# Comprehensive Guide to Stock Market Forecasting: Concepts, Data, Modeling, and Strategy

This guide outlines a phased approach to building a sophisticated stock market forecasting system. It integrates market basics, data handling, feature engineering, advanced modeling techniques (Machine Learning/Deep Learning), evaluation, and strategic decision-making, incorporating technical, fundamental, and macroeconomic perspectives. The goal is to predict future stock prices, volatility, and risk, translating these forecasts into actionable Buy/Sell/Hold decisions for multiple companies.

---

## Phase 1: Market Education & Key Concepts

### Core Market Data:

- **Price (Open, High, Low, Close - OHLC)**: Represents a stock's trading range and final value. 'Close' is the most common target for forecasting.
- **Volume**: Number of shares traded. Indicates interest, liquidity, and potential momentum.
- **Returns**: % change in price. Used to model risk and profit potential.

### Technical Analysis Terms & Indicators:

- **Moving Averages (SMA/EMA)**: Used to detect trends by smoothing out price data.
- **Volatility (Std Dev, ATR)**: Indicates uncertainty/risk.
- **RSI**: Momentum indicator (0-100). RSI > 70 is overbought; < 30 is oversold.
- **MACD**: Momentum indicator using two EMAs. Signals trend reversals.
- **Bollinger Bands**: Price envelopes indicating volatility.
- **Support & Resistance**: Historical price levels acting as floors/ceilings.
- **Breakout**: When price moves beyond support/resistance.
- **Hold**: Price maintains a breakout. Validates a new trend.
- **Trend**: General direction of price (up/down/sideways).

### Fundamental Analysis & Valuation:

- **Market Cap**: Company size = Price × Shares.
- **P/E Ratio**: Price / EPS. Indicates valuation.
- **EPS**: Net income per share. Profitability measure.
- **Dividend Yield**: Dividend / Price. Income-focused metric.
- **Valuation**: DCF, P/E, P/S, P/B used to assess intrinsic value.

- **Equity**: Ownership (Assets - Liabilities).

## Trading Actions:

- **Buy**: Purchase shares expecting rise.
- **Sell**: Exit position to take profit or cut loss.
- **Hold**: Maintain position pending further data.

## Economic & Market-Wide Indicators:

- **CPI/Inflation**: Cost of living changes. Affects rates.
- **Interest Rates**: Borrowing cost. Influences valuation.
- **Oil/Gold Prices**: Affect inflation/safe-haven flows.
- **VIX**: Market volatility expectation.
- **Indices (e.g., S&P 500)**: Overall market direction.
- **Bond Yields/Spreads**: Economic sentiment & recession signals.

## Other Key Terms:

- **Liquidity, Market Sentiment, Correction, Bear/Bull Market, Portfolio, Diversification, Earnings Report, IPO, M&A, Ticker, Exchange**

---

# Phase 2: Data Collection (Multi-Source)

## Company Data:

- **OHLCV**: From yfinance, Alpha Vantage, Finnhub.
- **Fundamentals**: EPS, P/E, Market Cap, ROE, etc. via yfinance, FMP API.
- **Technicals**: RSI, MACD via pandas-ta.
- **Earnings Dates**: nasdaq.com, fmpcloud.io.

## Macroeconomic Data:

- **CPI, Rates, Unemployment, Bond Yields**: `fredapi`.
- **Oil, Gold**: yfinance, quandl.
- **VIX, Indices**: yfinance, FRED.
- **Exchange Rates**: Alpha Vantage.

---

# Phase 3: Data Preprocessing & Alignment

- **Handle Missing Values**: Forward fill macro; drop where needed.
- **Align Frequencies**: Standardize time granularity (daily/weekly).
- **Merge & Sync**: Combine by date.

- **Normalize/Scale**: Use MinMax or Standard Scaler.
- **Target Variables**: Define what to predict (e.g., price(t+1), volatility(t+1), risk_score(t+1), Buy/Sell/Hold).

---

# Phase 4: Feature Engineering

## Features:

- **Lag Features**: Past values (e.g., Close(t-1), RSI(t-1)).
- **Rolling Stats**: MA, Std Dev, ATR.
- **Technical Indicators**: RSI, MACD, Bollinger Bands.
- **Fundamental Ratios**: Growth in EPS, ROE, etc.
- **Macro Inputs**: CPI, Interest, Oil, VIX (can be lagged).
- **Calendar Features**: Day of week, earnings proximity.
- **Cross-Asset**: S&P500, Oil prices.
- **Sentiment (Optional)**: Text-based scoring.

## Factor Types:

- **Single-Factor**: Uses only past price.
- **Multi-Factor**: Combines technical, fundamental, macro.

---

# Phase 5: Modeling (Forecasting)

## Targets:

- **Regression**: Future price, volatility, risk_score.
- **Classification**: Buy / Hold / Sell labels.

## Models:

**Traditional:**

- **ARIMA**: Univariate, linear.
- **VAR**: Multi-variate time series.

**ML:**

- **Linear Regression**: Baseline.
- **Random Forest/XGBoost/LightGBM**: Non-linear tabular models.

**DL:**

- **LSTM/GRU**: Sequence learning.
- **BiLSTM/CNN-LSTM**: Capture local + long dependencies.
- **Transformers (TFT/Helformer)**: Advanced sequence models.

**Hybrid:**

- **Prophet**: Time series with trend/seasonality + regressors.
- **LSTM + Macro**: Combines deep sequence + tabular.
- **Reinforcement Learning**: Decision-focused modeling.

---

# Phase 6: Evaluation & Performance Tracking

## Regression:

- **MAE**, **RMSE**, **MAPE**, **R² Score**

## Classification:

- **Accuracy**, **Precision**, **Recall**, **F1**, **Confusion Matrix**

## Financial:

- **Sharpe Ratio**: Return / Volatility
- **Max Drawdown**: Worst loss
- **Hit Rate**: Directional accuracy
- **Calmar Ratio**: Return / Drawdown

  Compare models across tasks/stocks with dashboards.

---

# Phase 7: Strategy Engine (Trading Logic)

## Generate Signals:

- Based on predicted returns, volatility, and risk.

## Logic:

- **Thresholds**: e.g., Return > 2% and Volatility < 1.5% → Buy
- **Classification**: Predict label directly
- **Risk-Adjusted Ranking**:

score = expected_return / expected_volatility

- **Rank** stocks and select top performers.

---

# Phase 8: Visualization, UI & Deployment

- **Dashboard**: Streamlit, Dash
- **Charts**: Plotly candlestick, overlay forecasts
- **Alerts**: Email, Telegram when thresholds hit
- **Backtesting**: Strategy simulator with historical data
- **Scheduler**: Automate data updates, model retraining
- **Optional API**: Serve forecasts or decisions via endpoint

---

# Final Notes:

This guide delivers a detailed framework for designing a fully-fledged, intelligent, multi-stock forecasting and strategy system. It connects financial theory, data science, and real-world trading logic in a coherent pipeline for maximum interpretability, extensibility, and performance.

# Comprehensive Guide to Stock Market Forecasting: Concepts, Data, Modeling, and Strategy

This guide outlines a phased approach to building a sophisticated stock market forecasting system. It integrates market basics, data handling, feature engineering, advanced modeling techniques (Machine Learning/Deep Learning), evaluation, and strategic decision-making, incorporating technical, fundamental, and macroeconomic perspectives. The goal is to predict future stock prices, volatility, and risk, translating these forecasts into actionable Buy/Sell/Hold decisions for multiple companies.

---

## 🔁 End-to-End Project Flowchart (Conceptual Overview)

[Raw Data Sources] → [Data Collection Scripts] → [Data Cleaning/Alignment] → [Feature Engineering] → [Modeling (ML/DL)] → [Evaluation] → [Trading Strategy Logic] → [Dashboard/Deployment/Backtesting]

---

## 🔀 ML Pipeline Diagram

Raw Data → Preprocessing → Feature Engineering → Model Training → Model Validation → Signal Generation → Strategy Logic → Alerts/UI/API

---

## Phase 1: Market Education & Key Concepts

**Indicator Formulas (Select Examples):**

- **Simple Moving Average (SMA)**: $SMA(t) = (P_1 + P_2 + ... + P_n) / n$
- **Exponential Moving Average (EMA)**: $EMA(t) = Price(t) \times \alpha + EMA(t-1) \times (1 - \alpha)$, where $\alpha = 2 / (n + 1)$
- **RSI** = $100 - (100 / (1 + RS))$, where RS = Avg. Gain / Avg. Loss
- **MACD** = $EMA(12) - EMA(26)$; Signal Line = EMA(9) of MACD
- **Bollinger Bands** = $SMA \pm 2 \times$ standard deviation
- **Volatility ($\sigma$)** = $\sqrt{(\Sigma(P_t - \mu)^2 / N)}$

**Evaluation Metric Formulas:**

- **MAE** = $\Sigma|y_i - \hat{y}_i| / n$
- **RMSE** = $\sqrt{(\Sigma(y_i - \hat{y}_i)^2 / n)}$

- **R² Score** = 1 - Σ(yi - ŷi)² / Σ(yi - ȳ)²
- **Sharpe Ratio** = (Rp - Rf) / σp
- **Max Drawdown** = (Peak - Trough) / Peak

---

# Phase 2: Data Collection (Multi-Source)

## Sources & Tools:

- **Yahoo Finance (yfinance)**:

```python
import yfinance as yf
df = yf.download("AAPL", start="2020-01-01", end="2023-12-31")
```

- **FRED (fredapi)**:

```python
from fredapi import Fred
fred = Fred(api_key='YOUR_API_KEY')
cpi = fred.get_series('CPIAUCSL')
```

- **Alpha Vantage**: Exchange rates, economic data (requires API key).
- **FinancialModelingPrep**: Fundamentals, earnings, financials.

## Mega Dataset Merge (Example):

```python
import pandas as pd
# Merge stock + macro by date
df_combined = stock_df.merge(cpi_df, on='Date', how='left')
df_combined = df_combined.fillna(method='ffill')
```

---

# Phase 3: Preprocessing & Alignment

Tasks:

- Remove duplicates, forward-fill macro data
- Normalize: `from sklearn.preprocessing import MinMaxScaler`
- Align all datasets by date
- Create targets: `price(t+1), volatility(t+1), label(t+1)`

---

# Phase 4: Feature Engineering

- Lag Features: Close(t-1), RSI(t-1), etc.
- Rolling Features: SMA, EMA, ATR

- Fundamental Growth: EPS(t) - EPS(t-1)
- Sentiment (optional)
- Macro (CPI, VIX, etc.)
- Cross-sector proxies (e.g., Oil for airlines)

---

# Phase 5: Modeling (Time-Series Forecasting)

## 📘 Model Matrix

| Model | Inputs | Output(s) | Factors | Target Type | Use Case | Application |
|---|---|---|---|---|---|---|
| ARIMA | Price | Price(t+1) | Single | Regression | Baseline Forecast | Price Trend |
| VAR | Price, Macro | Price(t+1) | Multi | Regression | Macro-Aware | Index Forecast |
| RF/XGB/LGBM | All tabular features | Price, Label | Multi | Both | Explainable | Risk Analysis |
| LSTM | Sequences (n x features) | Price | Multi | Regression | Sequence Trend | Volatility/Return |
| BiLSTM/CNN-LSTM | Enhanced LSTM | Volatility | Multi | Regression | Pattern-Based | Volatility Risk |
| Transformer (TFT) | Full multi-series | All Targets | Multi | All | SOTA Long Forecast | Trading Agent |
| Prophet + XGB | Price + macro | Price | Multi | Regression | Seasonal + Ext | Mid-Term Forecast |

---

# Phase 6: Evaluation & Tracking

Track metrics per model, per stock, per horizon:

- MAE, RMSE, MAPE, $R^2$
- Sharpe Ratio, Max Drawdown, Calmar Ratio
- Classification: Confusion Matrix, Precision, Recall
- Visualization: Heatmap or radar chart per model

---

# Phase 7: Strategy Engine (Buy/Sell/Hold Logic)

```
score = expected_return / expected_volatility
if score > 2:
    action = 'Buy'
elif score < -1:
    action = 'Sell'
else:
    action = 'Hold'
```

- Filter based on liquidity or macro sentiment
- Rank stocks using model forecasts and score thresholds

---

# Phase 8: Visualization & Deployment

- Dashboards: Streamlit or Dash
- Charts: Plotly (candlestick + overlays)
- Alerts: Telegram Bot or Email API
- Scheduler: `cron`, `Airflow`, or `Prefect`
- Deployment: Host model via FastAPI/Flask

---

This updated guide provides you with full infrastructure, formulas, modeling logic, evaluation tracking, and practical code examples, plus integration strategies for building a powerful, data-rich forecasting and strategy system across multiple stocks and market indicators.

---------------------------------------------------------------------------------------------------------------

# 📥 Data Gathering: Clear Explanation

### 🎯 Objective

Collect relevant historical data from different sources for multiple companies and macroeconomic indicators to build a unified dataset for modeling.

### 🧩 Data Types and Why They Matter:

| Type | Examples | Purpose |
| --- | --- | --- |
| Stock Price Data | Open, High, Low, Close, Volume (OHLCV) | Predict future prices and calculate returns, volatility, trends |

| Technical Indicators | RSI, MACD, Bollinger Bands | Input features for model, identify patterns |
| --- | --- | --- |
| Fundamentals | EPS, P/E, ROE, Dividend Yield | Signal financial health and valuation |
| Macroeconomic | CPI, Fed Rate, Oil, Gold, VIX | Affect entire market; used as macro context |
| Calendar Events | Earnings dates | Help explain anomalies or volatility spikes |

# 🧰 How to Collect the Data

### 1. Using yfinance for OHLCV:

python
CopyEdit

```python
import yfinance as yf
tickers = ['AAPL', 'TSLA', 'GOOGL']
data = {ticker: yf.download(ticker, start='2018-01-01',
end='2023-12-31') for ticker in tickers}
```

### 2. Using FRED API for Macroeconomic Indicators:

python
CopyEdit

```python
from fredapi import Fred
fred = Fred(api_key='YOUR_API_KEY')
cpi = fred.get_series('CPIAUCSL')  # Consumer Price Index
rate = fred.get_series('FEDFUNDS')  # Federal Funds Rate
```

### 3. Using FinancialModelingPrep:

API endpoint for EPS, P/E, etc.:

bash
CopyEdit

```bash
https://financialmodelingprep.com/api/v3/income-statement/AAPL?apikey=YOUR_API_KEY
```

# 🔗 Combining Datasets by Date

## ✨ Why Merge?

All time series must be aligned to the same calendar to ensure meaningful comparisons and to match inputs with outputs in supervised learning.

## 🛠️ Merge Code Example:

```python
CopyEdit
import pandas as pd

# Stock price (daily)
df_stock = data['AAPL'].reset_index()[['Date', 'Close', 'Volume']]
df_stock.rename(columns={'Close': 'AAPL_Close'}, inplace=True)

# CPI monthly (convert to daily)
cpi_df = cpi.reset_index()
cpi_df.columns = ['Date', 'CPI']
cpi_df['Date'] = pd.to_datetime(cpi_df['Date'])
cpi_df =
cpi_df.set_index('Date').resample('D').ffill().reset_index()

# Merge
df_merged = pd.merge(df_stock, cpi_df, on='Date', how='left')
```

## 🧼 Post-Merge Cleanup:

- **Forward fill** macro data
- Drop NA rows (or interpolate)
- Ensure `Date` is `datetime64` type
- Set `Date` as index for time series models

---------------------------------------------------------------------------------------------------------------

# Comprehensive Guide to Stock Market Forecasting: Concepts, Data, Modeling, and Strategy

This guide outlines a phased approach to building a sophisticated stock market forecasting system. It integrates market basics, data handling, feature engineering, advanced modeling techniques (Machine Learning/Deep Learning), evaluation, and strategic decision-making, incorporating technical, fundamental, and macroeconomic perspectives. The goal is to predict future stock prices, volatility, and risk, translating these forecasts into actionable Buy/Sell/Hold decisions for multiple companies.

---

## 🔁 End-to-End Project Flowchart (Conceptual Overview)

[Raw Data Sources] → [Data Collection Scripts] → [Data Cleaning/Alignment] → [Feature Engineering] → [Modeling (ML/DL)] → [Evaluation] → [Trading Strategy Logic] → [Dashboard/Deployment/Backtesting]

---

## 🔀 ML Pipeline Diagram

Raw Data → Preprocessing → Feature Engineering → Model Training → Model Validation → Signal Generation → Strategy Logic → Alerts/UI/API

---

## Phase 1: Market Education & Key Concepts

**Indicator Formulas (Select Examples):**

- **Simple Moving Average (SMA)**: $SMA(t) = (P_1 + P_2 + ... + P_n) / n$
- **Exponential Moving Average (EMA)**: $EMA(t) = Price(t) \times \alpha + EMA(t-1) \times (1 - \alpha)$, where $\alpha = 2 / (n + 1)$
- **RSI** = $100 - (100 / (1 + RS))$, where RS = Avg. Gain / Avg. Loss
- **MACD** = $EMA(12) - EMA(26)$; Signal Line = $EMA(9)$ of MACD
- **Bollinger Bands** = $SMA \pm 2 \times$ standard deviation
- **Volatility (σ)** = $\sqrt{\Sigma(P_t - \mu)^2 / N}$

**Evaluation Metric Formulas:**

- **MAE** = $\Sigma|y_i - \hat{y}_i| / n$
- **RMSE** = $\sqrt{\Sigma(y_i - \hat{y}_i)^2 / n}$

- **R² Score** = $1 - \Sigma(y_i - \hat{y}_i)^2 / \Sigma(y_i - \bar{y})^2$
- **Sharpe Ratio** = $(R_p - R_f) / \sigma_p$
- **Max Drawdown** = (Peak - Trough) / Peak

---

# Phase 2: Data Collection (Multi-Source)

## 🎯 Objective

Collect high-quality historical data from multiple sources for a range of companies and macroeconomic indicators. This serves as the foundational dataset for the rest of the modeling pipeline.

## 🧩 Data Categories & Purpose:

| Type | Examples | Purpose |
|---|---|---|
| Stock Price (OHLCV) | Open, High, Low, Close, Volume | Used to compute returns, trends, and for forecasting prices |
| Technical Indicators | RSI, MACD, Bollinger Bands | Derived features to detect momentum and volatility |
| Fundamentals | EPS, P/E, ROE, Dividends, Debt | Signals of company health and valuation |
| Macroeconomic | CPI, Interest Rate, Oil, Gold, VIX | External factors that affect overall market conditions |
| Events | Earnings calendar, ex-dividend dates | Timing impact on volatility and price reactions |

## 🔍 Data Sources & Tools:

- **Yahoo Finance (`yfinance`)**: Historical OHLCV prices

```
import yfinance as yf
tickers = ['AAPL', 'MSFT', 'TSLA']
data = {ticker: yf.download(ticker, start='2018-01-01', end='2023-12-31') for ticker in tickers}
```

- **FRED API (`fredapi`)**: CPI, interest rates, unemployment

```
from fredapi import Fred
fred = Fred(api_key='YOUR_API_KEY')
cpi = fred.get_series('CPIAUCSL')
rate = fred.get_series('FEDFUNDS')
```

- **FinancialModelingPrep / Alpha Vantage**:

- ○ Income statements:
  `https://financialmodelingprep.com/api/v3/income-statement/AAPL`
- ○ Fundamental metrics (EPS, P/E, etc.)

## 🛠️ Best Practices for Data Collection:

- Store raw API output in CSV/Parquet format
- Document update frequency (e.g., daily, monthly)
- Track source URLs or API calls
- Log timestamps for last refresh

---

## Sources & Tools:

- **Yahoo Finance (yfinance)**:

```
import yfinance as yf
df = yf.download("AAPL", start="2020-01-01", end="2023-12-31")
```

- **FRED (fredapi)**:

```
from fredapi import Fred
fred = Fred(api_key='YOUR_API_KEY')
cpi = fred.get_series('CPIAUCSL')
```

- **Alpha Vantage**: Exchange rates, economic data (requires API key).
- **FinancialModelingPrep**: Fundamentals, earnings, financials.

## Mega Dataset Merge (Example):

```
import pandas as pd
# Merge stock + macro by date
df_combined = stock_df.merge(cpi_df, on='Date', how='left')
df_combined = df_combined.fillna(method='ffill')
```

---

# Phase 3: Preprocessing & Alignment

## 🎯 Objective

Prepare the raw data for modeling by cleaning, aligning, synchronizing, and creating a unified dataset indexed by time.

## 🛠️ Tasks Overview:

- **Clean Missing Data**:

- ○ Forward-fill macroeconomic time series (e.g., CPI, interest)
- ○ Drop or interpolate gaps in stock data
- **Align Frequencies**:
  - ○ Resample macro data (e.g., CPI monthly → daily):

```
cpi_df = cpi.reset_index()
cpi_df.columns = ['Date', 'CPI']
cpi_df = cpi_df.set_index('Date').resample('D').ffill().reset_index()
```

- **Synchronize Timelines**:
  - ○ Merge datasets on `Date` using left joins
  - ○ Align time zones and formats
- **Merge Datasets Example**:

```
import pandas as pd
# Combine stock and macro data
df_stock = data['AAPL'].reset_index()[['Date', 'Close']].rename(columns={'Close': 'AAPL_Close'})
df_merged = pd.merge(df_stock, cpi_df, on='Date', how='left')
df_merged = df_merged.fillna(method='ffill')
```

- **Normalize/Scale Features**:

```
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
df_scaled = pd.DataFrame(scaler.fit_transform(df_merged.iloc[:, 1:]),
columns=df_merged.columns[1:])
df_scaled['Date'] = df_merged['Date']
```

- **Create Target Variables**:

```
df_scaled['Target_Price'] = df_scaled['AAPL_Close'].shift(-1)  # Predict tomorrow's price
df_scaled['Target_Label'] = (df_scaled['Target_Price'] > df_scaled['AAPL_Close']).astype(int)
# 1 = Buy, 0 = Hold/Sell
```

## ✅ Output of This Phase

- A clean, synchronized, scaled, and target-enhanced DataFrame
- All variables indexed by `Date` and ready for feature engineering

Tasks:

- Remove duplicates, forward-fill macro data
- Normalize: `from sklearn.preprocessing import MinMaxScaler`
- Align all datasets by date
- Create targets: `price(t+1)`, `volatility(t+1)`, `label(t+1)`

# Phase 4: Feature Engineering

- Lag Features: Close(t-1), RSI(t-1), etc.
- Rolling Features: SMA, EMA, ATR
- Fundamental Growth: EPS(t) - EPS(t-1)
- Sentiment (optional)
- Macro (CPI, VIX, etc.)
- Cross-sector proxies (e.g., Oil for airlines)

---

# Phase 5: Modeling (Time-Series Forecasting)

### 📘 Model Matrix

| Model | Inputs | Output(s) | Factors | Target Type | Use Case | Application |
|---|---|---|---|---|---|---|
| ARIMA | Price | Price(t+1) | Single | Regression | Baseline Forecast | Price Trend |
| VAR | Price, Macro | Price(t+1) | Multi | Regression | Macro-Aware | Index Forecast |
| RF/XGB/LGBM | All tabular features | Price, Label | Multi | Both | Explainable | Risk Analysis |
| LSTM | Sequences (n x features) | Price | Multi | Regression | Sequence Trend | Volatility/Return |
| BiLSTM/CNN-LSTM | Enhanced LSTM | Volatility | Multi | Regression | Pattern-Based | Volatility Risk |
| Transformer (TFT) | Full multi-series | All Targets | Multi | All | SOTA Long Forecast | Trading Agent |
| Prophet + XGB | Price + macro | Price | Multi | Regression | Seasonal + Ext | Mid-Term Forecast |

---

# Phase 6: Evaluation & Tracking

Track metrics per model, per stock, per horizon:

- MAE, RMSE, MAPE, R²
- Sharpe Ratio, Max Drawdown, Calmar Ratio

- Classification: Confusion Matrix, Precision, Recall
- Visualization: Heatmap or radar chart per model

---

# Phase 7: Strategy Engine (Buy/Sell/Hold Logic)

```
score = expected_return / expected_volatility
if score > 2:
    action = 'Buy'
elif score < -1:
    action = 'Sell'
else:
    action = 'Hold'
```

- Filter based on liquidity or macro sentiment
- Rank stocks using model forecasts and score thresholds

---

# Phase 8: Visualization & Deployment

- Dashboards: Streamlit or Dash
- Charts: Plotly (candlestick + overlays)
- Alerts: Telegram Bot or Email API
- Scheduler: `cron`, `Airflow`, or `Prefect`
- Deployment: Host model via FastAPI/Flask

---

This updated guide provides you with full infrastructure, formulas, modeling logic, evaluation tracking, and practical code examples, plus integration strategies for building a powerful, data-rich forecasting and strategy system across multiple stocks and market indicators.