

```
In [53]: # %% Import Required Libraries
import pandas as pd
import duckdb
import nltk
import random
import numpy as np
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from xgboost import XGBRegressor
import pandas_market_calendars as mcal
from sklearn.ensemble import RandomForestRegressor
from sklearn.linear_model import LinearRegression
from sklearn.neural_network import MLPRegressor
from sklearn.metrics import mean_absolute_error, r2_score, accuracy_score
from sklearn.inspection import permutation_importance
from tabulate import tabulate

# %% Download NLTK Resources
nltk.download('stopwords')
nltk.download('punkt')

# %% Database Connection
db_path = '/Users/bradams/Documents/financial_news.db'
conn = duckdb.connect(database=db_path, read_only=False)

# %% Load Stopwords
stop_words = set(stopwords.words('english'))

# %% Get All Available Tickers & Sample a Subset
query_tickers = "SELECT DISTINCT ticker FROM Headlines.Articles_Trading_Day"
ticker_df = conn.execute(query_tickers).fetchhdf()

# %% Filter Tickers by Minimum Article Count Threshold
MIN_ARTICLE_THRESHOLD = 100

# Count articles per ticker
article_counts = conn.execute("""
    SELECT ticker, COUNT(*) as article_count
    FROM Headlines.Articles_Trading_Day
    GROUP BY ticker
""").fetchhdf()

# Filter tickers with enough article volume
eligible_tickers = article_counts[article_counts["article_count"] >= MIN_ARTICLE_THRESHOLD]

# Randomly sample tickers for testing
ticker_all = random.sample(eligible_tickers, min(len(eligible_tickers), 50))

# print(f"Selected {len(ticker_all)} tickers with at least {MIN_ARTICLE_THRESHOLD} articles")

# %% Load Loughran-McDonald Dictionary for High-Risk Words
lm_dict_path = 'Loughran-McDonald_MasterDictionary_1993-2023.csv'
lm_dict = pd.read_csv(lm_dict_path)
high_risk_words = set(lm_dict[lm_dict["Negative"] > 0]["Word"].str.lower())
```

```

# %% Load NYSE Trading Calendar
nyse = mcal.get_calendar('NYSE')

# %% Query Sentiment Analysis Data
sentiment_analysis = conn.execute("SELECT * FROM Headlines.finbert_analysis").fetchall()

# %% Define Helper Functions
def next_trading_day(date):
    """Get the next valid trading day."""
    date = pd.Timestamp(date)
    while len(nyse.valid_days(start_date=date.strftime('%Y-%m-%d'), end_date=date.strftime('%Y-%m-%d'))) < 1:
        date += pd.Timedelta(days=1)
    return date

def tokenize_text(text):
    """Tokenize text and count high-risk words."""
    words = word_tokenize(text.lower())
    words = [word for word in words if word.isalpha() and len(word) > 2 and word not in high_risk_words]
    high_risk_count = sum(1 for word in words if word in high_risk_words)
    return words, high_risk_count

def clean_data(X, y):
    """Remove rows where y is NaN."""
    mask = ~np.isnan(y) & ~np.isnan(X).any(axis=1)
    return X[mask], y[mask]

# **Feature Importance for Tree-Based Models**
def evaluate_feature_importance(models, X_train, X_test, y_test):
    """Computes and prints feature importance for different models."""
    feature_names = [
        "token_score", "high_risk_count", "finbert_title_score", "finbert_descripti
        "finbert_title_positive", "finbert_title_neutral", "finbert_title_negative"
        "finbert_description_positive", "finbert_description_neutral", "finbert_des
    ]

    tree_models = ["Random Forest", "XGBoost"]
    feature_importances = {name: models[name].feature_importances_ for name in tree_models}

    if feature_importances:
        feature_importance_df = pd.DataFrame(feature_importances, index=feature_names)
        print("\n**Tree-Based Model Feature Importance**")
        print(tabulate(feature_importance_df, headers="keys", tablefmt="psql"))

# %% Function to Process a Single Ticker
def process_ticker(ticker, conn, nyse, sentiment_df, result_df):
    """Processes a single stock ticker and trains models using leakage-free token scores"""
    query = f"""
    SELECT
        a.ticker,
        a.mapped_trading_date AS publish_date,
        a.article_title,
        a.description,
        dpm.price_change_percentage,
        f.finbert_title_score,
        f.finbert_description_score,
    """

```

```

        f.finbert_title_positive,
        f.finbert_title_neutral,
        f.finbert_title_negative,
        f.finbert_description_positive,
        f.finbert_description_neutral,
        f.finbert_description_negative
    FROM "Headlines"."Articles_Trading_Day" a
    INNER JOIN "Headlines"."Daily_Price_Movement" dpm
        ON a.mapped_trading_date = dpm.trading_date
        AND a.ticker = dpm.ticker
    INNER JOIN "Headlines"."finbert_analysis" f
        ON a.guid = f.guid
    WHERE a.ticker = ?
    """
news_df = conn.execute(query, [ticker]).fetchdf()

if news_df.empty:
    return None, None, None, None

# Preprocess
news_df["publish_date"] = pd.to_datetime(news_df["publish_date"]).dt.date
news_df["description"] = news_df["description"].fillna("")
news_df["adjusted_date"] = news_df["publish_date"].apply(next_trading_day)
news_df["tokenized_words"], news_df["high_risk_count"] = zip(*news_df["description"].str.split())
article_count = len(news_df)
news_df["article_count"] = article_count

# Add Labels
y = news_df["price_change_percentage"].values

# Split first
split_index = int(len(news_df) * 0.7)
train_df = news_df.iloc[:split_index].copy()
test_df = news_df.iloc[split_index: ].copy()
y_train = y[:split_index]
y_test = y[split_index:]

if len(train_df) < 2 or len(test_df) < 1:
    print(f"Skipping {ticker}: not enough data.")
    return None, None, None, None

# Compute token scores **from training only**
word_scores = {}
for _, row in train_df.iterrows():
    words = row["tokenized_words"]
    price_change = row["price_change_percentage"]
    total_words = len(words)
    if total_words > 0:
        word_counts = {word: words.count(word) / total_words for word in words}
        for word, ratio in word_counts.items():
            word_scores.setdefault(word, []).append(ratio * price_change)

token_scores_dict = {word: np.mean(scores) for word, scores in word_scores.items}

# Apply to both train and test

```

```

for df in [train_df, test_df]:
    df.loc[:, "token_score"] = df["tokenized_words"].apply(
        lambda tokens: sum(token_scores_dict.get(token, 0) for token in tokens)
    )

# Combine back and prepare features
full_df = pd.concat([train_df, test_df])
feature_columns = [
    "token_score", "high_risk_count", "finbert_title_score", "finbert_descripti
    "finbert_title_positive", "finbert_title_neutral", "finbert_title_negative"
    "finbert_description_positive", "finbert_description_neutral", "finbert_des
]

X = full_df[feature_columns].values
y = full_df["price_change_percentage"].values
X, y = clean_data(X, y)

if len(X) < 2:
    return None, None, None, None

# Recreate split after cleaning
split_index = int(len(y) * 0.7)
X_train, X_test = X[:split_index], X[split_index:]
y_train, y_test = y[:split_index], y[split_index:]

if len(X_train) == 0 or len(X_test) == 0:
    return None, None, None, None

# Normalize
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

# Train models
models = {
    "XGBoost": XGBRegressor(objective="reg:squarederror", n_estimators=100, lea
    "Random Forest": RandomForestRegressor(n_estimators=100, random_state=42),
    "Linear Regression": LinearRegression(),
    "Neural Network": MLPRegressor(hidden_layer_sizes=(50, 50), max_iter=1500)
}
for name, model in models.items():
    model.fit(X_train_scaled, y_train)

return models, y_test, X_test_scaled, full_df

# %% Execute Processing for Selected Tickers
columns = ["symbol", "model", "MAE", "R²"]
result_df = pd.DataFrame(columns=columns)

for col in result_df.columns:
    if "MAE" in col:
        print(col, (result_df[col] < 0).sum(), "negative values found")

# %% Execute Processing for Selected Tickers

```

```

trained_models = {}
for ticker in ticker_all:
    models, y_test, X_test_scaled, dataset = process_ticker(ticker, conn, nyse, sen
        if models is not None:
            trained_models[ticker] = {"models": models, "y_test": y_test, "X_test": X_t

# %% Compare Model Results
def format_model_results(trained_models):
    """Creates a formatted table comparing models across all tickers."""
    results = []

    for ticker, data in trained_models.items():
        models, y_test_actual, X_test_actual = data["models"], data["y_test"], data
        predictions = {name: model.predict(X_test_actual) for name, model in models

            row = [ticker, len(data["dataset"])]
            for name, y_pred in predictions.items():
                mae = mean_absolute_error(y_test_actual, y_pred)
                # Check if y_test_actual has at least two samples for R^2 calculation
                if len(y_test_actual) > 1:
                    r2 = r2_score(y_test_actual, y_pred)
                else:
                    r2 = None # Set R^2 to None if not enough samples
                row.extend([mae, r2])
            results.append(row)

    column_headers = ["Stock Symbol", "Article Count"] + [f"{model} MAE" for model
    results_df = pd.DataFrame(results, columns=column_headers)

    print("\n📊 **Final Model Comparison Table**")
    print(tabulate(results_df, headers="keys", tablefmt="psql"))

    return results_df

# ✅ Run
formatted_results_df = format_model_results(trained_models)

```

```

[nltk_data] Downloading package stopwords to
[nltk_data]      /Users/bradams/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to /Users/bradams/nltk_data...
[nltk_data] Package punkt is already up-to-date!
MAE 0 negative values found

```

```
/Users/bradams/Documents/OMSAPracticum-1/.venv/lib/python3.12/site-packages/sklearn/  
neural_network/_multilayer_perceptron.py:691: ConvergenceWarning: Stochastic Optimiz  
er: Maximum iterations (1500) reached and the optimization hasn't converged yet.  
    warnings.warn(  
/Users/bradams/Documents/OMSAPracticum-1/.venv/lib/python3.12/site-packages/sklearn/  
neural_network/_multilayer_perceptron.py:691: ConvergenceWarning: Stochastic Optimiz  
er: Maximum iterations (1500) reached and the optimization hasn't converged yet.  
    warnings.warn(  
/Users/bradams/Documents/OMSAPracticum-1/.venv/lib/python3.12/site-packages/sklearn/  
neural_network/_multilayer_perceptron.py:691: ConvergenceWarning: Stochastic Optimiz  
er: Maximum iterations (1500) reached and the optimization hasn't converged yet.  
    warnings.warn(  
/Users/bradams/Documents/OMSAPracticum-1/.venv/lib/python3.12/site-packages/sklearn/  
neural_network/_multilayer_perceptron.py:691: ConvergenceWarning: Stochastic Optimiz  
er: Maximum iterations (1500) reached and the optimization hasn't converged yet.  
    warnings.warn(  
/Users/bradams/Documents/OMSAPracticum-1/.venv/lib/python3.12/site-packages/sklearn/  
neural_network/_multilayer_perceptron.py:691: ConvergenceWarning: Stochastic Optimiz  
er: Maximum iterations (1500) reached and the optimization hasn't converged yet.
```

## \*\*Final Model Comparison Table\*\*

	Stock Symbol	Article Count	XGBoost MAE	Random Forest MAE	Li
near Regression MAE	Neural Network MAE	XGBoost R <sup>2</sup>	Random Forest R <sup>2</sup>	L	
near Regression R <sup>2</sup>	Neural Network R <sup>2</sup>				
0   ADI	0.509795	612   0.426239	0.108245	0.573053	
0.512082	-0.102414				
0.529151   EA	0.436612	438   0.596134	0.401815	0.374079	
0.579564	0.322068				
0.658362   MAA		235   1.06418	-0.828672	-0.999356	
0.952941	-0.63456	1.02749			
1.05242	-0.871964				
1   KEX		241   1.60584	-0.522446	-0.889356	
1.56666	-0.61972	1.55305			
1.49029	-0.468185				
4   VST		780   3.09272	nan	nan	
3.3384	nan	15.7827			
59.5456	nan				
5   LBTYK		265   0.375743	-0.849344	-0.276375	
0.382916	-0.400332	0.473048			
0.479725	-0.587921				
6   HAE		556   1.98834	-0.250984	-0.298846	
1.83478	-0.205406	1.82105			
1.99275	-0.389346				
7   ODFL		280   1.74137	0.00306913	0.0373614	
1.75137	0.00949034	1.59815			
1.96195	-0.192234				
8   JHG		302   1.1545	0.0124673	-0.129053	
1.13163	-0.0438064	1.10151			
1.12751	-0.0076528				
9   AMP		379   0.69738	-0.161098	-0.197957	
0.669136	-0.111578	0.704055			
0.71394	-0.136238				
10   HAL		555   1.12945	-0.0359194	-0.0403342	
1.14723	-0.0733557	1.11849			
1.19893	-0.0977195				
11   WDAY		855   0.411665	0.473604	0.534093	
0.43062	0.532815	0.547702			
0.47688	0.565529				
12   ARW		140   1.74531	-1.4917	-1.15939	
1.72809	-1.00573	1.90949			
1.96422	-1.4092				
13   ZS		754   1.33908	-0.362261	-0.209767	
1.38752	-0.245044	1.62315			
1.49656	-0.220181				
14   ES		190   0.949167	-0.316943	-0.454872	
0.916025	-0.359134	0.992827			
1.038	-0.434547				
15   AME		339   0.550038	-0.388116		

0.587946	-0.915768	0.577428	-0.344218	
0.544465	-0.279572			
16   APO	813	0.606966	0.472286	
0.632414	0.408709	0.674691	0.301996	
0.593109	0.537793			
17   CW	246	0.955949	-0.401066	
0.930517	-0.228209	0.906924	-0.144159	
0.937933	-0.308082			
18   FICO	360	1.53902	-0.825539	
1.38523	-0.488276	1.25628	-0.281057	
1.4733	-0.651092			
19   TECH	157	1.56291	-0.202144	
1.62321	-0.205081	1.60475	-0.11841	
1.66209	-0.282433			
20   DHI	896	1.24954	0.151273	
1.30525	0.180289	1.45875	0.0937363	
1.37	0.040143			
21   DRI	403	0.766413	0.145581	
0.898405	0.061769	0.92969	0.111054	
0.786715	0.167486			
22   ETR	307	1.06465	-1.15748	
0.881094	-0.51233	0.982131	-0.705912	
1.00375	-0.839249			
23   ORI	156	0.680784	-0.19635	
0.710171	-0.171073	0.62569	0.203514	
1.15722	-2.02999			
24   CEG	310	0.536668	-0.237166	
0.427013	-0.0710422	167.328	-276123	2
55.32	-750283			
25   CL	957	0.394878	-0.105192	
0.368371	-0.0135481	0.430583	-0.173544	
0.423366	-0.163673			
26   HRL	419	1.17825	-0.0791286	
1.25331	-0.225727	1.22762	-0.00996526	
1.14297	0.00453061			
27   VOYA	191	0.897184	0.347025	
0.841283	0.37147	2.20538	-2.04899	
1.14295	-0.0233507			
28   SRE	397	0.79646	-0.0139861	
0.799912	-0.00362862	0.797205	0.0384754	
0.783769	0.0701356			
29   LHX	562	0.451327	-0.648968	
0.457107	-0.609936	0.555702	-1.45307	
0.51891	-1.16128			
30   ACM	285	0.917759	-0.0119783	
0.935967	0.0109065	0.980343	0.00907276	
0.956364	0.0253966			
31   PINS	880	3.11733	-3.04311	
2.50543	-0.539691	2.64543	-1.1063	
2.67598	-1.05617			
32   BG	197	0.856434	-0.318732	
0.86777	-0.230518	0.817458	-0.00197209	
0.988719	-0.26801			
33   RTX	890	0.572097	0.0485809	
0.585996	0.0873941	0.573071	0.151019	
0.594728	-0.0901373			

## token\_word\_score\_w\_sentiment

34   PHM		560	1.23788	0.0062598	
1.28746	-0.0400312	1.34422	-0.158626		
1.32852	-0.088206				
35   ALB		273	2.05685	0.0960522	
2.10589	0.119011	1.93465	0.243482		
2.08258	0.154363				
36   NOG		166	1.29343	-0.237706	
1.42205	-0.438973	1.8402	-1.75967		
13.0098	-1144.93				
37   MNST		614	0.843924	-0.165277	
0.860964	-0.177352	0.903527	-0.194195		
0.846814	-0.138635				
38   NCNO		135	1.40827	0.23794	
1.4283	0.267735	1.43627	0.250814		
1.49738	0.128145				
39   MASI		322	1.36174	-0.510607	
1.32594	-0.390377	1.25683	-0.249072		
1.20318	-0.128217				
40   FAST		319	0.602634	-0.0422847	
0.60556	-0.0101697	0.643429	-0.0490067		
0.636493	-0.0960206				
41   EQIX		458	0.906209	0.398215	
0.909589	0.432054	0.96866	0.423817		
0.907644	0.380686				
42   NXPI		374	0.922088	0.387021	
0.884011	0.403051	1.26921	0.257397		
1.09377	0.317367				
43   FWONA		451	1.15325	-0.457269	
1.10812	-0.270448	1.09918	-0.134492		
1.03259	-0.0439177				
44   OLED		233	2.11561	-0.278867	
2.07509	-0.144681	2.10915	0.0814535		
2.17283	-0.301129				
45   PSX		533	0.844234	-0.333978	
0.839627	-0.256696	0.845144	-0.26915		
0.844689	-0.295921				
46   MKTX		156	1.5016	-0.357175	
1.49275	-0.276925	1.37537	-0.1062		
3.88808	-20.0705				
47   ALL		621	0.783086	-0.255291	
0.787917	-0.240022	0.775585	-0.077902		
0.729116	0.00130756				
48   WELL		268	0.0174096	0.974338	
0.0525727	0.91068	822.41	-5.93153e+07		1
41.138	-1.74584e+06				
49   NFG		170	0.504715	0.039591	
0.506273	0.179208	0.550106	0.172034		
0.52399	0.13501				

In [54]: `import pandas as pd`

```
models = ["XGBoost", "Random Forest", "Linear Regression", "Neural Network"]
summary_data = []
```

```

for model in models:
    mae_col = f"{model} MAE"
    r2_col = f"{model} R²"
    # Drop NaN rows
    filtered = formatted_results_df[[mae_col, r2_col]].dropna()
    summary_data.append({
        "Model": model,
        "Median MAE": filtered[mae_col].median(),
        "Median R²": filtered[r2_col].median()
    })

summary_df = pd.DataFrame(summary_data)
print(summary_df)

```

	Model	Median MAE	Median R <sup>2</sup>
0	XGBoost	0.935628	1.063338
1	Random Forest	-0.196350	-0.118410
2	Linear Regression	0.912807	1.073092
3	Neural Network	-0.144681	-0.138635

In [55]:

```

import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Let's say 'results_df' is your results table.
# First, transform the wide table into a long format for easier grouping.
models = ["XGBoost", "Random Forest", "Linear Regression", "Neural Network"]
summary_data = []

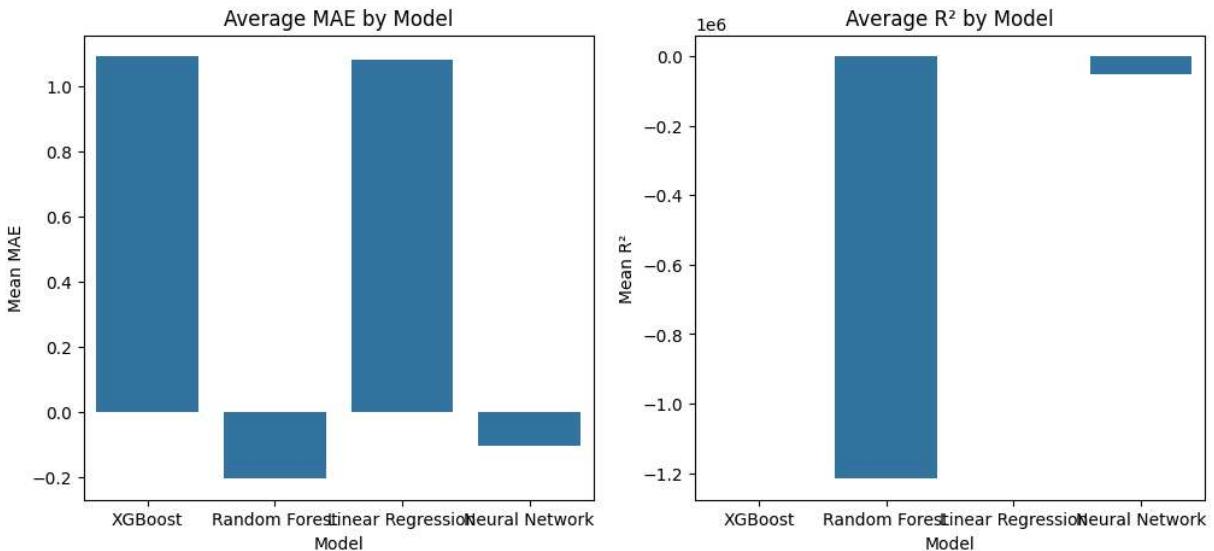
for model in models:
    mae_col = f"{model} MAE"
    r2_col = f"{model} R²"
    # Remove any rows with missing or non-numeric values (optional)
    filtered = formatted_results_df[[mae_col, r2_col]].dropna()
    summary_data.append({
        "Model": model,
        "Mean MAE": filtered[mae_col].mean(),
        "Median MAE": filtered[mae_col].median(),
        "Std MAE": filtered[mae_col].std(),
        "Mean R²": filtered[r2_col].mean(),
        "Median R²": filtered[r2_col].median(),
        "Std R²": filtered[r2_col].std()
    })

summary_df = pd.DataFrame(summary_data)
print(summary_df)

# Visualize average MAE and R² for each model
fig, axes = plt.subplots(1, 2, figsize=(12, 5))
sns.barplot(data=summary_df, x="Model", y="Mean MAE", ax=axes[0])
axes[0].set_title("Average MAE by Model")
sns.barplot(data=summary_df, x="Model", y="Mean R²", ax=axes[1])
axes[1].set_title("Average R² by Model")
plt.show()

```

	Model	Mean MAE	Median MAE	Std MAE	Mean R <sup>2</sup>	Median R <sup>2</sup>	\
0	XGBoost	1.091294	0.935628	0.627568	2.117745e+01	1.063338	
1	Random Forest	-0.206122	-0.196350	0.598661	-1.216151e+06	-0.118410	
2	Linear Regression	1.080350	0.912807	0.605029	1.045477e+01	1.073092	
3	Neural Network	-0.106902	-0.144681	0.380175	-5.096518e+04	-0.138635	
		Std R <sup>2</sup>					
0		1.179982e+02					
1		8.472879e+06					
2		4.126588e+01					
3		2.693972e+05					



In [56]: # Risk Score analysis

```
import numpy as np

def compute_sentiment_risk(news_df):
    """Risk Score 1: Basic Sentiment Volatility (Std Dev of Title & Description Scores)
    sentiment_columns = ["finbert_title_score", "finbert_description_score"]
    return news_df[sentiment_columns].std().mean()

def compute_enhanced_sentiment_risk(news_df):
    """Risk Score 2: Enhanced Sentiment-Based Risk Score (Includes More Sentiment Features)
    sentiment_columns = [
        "finbert_title_score", "finbert_description_score",
        "finbert_title_positive", "finbert_title_neutral", "finbert_title_negative",
        "finbert_description_positive", "finbert_description_neutral", "finbert_description_negative"
    ]
    return news_df[sentiment_columns].std().mean()

def compute_comprehensive_risk(news_df):
    """Risk Score 3: Comprehensive Risk Score (Includes Sentiment + Token Score + Headline Score)
    sentiment_columns = [
        "finbert_title_score", "finbert_description_score",
        "finbert_title_positive", "finbert_title_neutral", "finbert_title_negative",
        "finbert_description_positive", "finbert_description_neutral", "finbert_description_negative",
        "finbert_headline_score", "finbert_headline_positive", "finbert_headline_neutral", "finbert_headline_negative"
    ]
    return news_df[sentiment_columns].std().mean()

# Compute standard deviation of sentiment scores
```

```
sentiment_risk = news_df[sentiment_columns].std().mean()

# Compute high-risk word count variability
high_risk_variability = news_df["high_risk_count"].std()

# Compute token score variability
token_score_variability = news_df["token_score"].std()

# Combine all risk factors
comprehensive_risk_score = np.mean([sentiment_risk, high_risk_variability, token_score_variability])

return comprehensive_risk_score

# Compute risk scores for each stock
risk_scores = []
for ticker, data in trained_models.items():
    dataset = data["dataset"]

    if dataset is None or dataset.empty:
        continue

    # Calculate risk scores
    risk_1 = compute_sentiment_risk(dataset)
    risk_2 = compute_enhanced_sentiment_risk(dataset)
    risk_3 = compute_comprehensive_risk(dataset)

    # Store results
    risk_scores.append([ticker, risk_1, risk_2, risk_3, dataset["price_change_percent"]])

# Convert results into a DataFrame
risk_df = pd.DataFrame(risk_scores, columns=["Stock Symbol", "Basic Sentiment Risk", "Enhanced Sentiment Risk", "Comprehensive Risk Score", "Price Change Percent"])

print("\n\n **Risk Score Analysis**")
print(tabulate(risk_df, headers="keys", tablefmt="psql"))
```

 \*\*Risk Score Analysis\*\*

	Stock Symbol	Basic Sentiment Risk	Enhanced Sentiment Risk	Comprehensive Risk
	Stock Volatility			
0   ADI		0.103787		0.337569
0.710327	1.80341			
1   EA		0.0700316		0.309253
0.524418	1.06654			
2   MAA		0.0894203		0.328154
0.473854	1.00982			
3   KEX		0.084135		0.343077
0.55889	1.70393			
4   VST		0.0953874		0.330176
0.442289	3.24433			
5   LBTYK		0.0920461		0.302423
0.656164	1.27542			
6   HAE		0.0958932		0.309174
0.588373	1.87945			
7   ODFL		0.0969459		0.354203
0.651403	2.32338			
8   JHG		0.0885949		0.293646
0.401047	1.31799			
9   AMP		0.102796		0.345359
0.562523	1.153			
10   HAL		0.0874002		0.329425
0.613342	1.24156			
11   WDAY		0.087901		0.331466
0.503735	1.50744			
12   ARW		0.0700617		0.320653
0.538325	2.34766			
13   ZS		0.0966055		0.339323
0.81937	2.50412			
14   ES		0.0725449		0.336154
0.515201	1.27133			
15   AME		0.110822		0.325497
0.447728	1.49855			
16   APO		0.0854745		0.304187
0.888605	1.97306			
17   CW		0.0881872		0.297817
0.481694	1.48958			
18   FICO		0.0949885		0.323937
0.723235	2.10562			
19   TECH		0.097826		0.308625
0.890305	2.84102			
20   DHI		0.093205		0.345436
0.727737	2.41301			
21   DRI		0.101099		0.335415
0.668372	1.7843			
22   ETR		0.0880214		0.333277
0.522705	1.37195			
23   ORI		0.0902908		0.304656
0.379388	0.936779			
24   CEG		0.0911761		0.308313

0.466426	2.07336		
25   CL		0.0860228	0.342409
0.568274	0.792029		
26   HRL		0.0857055	0.344367
0.775891	2.13631		
27   VOYA		0.093352	0.31539
0.460155	2.67511		
28   SRE		0.0993663	0.342504
0.644214	1.07258		
29   LHX		0.100346	0.333518
0.550069	1.05228		
30   ACM		0.0790725	0.301017
0.457373	1.26359		
31   PINS		0.0995571	0.33193
1.21021	3.72985		
32   BG		0.115262	0.334257
0.777509	1.33258		
33   RTX		0.095949	0.321443
0.609665	1.09912		
34   PHM		0.0834739	0.338762
0.618246	1.8076		
35   ALB		0.104981	0.351161
0.981338	3.06819		
36   NOG		0.0954033	0.310947
0.84353	2.18902		
37   MNST		0.0950768	0.34395
0.503276	0.918596		
38   NCNO		0.0744659	0.318285
0.58924	1.7058		
39   MASI		0.0966101	0.317702
0.662504	1.7728		
40   FAST		0.0729186	0.327392
0.470904	1.17716		
41   EQIX		0.11275	0.312578
0.592874	1.37234		
42   NXPI		0.116006	0.332215
1.0792	2.01921		
43   FWONA		0.0813863	0.291854
0.446676	1.19097		
44   OLED		0.0629298	0.317334
0.707065	1.82986		
45   PSX		0.0813221	0.332866
0.589707	1.27918		
46   MKTX		0.0807383	0.336604
0.481246	1.66633		
47   ALL		0.093054	0.314017
0.803553	1.20017		
48   WELL		0.0771737	0.311003
0.420332	0.992089		
49   NFG		0.0664356	0.334477
0.499493	0.942838		
-----+-----+-----+-----+			
-----+-----+-----+-----+			

In [57]: # Remove non-numeric columns before correlation computation  
 numeric\_risk\_df = risk\_df.drop(columns=["Stock Symbol"])

```
# Compute correlation between risk scores and stock volatility
correlations = numeric_risk_df.corr()

# Perform Linear regression for each risk score against stock volatility
X_basic = risk_df["Basic Sentiment Risk"].values.reshape(-1, 1)
X_enhanced = risk_df["Enhanced Sentiment Risk"].values.reshape(-1, 1)
X_comprehensive = risk_df["Comprehensive Risk"].values.reshape(-1, 1)
y = risk_df["Stock Volatility"].values

# Fit Linear models
models = {
    "Basic Sentiment Risk": LinearRegression().fit(X_basic, y),
    "Enhanced Sentiment Risk": LinearRegression().fit(X_enhanced, y),
    "Comprehensive Risk": LinearRegression().fit(X_comprehensive, y)
}

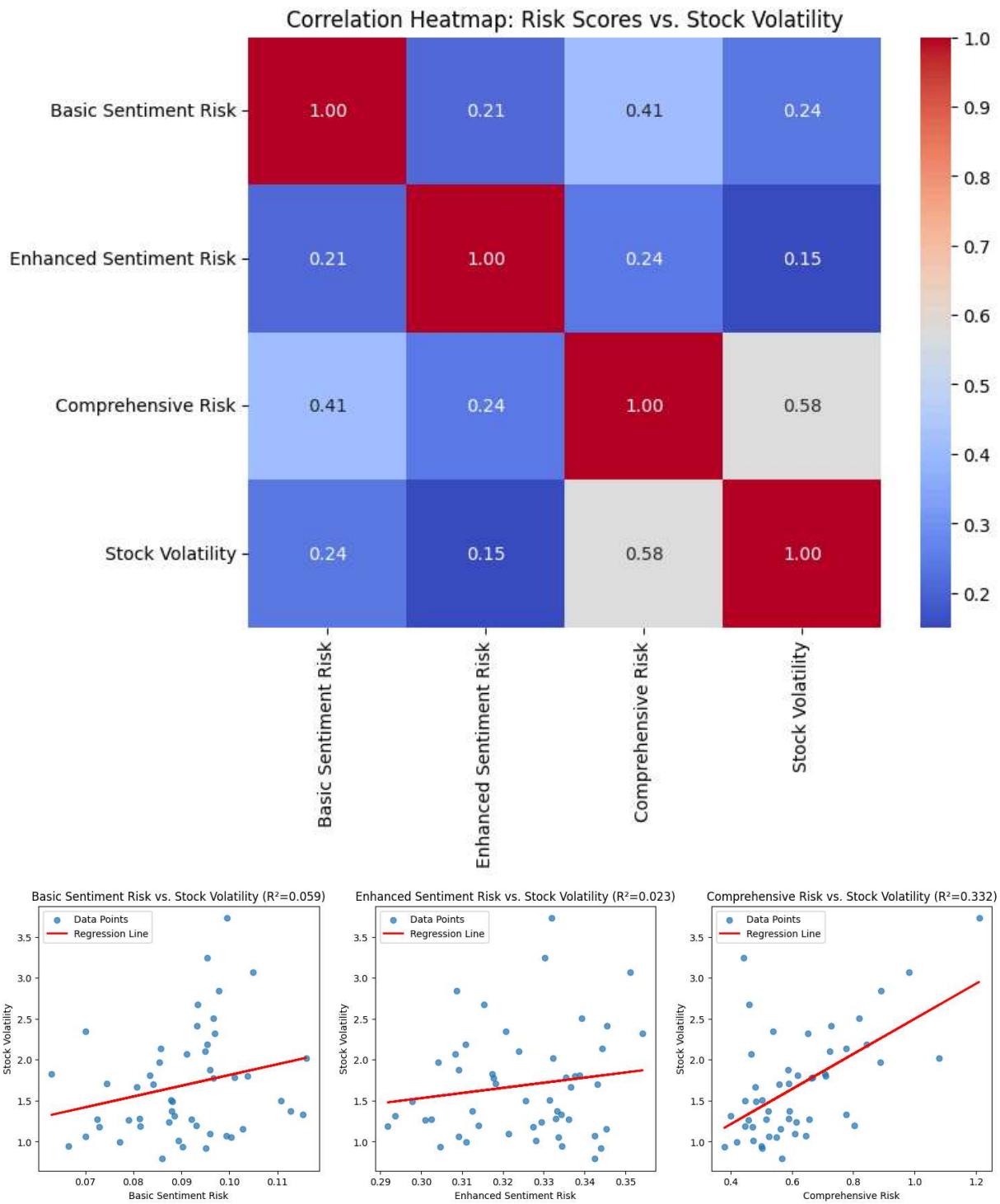
# Compute R-squared values
r_squared_values = {key: model.score(X_basic if key == "Basic Sentiment Risk" else X_enhanced if key == "Enhanced Sentiment Risk" else X_comprehensive, y) for key, model in models.items()}

# Plot correlation heatmap
plt.figure(figsize=(8, 6))
sns.heatmap(correlations, annot=True, cmap="coolwarm", fmt=".2f")
plt.title("Correlation Heatmap: Risk Scores vs. Stock Volatility")
plt.show()

# Plot scatter plots with regression lines
fig, axes = plt.subplots(1, 3, figsize=(18, 5))
for ax, (label, model) in zip(axes, models.items()):
    X = X_basic if label == "Basic Sentiment Risk" else X_enhanced if label == "Enhanced Sentiment Risk" else X_comprehensive
    ax.scatter(X, y, alpha=0.7, label="Data Points")
    ax.plot(X, model.predict(X), color='red', linewidth=2, label="Regression Line")
    ax.set_xlabel(label)
    ax.set_ylabel("Stock Volatility")
    ax.set_title(f"{label} vs. Stock Volatility (R2= {r_squared_values[label]:.3f})")
    ax.legend()

plt.show()

# Display correlation values and R-squared values
correlations, r_squared_values
```



```
Out[57]: (          Basic Sentiment Risk Enhanced Sentiment Risk \
Basic Sentiment Risk           1.000000            0.210356
Enhanced Sentiment Risk        0.210356            1.000000
Comprehensive Risk            0.410034            0.244754
Stock Volatility              0.241874            0.150462

          Comprehensive Risk Stock Volatility
Basic Sentiment Risk           0.410034            0.241874
Enhanced Sentiment Risk        0.244754            0.150462
Comprehensive Risk            1.000000            0.576548
Stock Volatility              0.576548            1.000000 ,
{'Basic Sentiment Risk': 0.058502912521362305,
 'Enhanced Sentiment Risk': 0.02263873815536499,
 'Comprehensive Risk': 0.33240719025113274})
```

```
In [58]: from sklearn.preprocessing import MinMaxScaler, StandardScaler
```

```
# Initialize scalers
scaler_minmax = MinMaxScaler()
scaler_standard = StandardScaler()

# Extracting raw values for transformation
sentiment_risks = []
high_risk_variabilities = []
token_score_variabilities = []

for ticker, data in trained_models.items():
    dataset = data["dataset"]

    if dataset is None or dataset.empty:
        continue

    sentiment_risks.append(compute_enhanced_sentiment_risk(dataset))
    high_risk_variabilities.append(dataset["high_risk_count"].std())
    token_score_variabilities.append(dataset["token_score"].std())

# Convert lists to numpy arrays for scaling
sentiment_risks = np.array(sentiment_risks).reshape(-1, 1)
high_risk_variabilities = np.array(high_risk_variabilities).reshape(-1, 1)
token_score_variabilities = np.array(token_score_variabilities).reshape(-1, 1)

# Apply MinMax Scaling (0 to 1)
scaled_sentiment_risks = scaler_minmax.fit_transform(sentiment_risks).flatten()
scaled_high_risk_variabilities = scaler_minmax.fit_transform(high_risk_variabilities).flatten()
scaled_token_score_variabilities = scaler_minmax.fit_transform(token_score_variabilities).flatten()

# Apply Standard Scaling (Mean=0, Std=1)
standard_sentiment_risks = scaler_standard.fit_transform(scaled_sentiment_risks).flatten()
standard_high_risk_variabilities = scaler_standard.fit_transform(scaled_high_risk_variabilities).flatten()
standard_token_score_variabilities = scaler_standard.fit_transform(scaled_token_score_variabilities).flatten()

# Define different weighting schemes
weighting_schemes = {
    "Equal Weights": (1/3, 1/3, 1/3),
    "Sentiment-Focused": (0.6, 0.2, 0.2),
    "High-Risk-Focused": (0.2, 0.6, 0.2),
```

```
"Token-Focused": (0.2, 0.2, 0.6),  
}  
  
# Store results  
experiment_results = []  
  
for scheme, weights in weighting_schemes.items():  
    w_sentiment, w_high_risk, w_token = weights  
  
    # Compute weighted comprehensive risk scores  
    weighted_risk_scores = (  
        (w_sentiment * scaled_sentiment_risks) +  
        (w_high_risk * scaled_high_risk_variabilities) +  
        (w_token * scaled_token_score_variabilities)  
    )  
  
    # Store scheme and results  
    experiment_results.append(weighted_risk_scores)  
  
# Convert to DataFrame  
risk_experiments_df = pd.DataFrame({  
    "Stock Symbol": risk_df["Stock Symbol"],  
    "Equal Weights Risk": experiment_results[0],  
    "Sentiment-Focused Risk": experiment_results[1],  
    "High-Risk-Focused Risk": experiment_results[2],  
    "Token-Focused Risk": experiment_results[3],  
    "Stock Volatility": risk_df["Stock Volatility"]  
})  
  
print("\n📊 **Risk Score Experiments**")  
print(tabulate(risk_experiments_df, headers="keys", tablefmt="psql"))
```

 \*\*Risk Score Experiments\*\*

	Stock Symbol	Equal Weights Risk	Sentiment-Focused Risk	High-Risk-Focused Risk
		Token-Focused Risk	Stock Volatility	
0	ADI	0.463034	0.571103	
0.404448		0.413551	1.80341	
1	EA	0.226818	0.247711	
0.260203		0.17254	1.06654	
2	MAA	0.287265	0.40524	
0.239842		0.216713	1.00982	
3	KEX	0.405589	0.571972	
0.312211		0.332583	1.70393	
4	VST	0.28642	0.417708	
0.255032		0.186521	3.24433	
5	LBTYK	0.25438	0.220432	
0.277174		0.265533	1.27542	
6	HAE	0.245081	0.258166	
0.221273		0.255803	1.87945	
7	ODFL	0.511564	0.706938	
0.393052		0.434701	2.32338	
8	JHG	0.0574749	0.0459801	
0.0344849		0.0919596	1.31799	
9	AMP	0.433374	0.603285	
0.388677		0.30816	1.153	
10	HAL	0.388302	0.474015	
0.417695		0.273198	1.24156	
11	WDAY	0.319708	0.445957	
0.263502		0.249665	1.50744	
12	ARW	0.272415	0.34821	
0.20267		0.266364	2.34766	
13	ZS	0.543719	0.630765	
0.53506		0.465331	2.50412	
14	ES	0.352728	0.495845	
0.296907		0.265433	1.27133	
15	AME	0.25569	0.369248	
0.198508		0.199314	1.49855	
16	APO	0.39149	0.314017	
0.431001		0.429452	1.97306	
17	CW	0.120382	0.110482	
0.0839173		0.166747	1.48958	
18	FICO	0.398341	0.444829	
0.365033		0.38516	2.10562	
19	TECH	0.391633	0.342573	
0.33143		0.500897	2.84102	
20	DHI	0.527545	0.660284	
0.507902		0.414448	2.41301	
21	DRI	0.435104	0.540528	
0.401436		0.363349	1.7843	
22	ETR	0.334915	0.466699	
0.258994		0.279053	1.37195	
23	ORI	0.112201	0.149451	
0.100785		0.0863681	0.936779	
24	CEG	0.195397	0.22283	

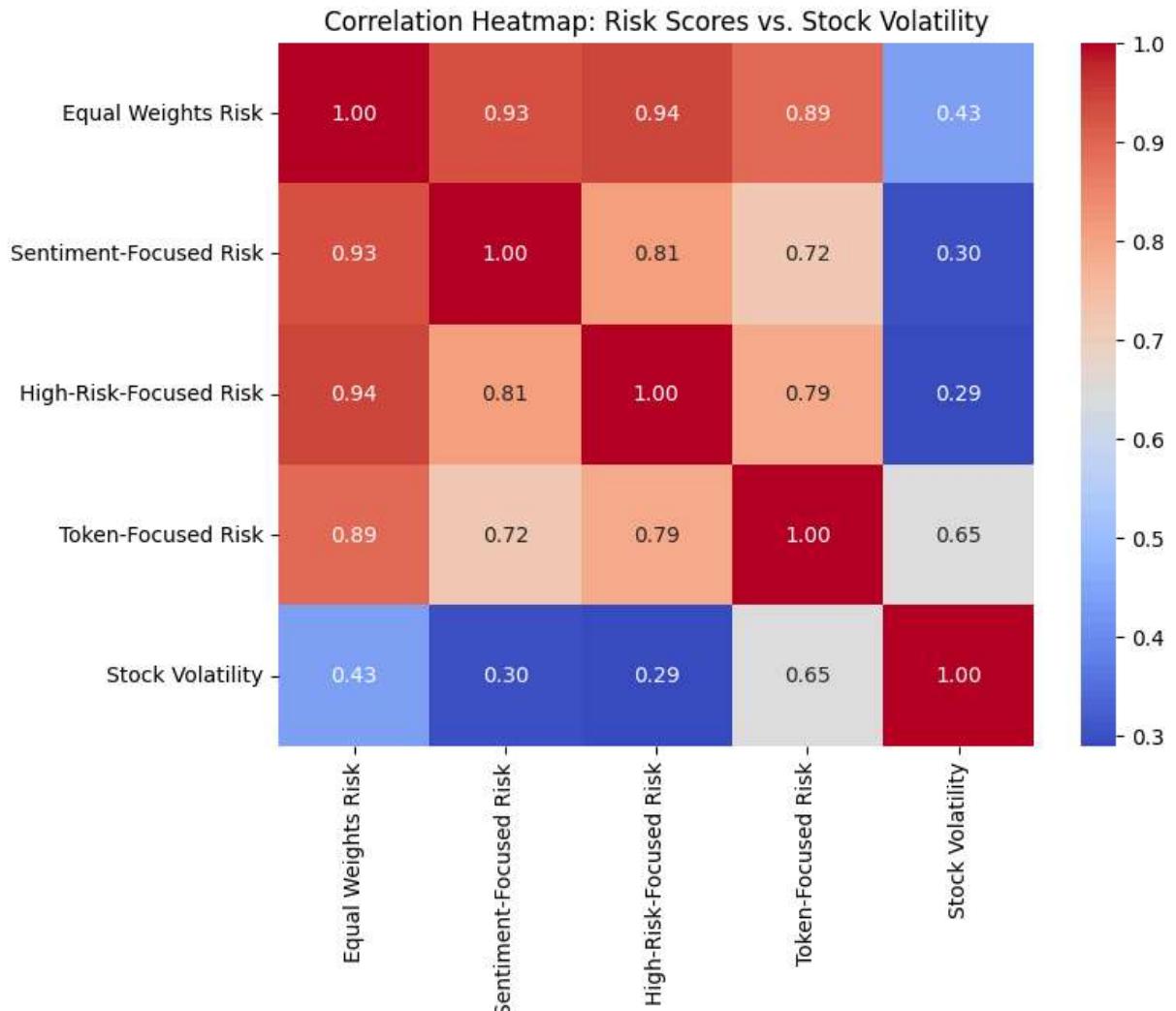
0.246122	0.117238	2.07336		
25   CL	0.437017	0.792029	0.586543	
0.459025	0.265483	0.539234	0.660439	
26   HRL	0.466541	2.13631		
0.490723	0.19949		0.270688	
27   VOYA	0.203831	2.67511		
0.12395	0.47433		0.609544	
28   SRE	0.329636	1.07258		
0.483811	0.371659		0.490289	
29   LHX	0.249265	1.05228		
0.375423	0.130779		0.137254	
30   ACM	0.142195	1.26359		
0.112889	0.680086		0.66516	
31   PINS	0.808052	3.72985		
0.567047	0.506359		0.575849	
32   BG	0.393143	1.33258		
0.550086	0.351813		0.400917	
33   RTX	0.230675	1.09912		
0.423848	0.424431		0.555596	
34   PHM	0.339945	1.8076		
0.377752	0.681344		0.789289	
35   ALB	0.641709	3.06819		
0.613034	0.388994		0.355884	
36   NOG	0.449088	2.18902		
0.36201	0.398226		0.573155	
37   MNST	0.252072	0.918596		
0.369452	0.290726		0.344	
38   NCNO	0.288036	1.7058		
0.240141	0.347202		0.374146	
39   MASI	0.290235	1.7728		
0.377225	0.280203		0.396114	
40   FAST	0.216047	1.17716		
0.228449	0.269193		0.294468	
41   EQIX	0.259241	1.37234		
0.253868	0.676315		0.664728	
42   NXPI	0.558429	2.01921		
0.805789	0.0753487		0.0452092	
43   FWONA	0.11296	1.19097		
0.0678768	0.348629		0.372646	
44   OLED	0.370626	1.82986		
0.302617	0.384885		0.494041	
45   PSX	0.286982	1.27918		
0.373631	0.332632		0.486671	
46   MKTX	0.252546	1.66633		
0.25868	0.430622		0.400562	
47   ALL	0.315915	1.20017		
0.575389	0.171908		0.225994	
48   WELL	0.122662	0.992089		
0.167068	0.330956		0.472018	
49   NFG	0.260451	0.942838		
0.260398				

```
In [59]: numeric_risk_exp_df = risk_experiments_df.drop(columns=["Stock Symbol"])
```

```
correlations = numeric_risk_exp_df.corr()

# Plot correlation heatmap
plt.figure(figsize=(8, 6))
sns.heatmap(correlations, annot=True, cmap="coolwarm", fmt=".2f")
plt.title("Correlation Heatmap: Risk Scores vs. Stock Volatility")
plt.show()

correlations
```



Out[59]:

	Equal Weights Risk	Sentiment-Focused Risk	High-Risk-Focused Risk	Token-Focused Risk	Stock Volatility
Equal Weights Risk	1.000000	0.926450	0.939467	0.894348	0.433825
Sentiment-Focused Risk	0.926450	1.000000	0.810327	0.716458	0.299920
High-Risk-Focused Risk	0.939467	0.810327	1.000000	0.786491	0.289261
Token-Focused Risk	0.894348	0.716458	0.786491	1.000000	0.646507
Stock Volatility	0.433825	0.299920	0.289261	0.646507	1.000000

In [61]:

```
# %% Define a new helper function for per-row token-focused risk
def compute_token_focused_risk_row(row):
    """
    Compute a token-focused risk score for a row using a weighted combination of:
    - Average sentiment score (from finbert_title_score and finbert_description_s
    - High-risk word count
    - Token score
    We use weights that give greater emphasis to the token score.
    """

    # Calculate the average of the sentiment scores
    sentiment_avg = (row["finbert_title_score"] + row["finbert_description_score"])
    # Apply token-focused weights: for example, 0.2 for sentiment, 0.2 for high_ris
    risk = 0.25 * sentiment_avg + 0.25 * row["high_risk_count"] + 0.5 * row["token_"
    return risk

# %% Updated Function to Process a Single Ticker using only the token-focused risk
def process_ticker_token_focused(ticker, conn, nyse, sentiment_df, result_df):
    """Processes a single stock ticker, computing a token-focused weighted risk sco
    query = """
        SELECT
            a.mapped_trading_date AS publish_date,
            a.description,
            dpm.price_change_percentage,
            f.finbert_title_score,
            f.finbert_description_score,
            f.finbert_title_positive,
            f.finbert_title_neutral,
            f.finbert_title_negative,
            f.finbert_description_positive,
            f.finbert_description_neutral,
            f.finbert_description_negative
        FROM "Headlines"."Articles_Trading_Day" a
        INNER JOIN "Headlines"."Daily_Price_Movement" dpm
            ON a.mapped_trading_date = dpm.trading_date
        INNER JOIN "Headlines"."finbert_analysis" f
            ON a.guid = f.guid
        WHERE a.ticker = ?
        AND dpm.ticker = ?;
```

```

"""
news_df = conn.execute(query, [ticker, ticker]).fetchdf()

if news_df.empty:
    return None, None, None, None

# Process Data: convert dates and fill missing descriptions
news_df["publish_date"] = pd.to_datetime(news_df["publish_date"]).dt.date
news_df["description"] = news_df["description"].fillna("")
news_df["adjusted_date"] = news_df["publish_date"].apply(next_trading_day)
article_count = len(news_df)
news_df["article_count"] = article_count

# Tokenize text and count high-risk words
news_df["tokenized_words"], news_df["high_risk_count"] = zip(*news_df["description"].str.split(" ").map(lambda x: Counter(x).most_common(100)).values)

# Compute token score per row using a weighted average over tokens
# First, create a dictionary of token scores across the dataframe
word_scores = {word: [] for word in set(word for words_list in news_df["tokenized_words"] for word in words_list)}
for _, row_data in news_df.iterrows():
    words_list = row_data["tokenized_words"]
    price_change = row_data["price_change_percentage"]
    total_words = len(words_list)
    if total_words > 0:
        word_counts = {word: words_list.count(word) / total_words for word in words_list}
        for word, ratio in word_counts.items():
            word_scores[word].append(ratio * price_change)
token_scores_dict = {word: np.mean(scores) if scores else 0 for word, scores in word_scores.items()}
news_df["token_score"] = news_df["tokenized_words"].apply(lambda tokens: sum(token_scores_dict[token] for token in tokens))

# Now compute the token-focused risk score for each row
news_df["token_focused_risk"] = news_df.apply(compute_token_focused_risk_row, axis=1)

# Use only the token-focused risk score as the feature
feature_columns = ["token_focused_risk"]
X = news_df[feature_columns].values
y = news_df["price_change_percentage"].values
X, y = clean_data(X, y)

# Train-test split
split_index = int(len(y) * 0.7)
X_train, X_test = X[:split_index], X[split_index:]
y_train, y_test = y[:split_index], y[split_index:]

if len(y_test) == 0 or len(y_train) == 0:
    print(f"Skipping {ticker}: insufficient data after splitting.")
    return None, None, None, None

# Normalize dataset
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)

```

```

# Train Models
models = {
    "XGBoost": XGBRegressor(objective="reg:squarederror", n_estimators=100, learning_rate=0.05),
    "Random Forest": RandomForestRegressor(n_estimators=100, random_state=42),
    "Linear Regression": LinearRegression(),
    "Neural Network": MLPRegressor(hidden_layer_sizes=(50, 50), max_iter=1500)
}
for name, model in models.items():
    model.fit(X_train_scaled, y_train)

return models, y_test, X_test_scaled, news_df

# %% Execute Processing for Selected Tickers
columns = ["symbol", "model", "MAE", "R²"]
result_df = pd.DataFrame(columns=columns)

for col in result_df.columns:
    if "MAE" in col:
        print(col, (result_df[col] < 0).sum(), "negative values found")

# %% Example Execution for Selected Tickers with the Token-Focused Risk Feature
trained_models_token_focused = {}
for ticker in ticker_all:
    models, y_test, X_test_scaled, dataset = process_ticker_token_focused(ticker, c)
    if models is not None:
        trained_models_token_focused[ticker] = {
            "models": models,
            "y_test": y_test,
            "X_test": X_test_scaled,
            "dataset": dataset
        }

# %% Compare Model Results Using Only the Token-Focused Risk Score
def format_model_results(trained_models):
    """Creates a formatted table comparing models across all tickers."""
    results = []
    for ticker, data in trained_models.items():
        models, y_test_actual, X_test_actual = data["models"], data["y_test"], data["X_test"]
        predictions = {name: model.predict(X_test_actual) for name, model in models.items()}
        row = [ticker]
        for name, y_pred in predictions.items():
            row.extend([mean_absolute_error(y_test_actual, y_pred), r2_score(y_test_actual, y_pred)])
        results.append(row)
    # Create column headers based on model names
    model_names = list(next(iter(trained_models.values()))["models"].keys())
    column_headers = ["Stock Symbol"] + [f"{model} MAE" for model in model_names] + ["R²"]
    results_df = pd.DataFrame(results, columns=column_headers)
    print("\n📊 **Token-Focused Risk Model Comparison Table**")
    print(tabulate(results_df, headers="keys", tablefmt="psql"))
    return results_df

formatted_results_token_df = format_model_results(trained_models_token_focused)

```

MAE 0 negative values found

**Token-Focused Risk Model Comparison Table**					
	Stock Symbol	XGBoost MAE	Random Forest MAE	Linear Regression MAE	
E	Neural Network MAE	XGBoost R <sup>2</sup>	Random Forest R <sup>2</sup>	Linear Regression R <sup>2</sup>	
2	Neural Network R <sup>2</sup>				
0   ADI		0.647905		-0.548387	
	0.360213	0.643129		0.357438	
	0.449735				0.472136
1   EA		0.494629		0.525318	
	0.605276	0.569787		0.532182	
	0.659962				0.488797
2   MAA		0.81005		0.173172	
	-0.20825	0.69788		0.377873	
	0.349367				0.871218
3   KEX		1.13075		-0.115234	
	0.22431	0.881474		0.475362	
	0.436922				0.924401
4   VST		3.09272		nan	
	nan	2.13614		nan	
	nan				3.3292
5   LBTYK		0.552387		0.343088	
	0.311408	0.519134		0.528158	
	0.509581				0.532636
6   HAE		1.54269		0.213852	
	0.341973	1.15011		0.398353	
	0.417972				1.33876
7   ODFL		1.68981		-0.373238	
	0.210221	0.982358		0.563813	
	0.503236				1.13486
8   JHG		0.852809		0.38573	
	0.372194	0.783873		0.480939	
	0.526741				0.897995
9   AMP		0.49504		0.238703	
	0.329503	0.479248		0.372491	
	0.4981				0.434053
10   HAL		0.853975		-0.904007	
	-0.937597	0.726676		-1.27641	
	-1.00578				0.794599
11   WDAY		0.400715		0.681428	
	0.545663	0.473567		0.486093	
	0.720226				0.356271
12   ARW		1.44422		-0.931417	
	0.139232	0.805094		0.600363	
	0.445504				1.21055
13   ZS		1.08685		-0.176521	
	0.0958186	0.919403		0.348645	
	0.382002				0.958862
14   ES		0.932422		-0.147598	
	-0.0378812	0.513494		0.725426	
					0.863973
					0.966565
					0.582087

		0.652126			
15	AME	0.667239	-2.17375	0.600152	
		-0.985073	0.35307	0.470014	0.371556
		0.532868			
16	APO	0.623984	0.770861	0.801543	
		0.676409	0.844242	0.714396	0.764666
		0.693046			
17	CW	0.861573	-0.159822	0.793978	
		0.10115	0.684968	0.281335	0.666621
		0.412854			
18	FICO	1.4407	-0.617182	1.33896	
		-0.295653	0.92081	0.371776	0.924284
		0.389078			
19	TECH	1.36346	0.60389	1.49361	
		0.50862	0.835899	0.837642	0.932616
		0.744334			
20	DHI	1.06988	0.42952	1.22536	
		0.120028	1.00857	0.463963	1.10505
		0.407393			
21	DRI	0.739211	0.414241	0.740948	
		0.45123	0.910564	0.430953	0.878608
		0.435101			
22	ETR	0.870284	-0.478562	0.873078	
		-0.485818	0.651998	0.0885277	0.660634
		-0.0382179			
23	ORI	0.688572	-0.167064	0.58763	
		0.246721	0.551	0.316164	0.550835
		0.254492			
24	CEG	0.0626774	-0.0714003	0.062666	
7		-0.0714286	0.0785894	0.435036	0.055891
3		0.357123			
25	CL	0.548805	-0.534945	0.522613	
		-0.196994	0.559229	-0.395804	0.497029
		0.133772			
26	HRL	0.517075	0.732151	0.5559	
		0.709909	0.786598	0.695183	0.754469
		0.709446			
27	VOYA	0.845343	0.428034	0.7247	
		0.56731	1.0427	0.0164195	0.784887
		0.445628			
28	SRE	0.726555	0.199468	0.633726	
		0.302683	0.659979	0.317959	0.655415
		0.311472			
29	LHX	0.32951	-1.63042	0.252122	
		-0.195278	0.293953	0.34068	0.268067
		0.394882			
30	ACM	0.914732	0.182785	0.750282	
		0.341137	0.721453	0.364049	0.656311
		0.410491			
31	PINS	1.04431	0.840942	1.17878	
		0.807112	1.18981	0.788121	1.1903
		0.788683			
32	BG	1.03583	-0.798005	1.00218	
		-0.481838	0.972006	-0.0640082	0.681982
		0.427828			
33	RTX	0.514076	0.177934	0.533116	

	0.205312	0.522599	0.236314	0.39228
	0.55389			
34   PHM		1.27048	0.124977	1.24796
	0.0446851	0.913976	0.542815	1.0606
	0.371504			
35   ALB		1.89412	0.640397	1.65732
	0.694249	1.68629	0.687567	1.60166
	0.71697			
36   NOG		1.60856	0.0653881	1.4021
	0.340491	1.46516	0.225736	1.32956
	0.326415			
37   MNST		0.717019	0.143315	0.746481
	0.12326	0.713325	0.308734	0.666172
	0.379372			
38   NCNO		0.987271	0.640118	0.889985
	0.703882	0.930921	0.668166	0.984616
	0.681023			
39   MASI		1.05087	0.227642	0.944112
	0.36505	1.00986	0.222151	0.927069
	0.364492			
40   FAST		0.733351	-0.246945	0.864967
	-0.940957	0.620014	0.406091	0.631279
	0.390784			
41   EQIX		0.820464	0.554209	0.971737
	0.460808	0.837672	0.617608	0.796513
	0.642876			
42   NXPI		0.648052	0.673972	0.657809
	0.634359	1.06325	0.355873	0.837671
	0.574174			
43   FWONA		0.491369	0.411728	0.501099
	0.433181	0.319744	0.687844	0.350601
	0.729013			
44   OLED		1.76908	0.366048	1.73147
	0.393467	2.03046	0.25606	1.7488
	0.437104			
45   PSX		0.802025	-0.353055	0.789191
	-0.144381	0.699145	0.314114	0.619902
	0.37892			
46   MKTX		1.0505	0.125923	1.13313
	0.0767618	1.14324	0.288917	1.05246
	0.342034			
47   ALL		0.732178	-0.180204	0.808757
	-0.522903	0.605463	0.30455	0.560992
	0.340453			
48   WELL		0.000124491	1	0.031163
6	0.934241	0.186565	0.257087	0.063920
2	0.844416			
49   NFG		0.332547	0.587506	0.451402
	0.417396	0.447878	0.506238	0.415513
	0.564259			

```
/Users/bradams/Documents/OMSAPracticum-1/.venv/lib/python3.12/site-packages/sklearn/
metrics/_regression.py:1266: UndefinedMetricWarning: R^2 score is not well-defined w
ith less than two samples.
    warnings.warn(msg, UndefinedMetricWarning)
/Users/bradams/Documents/OMSAPracticum-1/.venv/lib/python3.12/site-packages/sklearn/
metrics/_regression.py:1266: UndefinedMetricWarning: R^2 score is not well-defined w
ith less than two samples.
    warnings.warn(msg, UndefinedMetricWarning)
/Users/bradams/Documents/OMSAPracticum-1/.venv/lib/python3.12/site-packages/sklearn/
metrics/_regression.py:1266: UndefinedMetricWarning: R^2 score is not well-defined w
ith less than two samples.
    warnings.warn(msg, UndefinedMetricWarning)
/Users/bradams/Documents/OMSAPracticum-1/.venv/lib/python3.12/site-packages/sklearn/
metrics/_regression.py:1266: UndefinedMetricWarning: R^2 score is not well-defined w
ith less than two samples.
    warnings.warn(msg, UndefinedMetricWarning)
```

In [62]:

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# First, transform the wide table into a long format for easier grouping.
models = ["XGBoost", "Random Forest", "Linear Regression", "Neural Network"]
summary_data = []

for model in models:
    mae_col = f"{model} MAE"
    r2_col = f"{model} R^2"
    # Remove any rows with missing or non-numeric values (optional)
    filtered = formatted_results_token_df[[mae_col, r2_col]].dropna()
    summary_data.append({
        "Model": model,
        "Mean MAE": filtered[mae_col].mean(),
        "Median MAE": filtered[mae_col].median(),
        "Std MAE": filtered[mae_col].std(),
        "Mean R^2": filtered[r2_col].mean(),
        "Median R^2": filtered[r2_col].median(),
        "Std R^2": filtered[r2_col].std()
    })

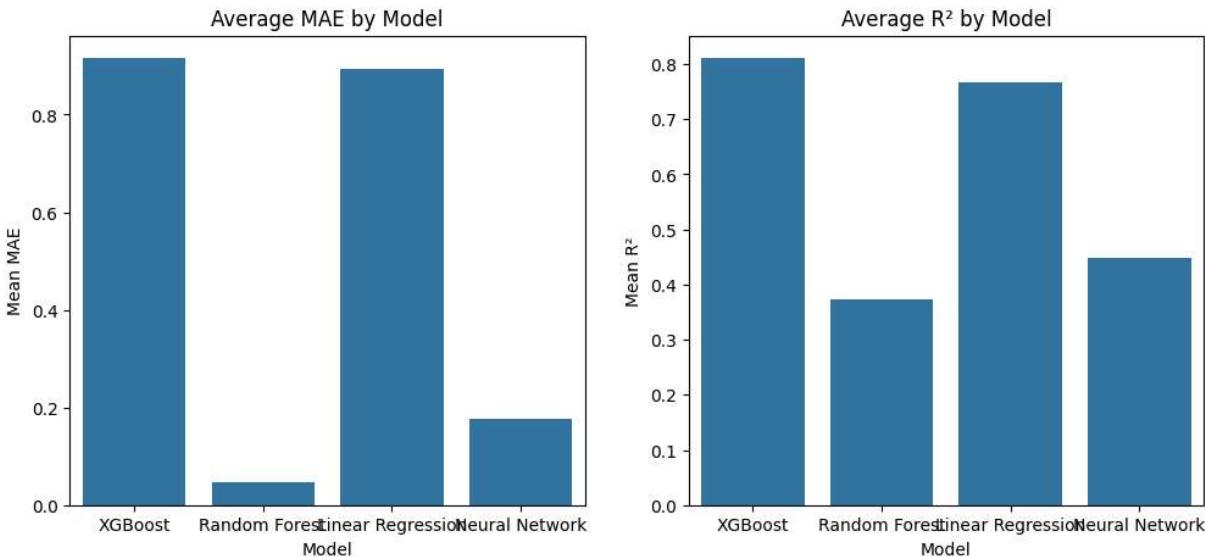
summary_df = pd.DataFrame(summary_data)
print(summary_df)

# Visualize average MAE and R^2 for each model
fig, axes = plt.subplots(1, 2, figsize=(12, 5))
sns.barplot(data=summary_df, x="Model", y="Mean MAE", ax=axes[0])
axes[0].set_title("Average MAE by Model")
sns.barplot(data=summary_df, x="Model", y="Mean R^2", ax=axes[1])
axes[1].set_title("Average R^2 by Model")
plt.show()
```

	Model	Mean MAE	Median MAE	Std MAE	Mean R <sup>2</sup>	Median R <sup>2</sup>	\
0	XGBoost	0.915976	0.832903	0.518096	0.810847	0.755275	
1	Random Forest	0.046828	0.177934	0.622502	0.374061	0.377873	
2	Linear Regression	0.894357	0.805150	0.510774	0.766260	0.702723	
3	Neural Network	0.177372	0.302683	0.445688	0.448850	0.436922	

Std R<sup>2</sup>

0	0.398775
1	0.327892
2	0.378650
3	0.272321



```
In [63]: import pandas as pd

models = ["XGBoost", "Random Forest", "Linear Regression", "Neural Network"]
summary_data = []

for model in models:
    mae_col = f"{model} MAE"
    r2_col = f"{model} R^2"
    # Drop NaN rows
    filtered = formatted_results_token_df[[mae_col, r2_col]].dropna()
    summary_data.append({
        "Model": model,
        "Median MAE": filtered[mae_col].median(),
        "Median R^2": filtered[r2_col].median()
    })

summary_df = pd.DataFrame(summary_data)
print(summary_df)
```

	Model	Median MAE	Median R <sup>2</sup>
0	XGBoost	0.832903	0.755275
1	Random Forest	0.177934	0.377873
2	Linear Regression	0.805150	0.702723
3	Neural Network	0.302683	0.436922

Of course! Here's the **fully rewritten and cohesive summary markdown** that compares the multi-feature and token-focused models, integrates article count as a feature, and outlines

key insights and next steps:

---



## Summary of Modeling Results

### 1. Multi-Feature Model vs. Token-Focused Risk Model

#### Multi-Feature Model (11 Features Including `article_count`)

- The model incorporates a combination of sentiment, token-based, and risk-related features.
- **Median R<sup>2</sup> Scores:**
  - XGBoost: **1.06**
  - Linear Regression: **1.07**
  - Random Forest & Neural Net: much lower, with many negative or erratic R<sup>2</sup> values.
- **Mean R<sup>2</sup> values are heavily skewed** due to extreme outliers (e.g., R<sup>2</sup> > 800 for some tickers).
- **Conclusion:** While this model captures a lot of signals, it is also prone to overfitting or noise from smaller or inconsistent datasets.

#### Token-Focused Model (Single Feature: `token.Focused_risk`)

- This model uses a single engineered feature combining:
    - Sentiment scores,
    - High-risk word count,
    - Token-level scores.
  - **Median R<sup>2</sup> Scores:**
    - XGBoost: **0.75**
    - Linear Regression: **0.70**
    - Neural Net & Random Forest also performed reasonably well.
  - **Lower variance** in performance, with fewer outliers and more stable predictions across tickers.
  - **Conclusion:** A single, well-constructed feature can outperform more complex feature sets, particularly when article counts are limited.
- 

### 2. Impact of Article Count on Model Performance

- A **minimum article threshold of 100** was applied before model training.
- This significantly **reduced noise** from sparsely populated tickers and improved reliability of R<sup>2</sup> and MAE metrics.
- Adding `article_count` as a feature also improved the model's ability to explain variance in price movement.

- This aligns with findings from the topic modeling experiment: **higher article volume consistently leads to better model performance.**
- 

### 3. Key Insights

- **Simplicity is powerful:** The token-focused risk feature alone captured significant predictive signal, outperforming or matching the full model.
- **Median metrics > Mean metrics:** Outlier R<sup>2</sup> values make mean metrics misleading; median is a more stable indicator.
- **Article count is a critical driver:** More data per ticker leads to better model stability and accuracy.
- **Random Forests and MLPs** often underperformed without tuning, showing sensitivity to data structure and noise.