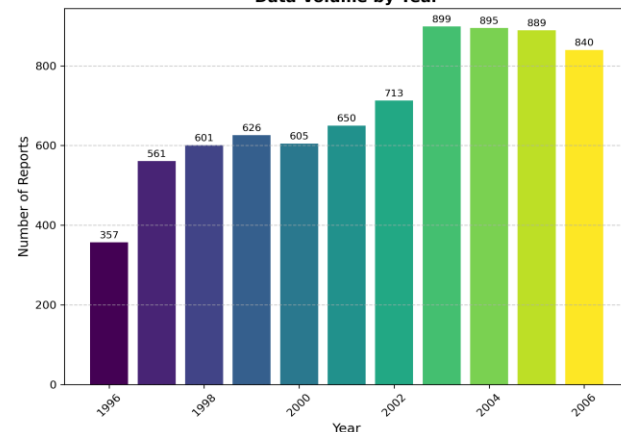
Data Exploration

Data Scale

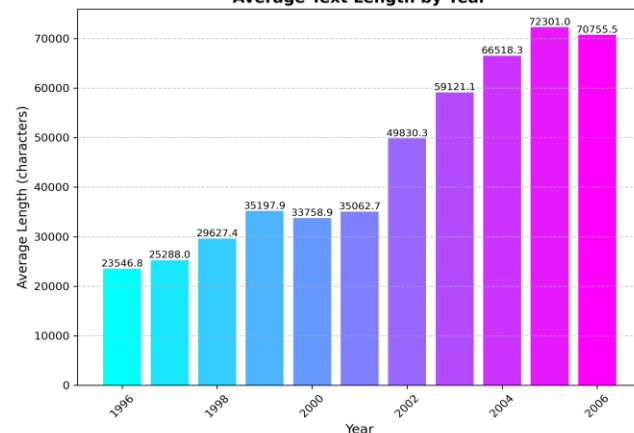
Lexical Diversity

Sentiment Analysis

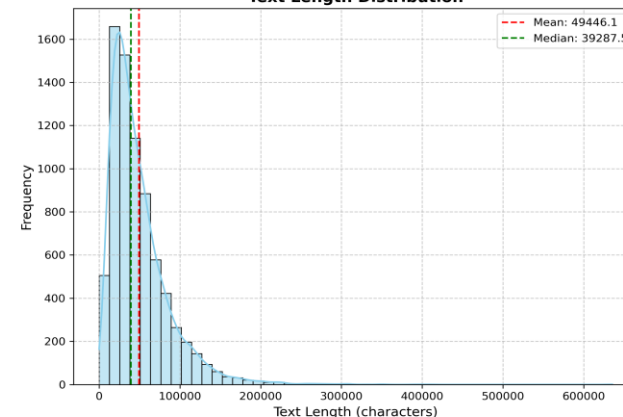
Data Volume by Year



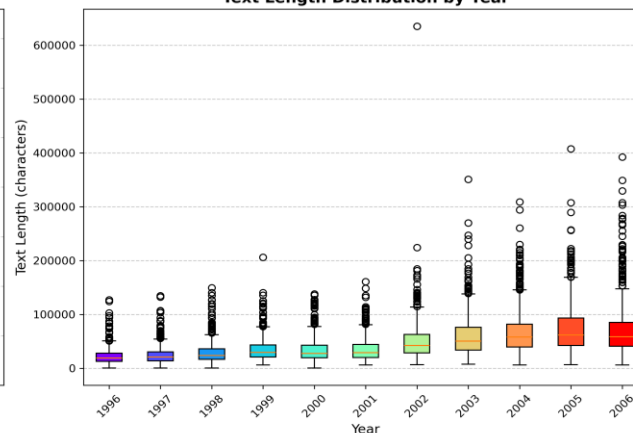
Average Text Length by Year



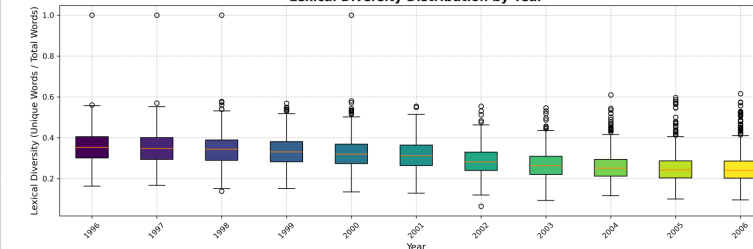
Text Length Distribution



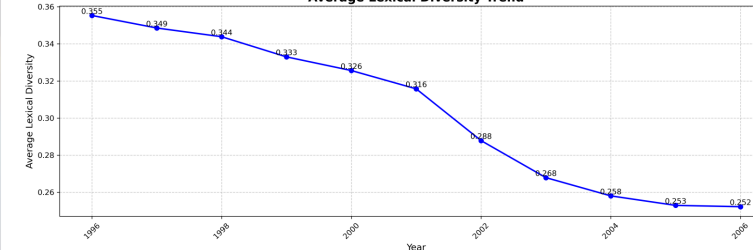
Text Length Distribution by Year



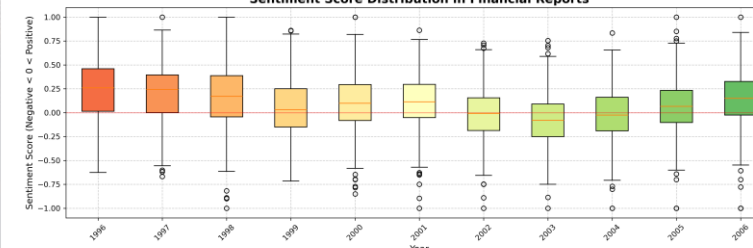
Lexical Diversity Distribution by Year



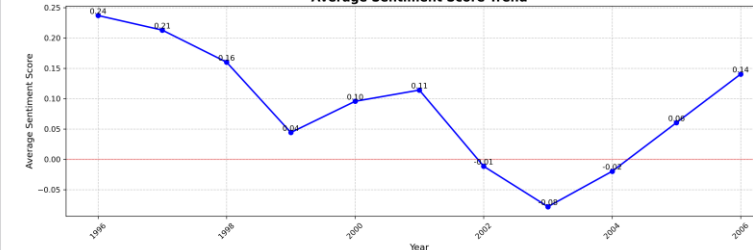
Average Lexical Diversity Trend



Sentiment Score Distribution in Financial Reports

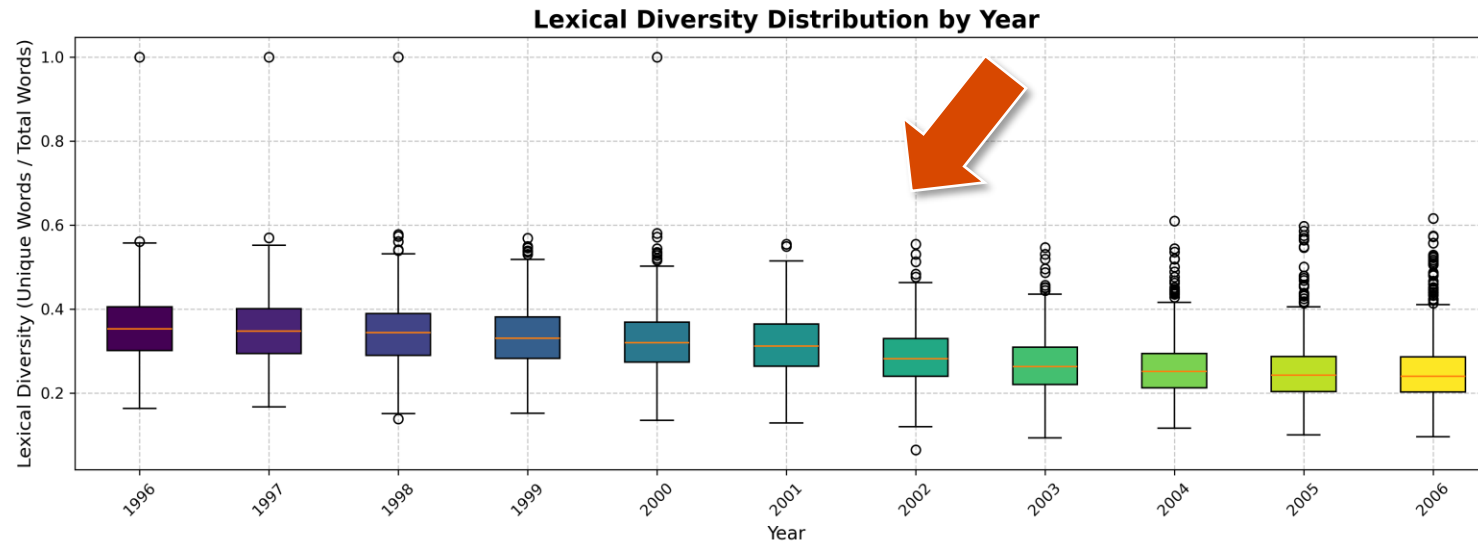


Average Sentiment Score Trend



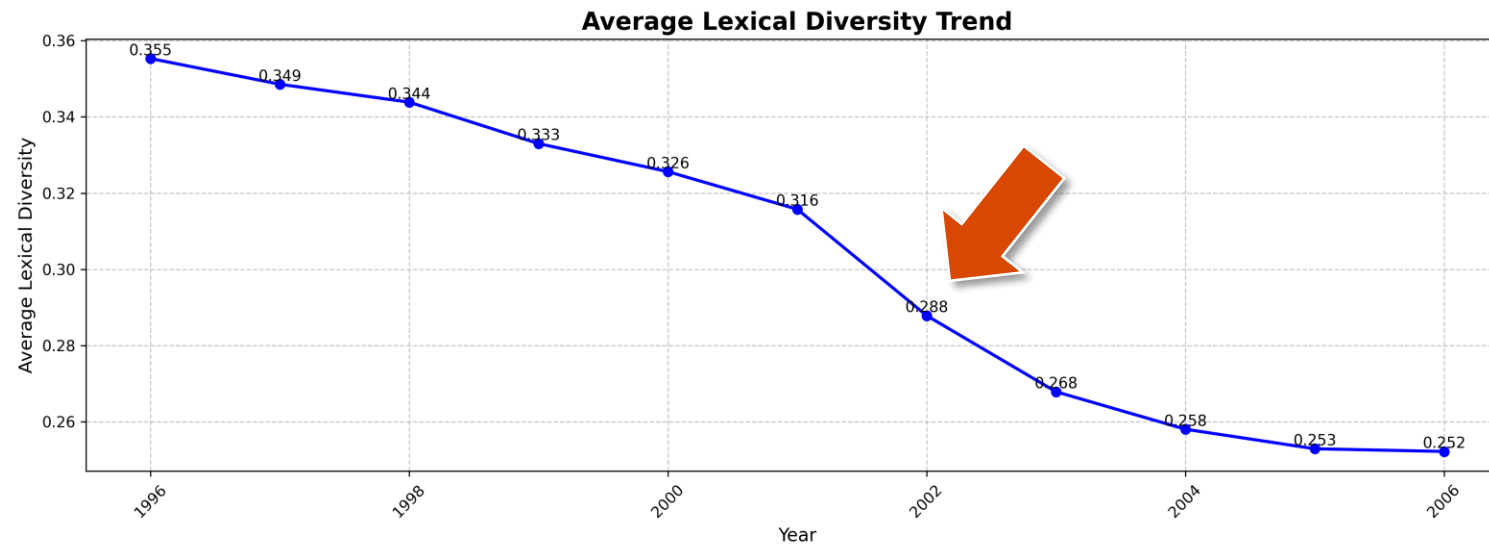
Data Exploration

Lexical Diversity & Sentimental Analysis

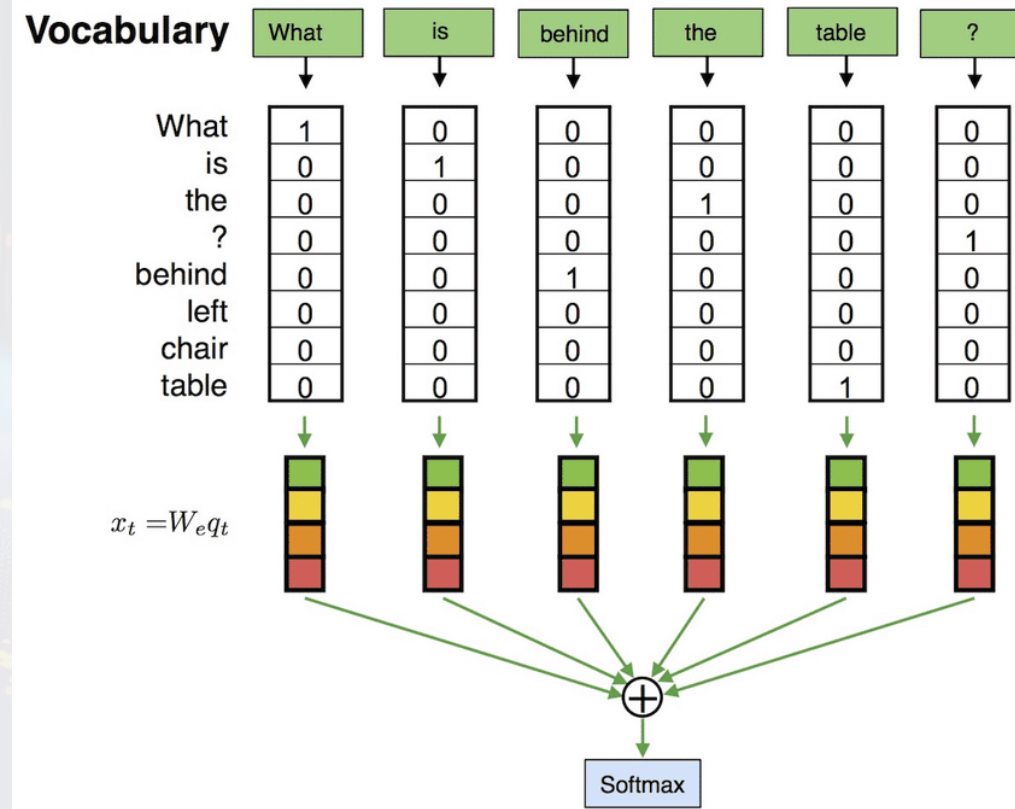
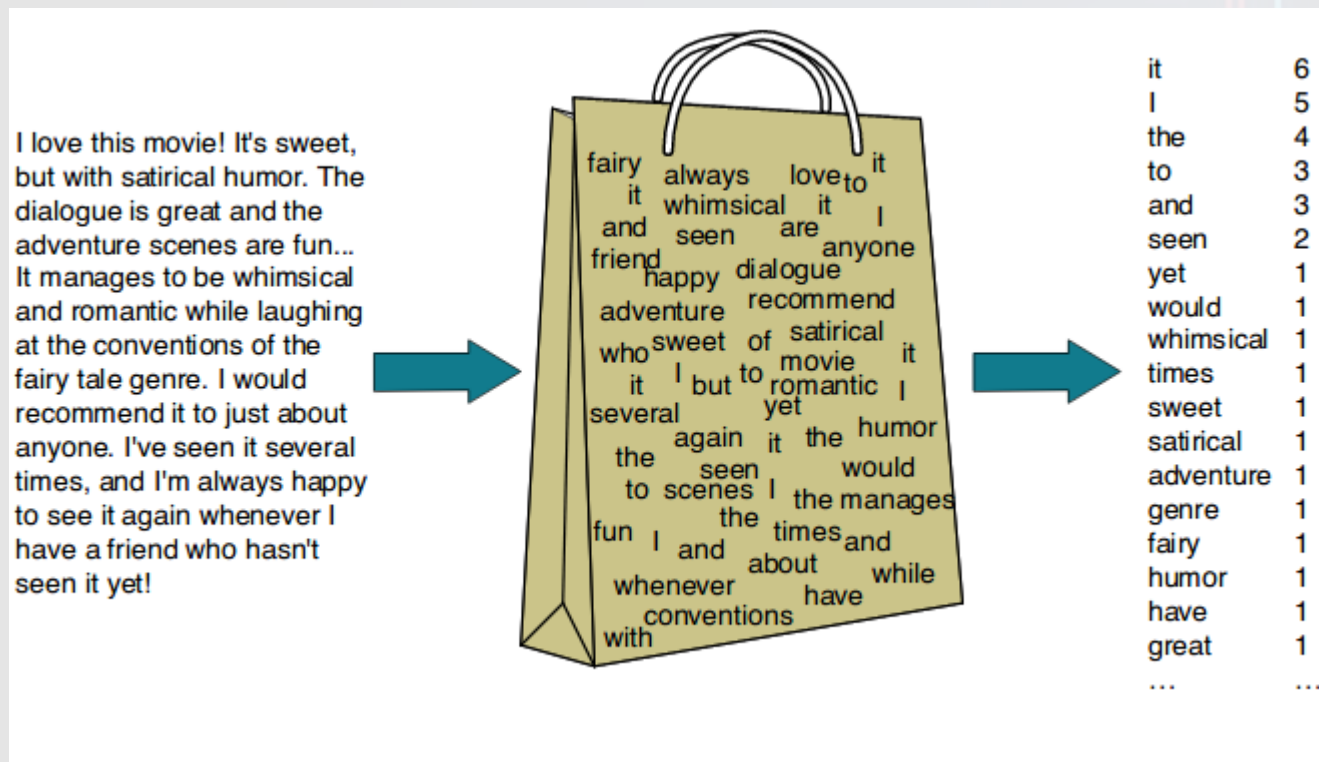


In July 2002, the United States enacted the *Sarbanes–Oxley Act*, imposing stricter disclosure standards for public companies.

Given the exogenous shock nature of the policy intervention, the project decided to divide the experiment into two time intervals, nearly **pre- and post-policy**, respectively



Bag of Words (BOW) Model



[NLP: Bag of words and TF-IDF explained!](#)
by Koushik kumar | Medium

[Bag of Words meets Random Forest](#)
| Kaggle

Feature Engineering – Text Analysis

Statistical Features Extraction

Term Frequency (TF)

- Frequency of words in the current document
- In use, we apply a log transform to mitigate outlier effects

$$h_j(d) = \log(1 + \text{freq}(x_j; d))$$

Bigram - TF

- Extract word frequencies of words and neighbour binary phrases and apply log transform.
- Capture common phrases and neighbour word relationships in text (e.g. 'net loss').

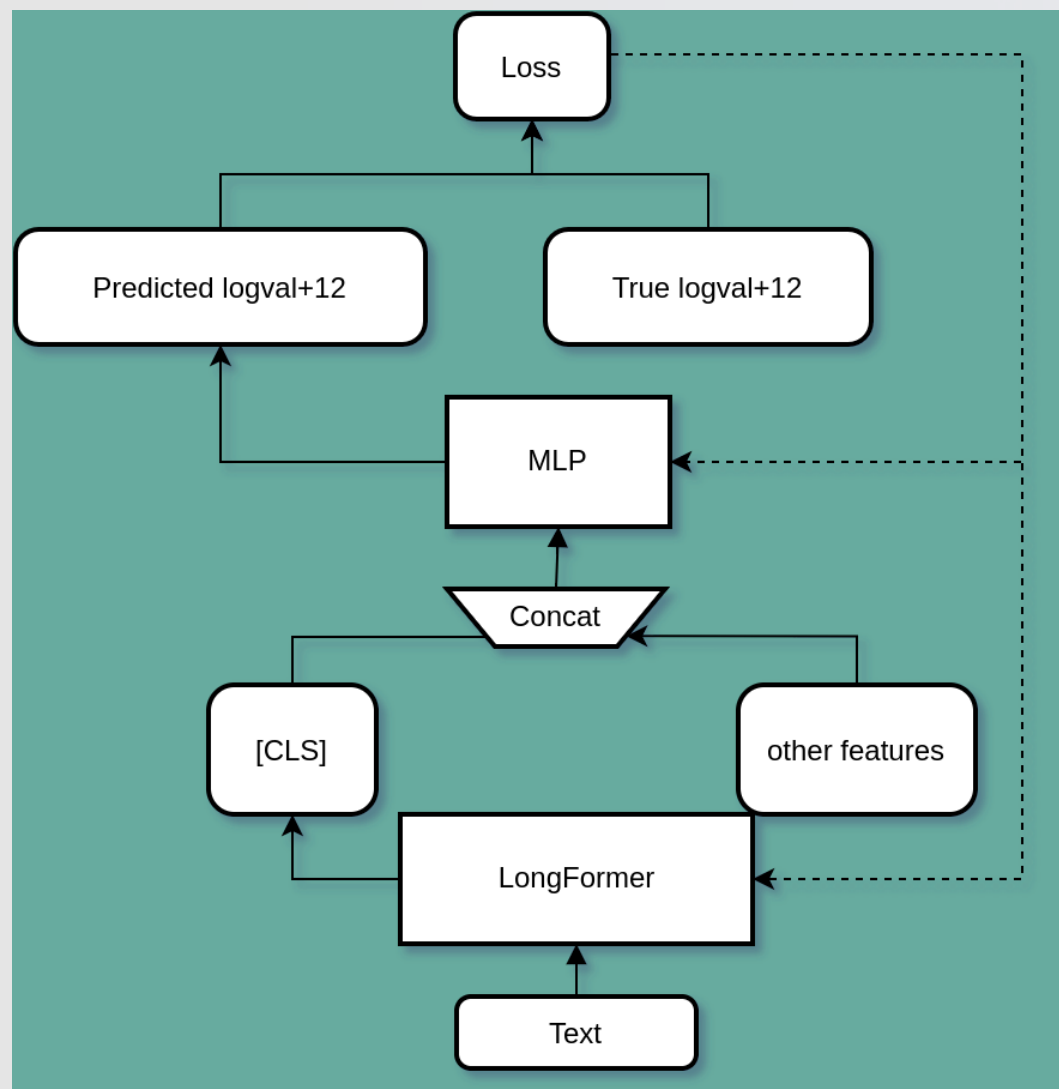
TF - IDF

- Builds on TF by incorporating Inverse Document Frequency (IDF)
- Reduce the weight of common words
- Emphasize rare words

Bigram - TF - IDF

- Builds on Bigram-TF by adding IDF weighting
- Balances phrase frequency with corpus-wide distribution
- Increases discrimination for less frequent phrases

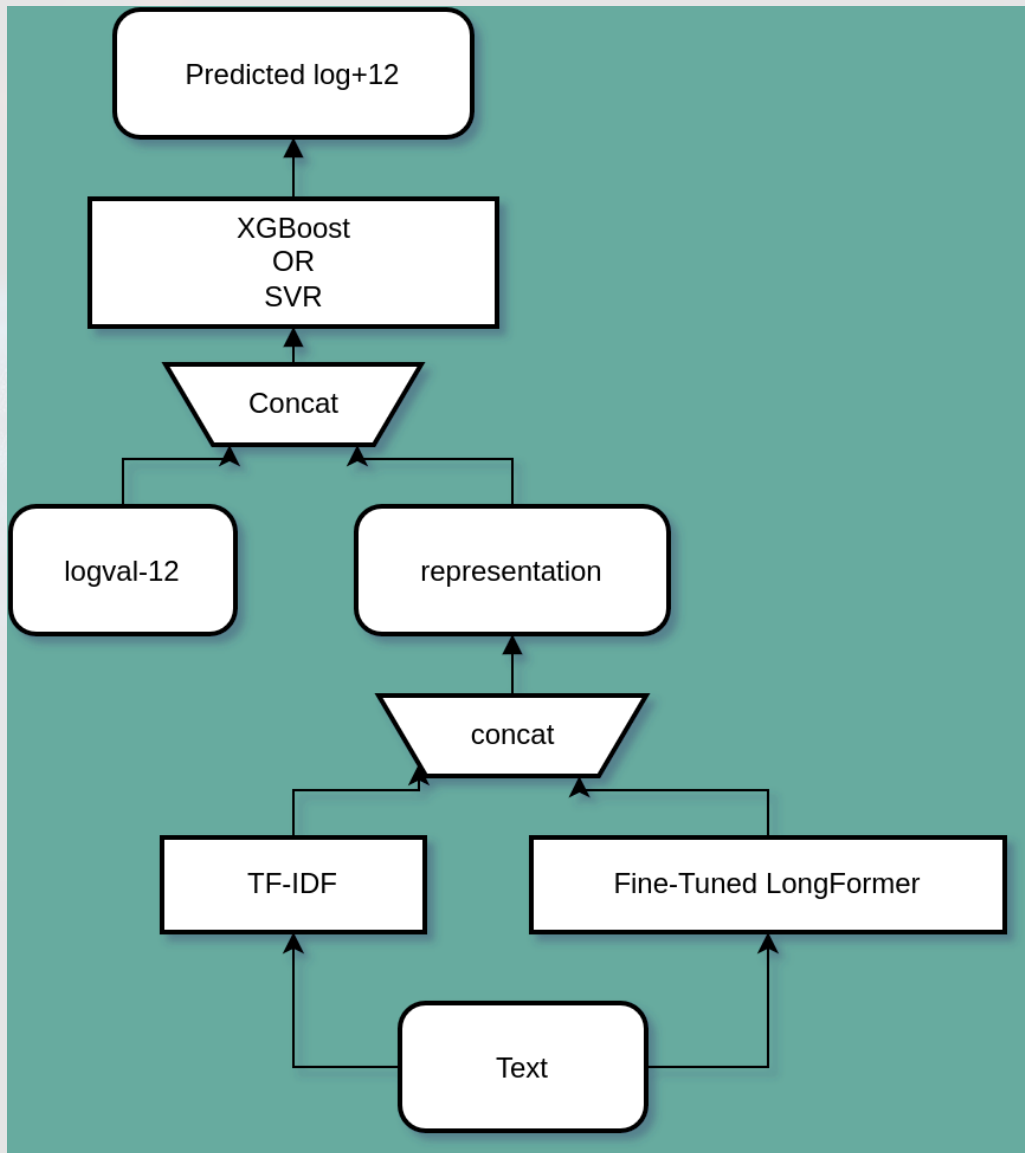
Feature Engineering – Text Analysis



This is a pre-trained model LongFormer fine-tuning flowchart.

It starts with text data, which is processed by the LongFormer model to produce a [CLS] vector. This vector is then concatenated with other useful features (such as historical volatility logval-12), and the combined data input to an MLP (Multi-Layer Perceptron) to generate the Predicted logval+12 (future volatility). The model calculates the MSE Loss by comparing the Predicted logval+12 with the True logval+12, which is used to update parameters in fine-tuning process (both MLP and LongFormer).

Framework & Workflow



In our work framework, we start with cleaned text data, which is processed by TF-IDF and Fine-Tuned LongFormer. The results are concatenated to form a representation.

This representation is then combined with logval-12 (historical volatility), and finally, XGBoost or SVR model is used to predict the logval+12 (future volatility) value.

We have tried many times in this experiment, combined different features and selected different machine learning models. Our aim is to explore the differences in predictive regression accuracy among different features and different models combinations and to analyze the reasons for these differences.

The results will be shown in the following slides, which mark the predictive machine learning model with the highest accuracy and its corresponding features.

Training Experiment Design

Training & Test Dataset

- To predict 2006 vol, using dataset from 2001 to 2005
- Similary, using dataset from 1996-2000 to predict 2001 vol

Basic data clearing

- Remove stopwords, Stemming, and convert text to lowercase for traditional stastistical features
- There is no need to do any preprocess to the text which used for pre-trained model

Xgboost Hyperparameters

max_depth=6
min_child_weight=1
min_loss=0
learning rate=0.1
estimators=100
L2 penalty=1

SVR Hyperparameters

kernel = rbf
error penalty = 1.0
tolerance = 0.001

Before the *Sarbanes–Oxley Act*

feature_type	model_type	MSE	R ²
bigram_tfidf_with_logv12	XGBoost	0.18017763	0.5610444
tf_with_logv12	XGBoost	0.18414356	0.55138245
bigram_tf_with_logv12	XGBoost	0.18419326	0.55126137
longformer_bigram_with_logv12	XGBoost	0.18487247	0.54960667
tfidf_with_logv12	XGBoost	0.18668545	0.5451898
bigram_tfidf_with_logv12	SVR	0.19303066	0.52973137
bigram_tfidf	SVR	0.19695076	0.52018107
bigram_tf_with_logv12	SVR	0.19697474	0.52012265
logv_minus_12_only (baseline)	XGBoost	0.19855868	0.51626381
bigram_tf	SVR	0.20185373	0.50823627
bigram_tfidf	XGBoost	0.20215232	0.50750885
tfidf_with_logv12	SVR	0.20340644	0.5044535
tfidf	XGBoost	0.20359978	0.50398249
longformer_bigram	XGBoost	0.20497933	0.50062157
tf	XGBoost	0.20634332	0.49729856
bigram_tf	XGBoost	0.20679932	0.49618765
tfidf	SVR	0.20871195	0.49152802
tf_with_logv12	SVR	0.21012691	0.48808084
tf	SVR	0.21907049	0.46629214
longformer_only	XGBoost	0.32508698	0.20801076
longformer_only	SVR	0.4138766	-0.0083019
longformer_bigram_with_logv12	SVR	0.50983793	-0.2420865
longformer_bigram	SVR	0.509838	-0.2420867
logv_minus_12_only (baseline)	SVR	2.41380051	-4.8805925

Feature	Importance
Historical vol logv12	80.85540771
estat	24.88868141
argentina	10.36391068
tangibl	9.533727646
quantit qualit	8.494257927
review	8.422605515
net loss	7.985421181
energi	7.944610596
seed	7.848720074
equip million	7.541954994
oper effici	7.292068481
corn	7.071839809
rate primarili	7.056113243
suppli chain	6.94379425
feed	6.826294422
offset effect	6.817842007
fx	6.790940285
tax per	6.622413635
apb	6.479614258
texa	6.402114868

After the *Sarbanes–Oxley Act*

feature_type	model_type	MSE	R ²	Feature	Importance
bigram_tf_with_logv12	XGBoost	0.127489606	0.558582722	Historical vol logv12	119.4632797
longformer_bigram_with_logv12	XGBoost	0.127771108	0.557608055	divert	20.41123962
bigram_tfidf_with_logv12	XGBoost	0.130517287	0.548099743	gener administr	20.03438187
tf_with_logv12	XGBoost	0.133240954	0.538669375	actual futur	15.14949512
tfidf_with_logv12	XGBoost	0.135382626	0.531254098	length	13.87048531
logv_minus_12_only (baseline)	XGBoost	0.143710815	0.502418755	act amend	13.25977898
bigram_tf	XGBoost	0.165544409	0.426822585	initi recognit	12.5438633
longformer_bigram	XGBoost	0.166453559	0.423674765	cash proce	11.82049561
bigram_tfidf_with_logv12	SVR	0.166521462	0.423439658	stock may	10.74275684
bigram_tfidf	SVR	0.170356544	0.410161154	fda approv	10.69298553
bigram_tfidf	XGBoost	0.170797915	0.408632958	product shipment	10.28904819
bigram_tf_with_logv12	SVR	0.173035605	0.400885229	financi account	9.902664185
tf	XGBoost	0.17562322	0.391925926	requir capit	9.301649094
tfidf	XGBoost	0.177506434	0.385405527	impact inflat	9.010601044
bigram_tf	SVR	0.180300338	0.375731973	interest entiti	8.736427307
tfidf_with_logv12	SVR	0.186845658	0.353069599	reit	8.259461403
tfidf	SVR	0.194527815	0.326471062	increas billion	8.059524536
tf_with_logv12	SVR	0.19598403	0.321429093	regularli review	7.951934814
tf	SVR	0.212755473	0.263360011	net loss	7.933226585
longformer_only	XGBoost	0.319589357	-0.106539341	coast	7.848445892
longformer_bigram	SVR	0.322647628	-0.117128234		
longformer_bigram_with_logv12	SVR	0.322647628	-0.117128235		
longformer_only	SVR	0.421106638	-0.4580306		
logv_minus_12_only (baseline)	SVR	2.704373234	-8.363563935		

Conclusion

This project draws the following three core conclusions from a systematic model comparison study:

One, in terms of financial text feature engineering, the hybrid feature that incorporates Word Frequency-Integrated Dual Grammar TF-IDF and traditional market factors (historical volatility) demonstrates superior predictive capability, which is significantly better than semantically dense features based on fine-tuning of pre-trained models.

Second, the integrated learning approach has significant advantages in this study. Specifically, the Xgboost model significantly outperforms the SVR model in terms of accuracy and becomes the preferred model in this study.

Third, financial text analysis has obvious domain specificity. Context agnostic textual statistical features may have an advantage over context based embeddings in terms of interpretability of risk representations. However, during the data preprocessing stage, ignoring the coherence of the entire text and the semantic relevance of the context results in the failure of pre - trained models to achieve significantly better performance than traditional methods in prediction tasks through extracting context deep embeddings. This also shows that when extracting keywords from financial reports, even without focusing on contextual information, a very high accuracy rate can still be achieved.

Interpretability Challenges



Deep Learning Models

Interpretability challenges with deep learning models like BERT.

Understanding the decision- making process of complex models is difficult.

Future Exploration

Future exploration includes integration of interpretability frameworks (e.g., SHAP values).

This can help in explaining the model's predictions and improving trust.

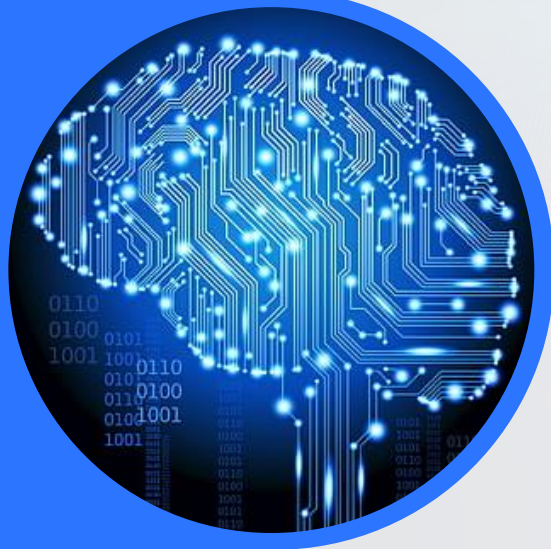
Potential for Expansion

Potential expansion with larger and newer datasets.

Continuous improvement of models with more data and advanced techniques.



References



Ali, Amal Al, et al. "A powerful predicting model for financial statement fraud based on optimized XGBoost ensemble learning technique." *Applied Sciences* 13.4 (2023): 2272.

Rawte, Vipula D., Mohammed J. Zaki, and Aparna Gupta. "FETILDA: An Effective Framework For Fin-tuned Embeddings For Long Financial Text Documents." *arXiv e-prints* (2022): arXiv-2206.

Chang, Ariana, Tian-Shyug Lee, and Hsiu-Mei Lee. "Applying sustainable development goals in financial forecasting using machine learning techniques." *Corporate social responsibility and environmental management* 31.3 (2024): 2277-2289.

Oukhouya, Hassan, and Khalid El Himdi. "Comparing machine learning methods—svr, xgboost, lstm, and mlp—for forecasting the moroccan stock market." *Computer Sciences & Mathematics Forum*. Vol. 7. No. 1. MDPI, 2023.

Kogan, S., Levin, D., Routledge, B. R., Sagi, J. S., & Smith, N. A. (2009). Predicting Risk from Financial Reports with Regression.



Thank you

