

## Our Project Overview

44

01

### Main goal of our project

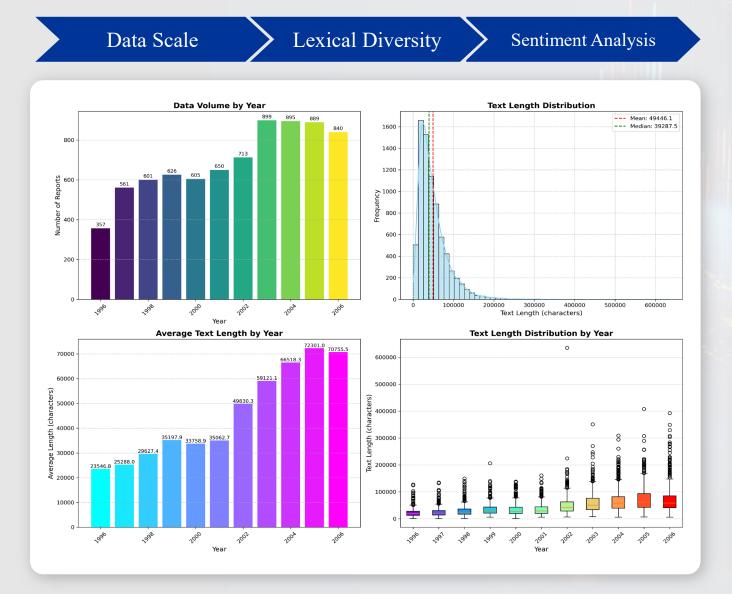
Taking corporate annual financial reports as the object of analysis, this study focuses on exploring the predictive ability of textual features on future stock price volatility. Based on empirical tests of machine learning regression models, we systematically evaluate the differences in the predictive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

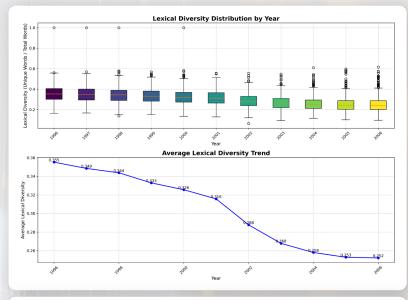
02

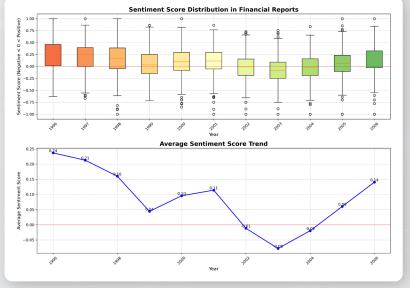
#### Conclusion

On the dataset studied in this project, the ensemble learning algorithm XGboost has the highest prediction accuracy, the effectiveness of textual information in interpreting volatility exhibits feature type heterogeneity, and the prediction accuracy of volatility based on keyword sparse features is higher than that based on semantically embedding from pre-trained models. This may imply that surface lexical patterns may be more powerful risk indicators than deep semantic representations in financial text analysis scenarios.

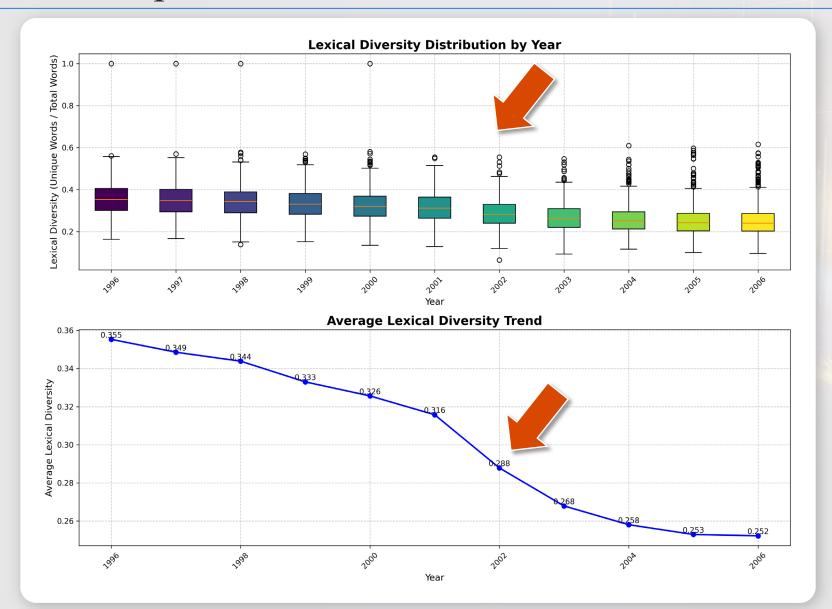
# Data Exploration







# **Data Exploration**

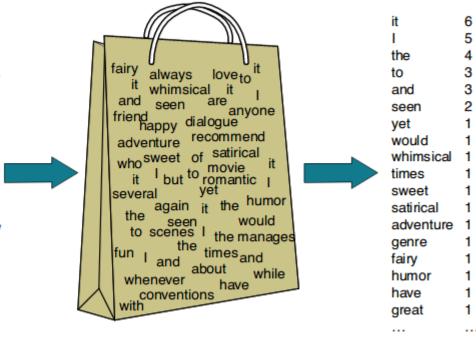


In July 2002, the United States enacted the *Sarbanes–Oxley Act*, imposing stricter disclosure standards for public companies.

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

I love this movie! It's sweet. but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



Vocabulary What table behind What is the behind left 0 chair 0 0 0 table 0 0  $x_t = W_e q_t$ Softmax

NLP: Bag of words and TF-IDF explained! by Koushik kumar | Medium



늘 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 applies a <u>log transform</u> to mitigate outlier effects

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

#### 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').

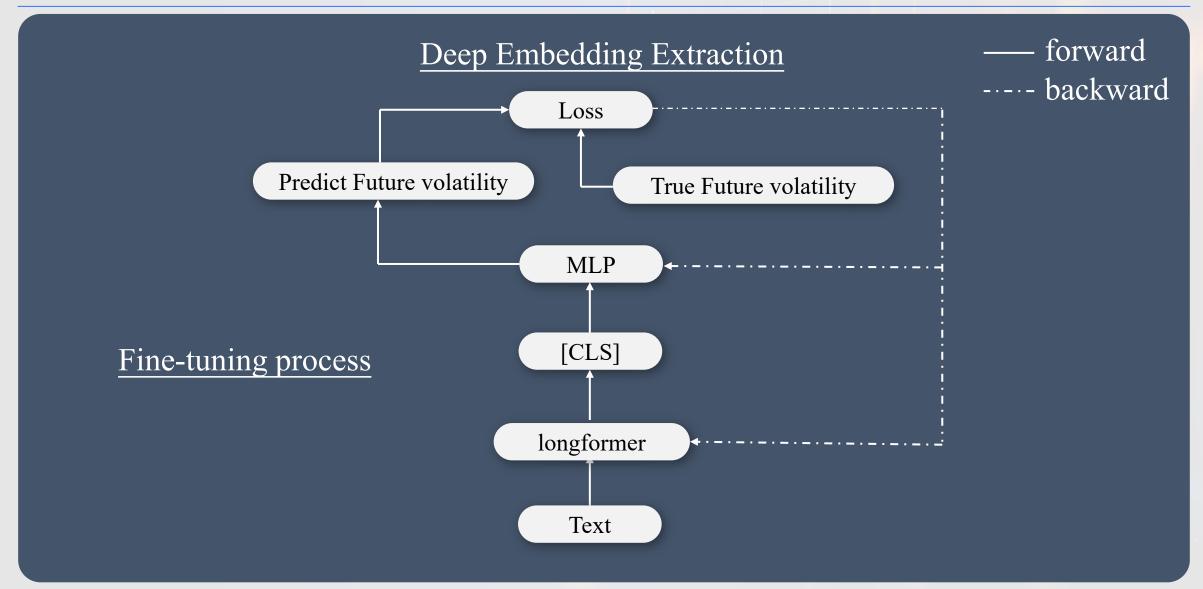
#### TF - IDF

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

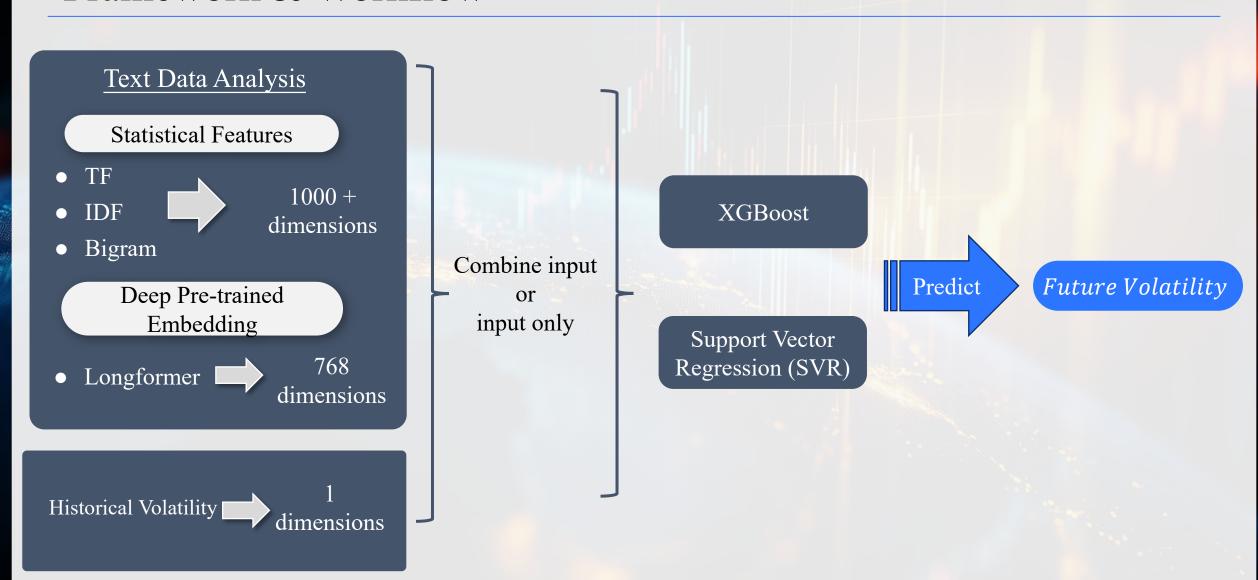
#### Bigram - TF - IDF

- 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



### Framework & Workflow



# 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

# Xgboost Hyperparameters

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

### Basic data clearing

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

### SVR Hyperparameters

kernel = rbf error penalty = 1.0 tolerance = 0.001

### Before the Sarbanes-Oxley Act

feature_type	model_type	MSE	$\mathbb{R}^2$		
bigram_tfidf_with_logv12	XGBoost	0.18017763	0.5610444	Feature	Importance
tf_with_logv12	XGBoost	0.18414356	0.55138245	Historical vol logv12	80.85540771
bigram_tf_with_logv12	XGBoost	0.18419326	0.55126137	estat	24.88868141
longformer_bigram_with_logv12	XGBoost	0.18487247	0.54960667	argentina	10.36391068
tfidf_with_logv12	XGBoost	0.18668545	0.5451898	tangibl	9.533727646
bigram_tfidf_with_logv12	SVR	0.19303066	0.52973137	quantit qualit	8.494257927
bigram_tfidf	SVR	0.19695076	0.52018107	review	8.422605515
bigram_tf_with_logv12	SVR	0.19697474	0.52012265	net loss	7.985421181
	XGBoost	0.10055060	0.51626381	energi	7.944610596
logv_minus_12_only (baseline)		0.19855868		seed	7.848720074
bigram_tf bigram_tfidf	SVR XGBoost	0.20185373 0.20215232	0.50823627 0.50750885	equip million	7.541954994
tfidf_with_logv12	SVR	0.20215252	0.5044535	8	
tfidf	XGBoost	0.20359978	0.50398249	oper effici	7.292068481
				corn	7.071839809
longformer_bigram	XGBoost	0.20497933	0.50062157	rate primarili	7.056113243
tf	XGBoost	0.20634332	0.49729856	-	
bigram_tf	XGBoost	0.20679932	0.49618765	suppli chain	6.94379425
tfidf	SVR	0.20871195	0.49152802	feed	6.826294422
tf_with_logv12	SVR	0.21012691	0.48808084	offset effect	6.817842007
tf	SVR	0.21907049	0.46629214	fx	6.790940285
longformer_only	XGBoost	0.32508698	0.20801076	IX	
longformer_only	SVR	0.4138766	-0.0083019	tax per	6.622413635
longformer_bigram_with_logv12	SVR	0.50983793	-0.2420865	apb	6.479614258
longformer_bigram	SVR	0.509838	-0.2420867	texa	6.402114868
logv_minus_12_only (baseline)	SVR	2.41380051	-4.8805925		

### After the Sarbanes–Oxley Act

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

### 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 + Bigram) 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. Surface-level textual statistical features may have an advantage over deep semantic embeddings in terms of interpretability of risk representations. This finding provides a new possibility hypothesis for textual analysis in the financial quantitative domain, which has important research value and practical significance.

# 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.

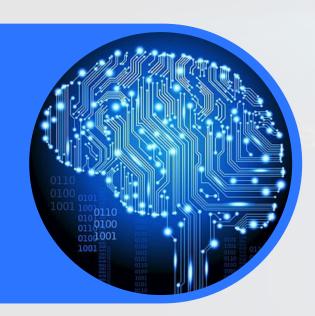
Continuous improvement of models with more data and advanced techniques.





## 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.

