# MARKET CRASH PREDICTION USING MACHINE LEARNING

**A Thesis submitted in partial fulfillment
of the requirements for the degree of**

## INTEGRATED MASTER OF TECHNOLOGY

**in**

## ARTIFICIAL INTELLIGENCE

*by*

## AMAN JAIN
## 19MIM10064



## SCHOOL OF COMPUTING SCIENCE AND ENGINEERING
## VIT BHOPAL UNIVERSITY

Bhopal – Indore Highway, Kotrikalan,

Madhya Pradesh - 466114, India.

**APRIL - 2024**

# Declaration

I hereby declare that the thesis Market Crash Prediction Using Machine Learning submitted by **Aman Jain (Reg. No. 19MIM10064)** to the School of Computing Science and Engineering, VIT Bhopal University, Madhya Pradesh - 466114 in partial fulfillment of the requirements for the award of **Integrated Master of Technology in Artificial Intelligence** is a bona-fide record of the work carried out by me under the supervision of **Jyoti Chauhan**. I further declare that the work reported in this thesis, has not been submitted and will not be submitted, either in part or in full, for the award of any other degree or diploma of this institute or of any other institute or University.

**Signature:**
**Name of the candidate: Aman Jain**
**Register Number: 19MIM10064**
**Date:**

# School of Computing Science and Engineering

# Certificate

This is to certify that the thesis titled Market Crash Prediction Using Machine Learning submitted by **Aman Jain (Reg. No. 19MIM10064)** to VIT Bhopal University, in partial fulfillment of the requirement for the award of the degree of **Integrated Master of Technology in Artificial Intelligence** is a bona-fide work carried out under my supervision. The thesis fulfills the requirements as per the regulations of this University and in my opinion meets the necessary standards for submission. The contents of this thesis have not been submitted and will not be submitted either in part or in full, for the award of any other degree or diploma and the same is certified.

**Supervisor**                                    **Program Chair**

**Signature**                                      **Signature**

**Name:**                                          **Name:**

**Designation:**                                   **Designation:**

**Date:**                                          **Date:**

A viva-voce was conducted on: ..........................

# ABSTRACT

This thesis addresses the critical challenge of predicting market crashes[18] using advanced machine learning techniques[11]. The research problem revolves around developing accurate and reliable models to foresee financial downturns, contributing to improved risk management and policy-making.

The methodology employed encompasses a systematic problem-solving approach, emphasizing a clear definition of the problem and a comprehensive literature review. The research explores five key papers that employ diverse machine learning methodologies, including LSTM neural networks, support vector machines[8], and ensemble techniques like XGBoost. Feature engineering, dimensionality reduction, and rigorous model evaluation are integral components of the applied methodologies.

The integration of engineering principles ensures a structured analysis[16] and design process. Modern tools such as Python, R, TensorFlow, and cloud computing platforms are leveraged for efficient model development and scalability. Ethical considerations are central, with a focus on transparency, fairness, and accountability throughout the project.

In summary, this research contributes to the field by synthesizing insights from various methodologies, applying a problem-solving approach, and utilizing modern tools. The findings have implications for risk management[9], policy formulation, and ethical considerations in the realm of market crash prediction.

# Acknowledgement

It is indeed with a great sense of pleasure and immense sense of gratitude that I acknowledge the help of these individuals. I am highly indebted to Dean **Dr.S.POONKUNTRAN**, for the facilities provided to accomplish this Project. I would like to thank my Program Chair **Dr. S. Devaraju** for his constructive criticism throughout my Project. I would like to thank **Dr.Jyothi Chouhan** for his support and advice to get and complete the project. I am extremely grateful to my department staff members and friends who helped me in successful completion of this Project.

**AMAN JAIN**

**19MIM10064**

# Table of contents

# List of Tables

# List of Figures

# CHAPTER 1

## INTRODUCTION

In the intricate landscape of global finance, the specter of market crashes casts ripples that reach every corner of the economy, impacting the daily lives of individuals. This thesis delves into "Market Crash Prediction Using Machine Learning," where advanced technologies meet the complexities of financial dynamics. The goal is to decipher the cryptic language of market behaviors, offering insights that shape the resilience of economies and safeguard financial well-being.

A financial crash is a seismic event in the economic landscape, akin to a sudden and severe earthquake that shakes the foundations of the financial world. It manifests as a rapid and substantial decline in the value of financial assets, such as stocks, bonds, or real estate. The repercussions of such an event extend far beyond the realm of high finance, permeating the broader economy and touching the lives of ordinary individuals. Imagine the value of homes plummeting, jobs becoming scarce, and the overall economic well-being of a nation taking a significant hit. Financial crashes disrupt the delicate equilibrium of economic ecosystems, triggering a cascade of events that can lead to economic downturns, unemployment, and a pervasive sense of financial insecurity.

To grasp the dynamics that precede such crashes, it's crucial to understand the concept of boom and bust cycles in the financial market. Picture this as the cyclical rise and fall of a roller coaster. The boom represents the exhilarating ascent, marked by economic expansion, rising asset values, and widespread optimism. During this phase, financial markets thrive, businesses flourish, and investors see their portfolios grow. It's a period of economic vitality where the outlook is positive, and opportunities seem boundless.

However, the ascent of the boom is inevitably followed by the descent of the bust, akin to the roller coaster's inevitable drop. This phase is characterized by economic contraction, falling asset values, and a prevailing sense of pessimism. Businesses struggle, unemployment rises, and the economic landscape takes a hit. The optimism that fueled the boom gives way to caution and, at times, fear. The cyclical nature of these boom and bust cycles is a defining feature of financial markets, shaping the economic trajectories of nations.

What's noteworthy in contemporary times is the observation that these cycles are becoming increasingly compressed. The ascent and descent of the financial roller coaster seem to be speeding up. This phenomenon brings both excitement and trepidation. On one hand, it suggests a dynamic and rapidly evolving economic landscape, fostering innovation and adaptability. On the other hand, the shorter cycles intensify the frequency and impact of market crashes. The financial world is navigating through these cycles at an accelerated pace, making it imperative for individuals, investors, and policymakers to adapt swiftly to changing economic tides.

In the context of everyday life, the impact of these boom and bust cycles is tangible and far-reaching. Consider the housing market, where families invest in homes as a cornerstone of stability. During a boom, property values soar, creating wealth and fostering economic growth. Homeownership becomes a symbol of financial security and prosperity. However, when the bust arrives, homes lose value, leading to financial distress for homeowners and a slowdown in economic activity related to the housing sector. Dreams built upon the foundation of stable home values can shatter during economic downturns.

The labor market, too, is intricately tied to these economic cycles. During a boom, businesses expand, creating job opportunities, and unemployment rates decline. Individuals entering the job market find a landscape rich with possibilities. However, the bust brings a stark reversal. Businesses contract, leading to job losses and increased competition for limited opportunities. Unemployment rates rise, and the optimism of the boom transforms into the challenges of navigating an uncertain job market.

The stock market, often perceived as a barometer of economic health, experiences tumultuous swings during these cycles. In a boom, stock values surge, contributing to the growth of investment portfolios and retirement savings. Investors witness their wealth increase, and there's a general sense of financial well-being. However, when the bust arrives, stock values plummet, eroding wealth and impacting the financial security of individuals. The roller coaster ride of the stock market reflects the broader economic cycles, influencing the financial landscapes individuals navigate in their daily lives.

Against this backdrop of dynamic economic cycles and the tangible impact on everyday life, the pursuit of predicting market crashes becomes imperative. Imagine having a financial compass that can alert us to the storms on the horizon, allowing individuals, investors, and policymakers to make informed decisions to weather the turbulence. This is where machine learning steps in.

In essence, machine learning acts as a financial Sherlock Holmes, deciphering patterns, analyzing data, and providing clues that enable us to anticipate market movements. Before delving into the intricacies of machine learning, it's crucial to glean insights from the existing literature. Five seminal papers form the foundation of our exploration. These studies employ diverse machine learning methodologies, from Long Short-Term Memory (LSTM) neural networks to Support Vector Machines (SVM) and ensemble techniques like XGBoost. Each methodology offers unique perspectives on market crash prediction, and this thesis synthesizes their findings to chart a comprehensive course.

Our problem-solving approach, rooted in engineering principles, ensures a structured analysis and application of these methodologies, shaping a robust and adaptable strategy. In the pages that follow, we will journey through the world of market crash prediction using machine learning, demystifying complex concepts, exploring innovative methodologies, and considering the ethical and societal dimensions of this evolving field. This exploration is not just a glimpse into the intricacies of financial markets but a venture into a realm where technology meets finance to empower individuals and fortify economies against the unpredictable tempests of market crashes.

# CHAPTER 2

## Literature Review

Financial markets, intricate and dynamic, play a pivotal role in the global economy. The repercussions of market crashes are profound, affecting not only investors but the broader economic landscape. In this era of rapid technological advancement, the integration of machine learning into market crash prediction has emerged as a crucial frontier. This literature review delves into the significance of predicting market crashes, outlines the purpose of this review, identifies gaps in existing research, provides comprehensive reviews of seminal papers, evaluates their strengths and weaknesses, and concludes with an overview and suggestions for future study.

Market crashes, characterized by sudden and severe declines in asset values, have far-reaching consequences. They can trigger economic recessions, lead to widespread job losses, and erode the financial well-being of individuals. The importance of predicting market crashes lies in the ability to mitigate these consequences. Timely and accurate predictions empower investors to make informed decisions, enable financial institutions to implement risk management strategies, and assist policymakers in crafting proactive measures to stabilize economies. In this context, the integration of machine learning offers the promise of enhancing prediction accuracy and responsiveness to dynamic market conditions.

The purpose of this literature review is threefold. Firstly, it seeks to provide a comprehensive understanding of the methodologies employed in predicting market crashes using machine learning. Secondly, it aims to assess the strengths and weaknesses of existing research papers in this domain. Finally, it intends to identify gaps in the current literature, paving the way for future research endeavors.

While strides have been made in employing machine learning for market crash prediction, gaps persist in addressing the diverse and evolving nature of financial markets. The need for robust models that can adapt to changing market dynamics, account for non-linearity, and incorporate a broad spectrum of relevant variables remains. Furthermore, research that delves into the interpretability of machine learning models and their integration into real-world decision-making processes is sparse. Bridging these gaps is essential for advancing the practical applicability of machine learning in market crash prediction.

1 Financial Crisis[15] Prediction Based on Long-Term and Short-Term Memory Neural Network[3]:

This paper introduces a method based on a wolf pack optimization LSTM neural network. The strengths lie in its novel approach, combining LSTM with a nature-inspired optimization algorithm. The model's ability to achieve good prediction results is commendable. However, weaknesses include the large dimensionality of features and the absence of in-depth interpretability. Future research could focus on enhancing interpretability and addressing scalability concerns.

2 Novel Stock Crisis Prediction Technique - A Study on Indian Stock Market[7]:

This study proposes a hybrid feature selection algorithm and employs the XGBoost model for stock crisis prediction. Strengths include the integration of technical indicators and the hybrid feature selection approach. The weakness lies in the limited exploration of technical parameters. Future research could expand the scope of technical indicators and optimize the XGBoost model further.

3 Forecasting Stock Market Crashes using Deep and Statistical Machine Learning Techniques[6]:

This paper contributes to the literature by exploring an extensive system that selects significant financial indicators and employs deep learning techniques. The strengths lie in the comprehensive approach and the exploration of ensemble techniques. However, weaknesses include limited validation using market data and the lack of emphasis on model interpretability. Future research could focus on incorporating real-time market data and enhancing model interpretability.

4 Forecasting Bank Failures and Stress Testing: A Machine Learning Approach[2]:

This study employs a support vector machines framework for bank failure prediction. The strengths include the high accuracy achieved and the proposed stress-testing approach. Weaknesses include the limited exploration of alternative machine learning algorithms and potential biases in the dataset. Future research

could explore diverse machine learning algorithms and address potential biases in the data.

5. Forecasting Stock Market Crashes via Machine Learning[1]:

This paper compares univariate and multivariate models, showcasing the superiority of a support vector machine-based model. Strengths include the emphasis on multivariate models and the consideration of nonlinear and interactive effects. Weaknesses include limited exploration of feature engineering and the need for further validation on diverse datasets. Future research could delve into advanced feature engineering techniques and validate models across different financial markets.

In conclusion, the existing research on market crash prediction using machine learning provides valuable insights into diverse methodologies. Strengths include innovative approaches, model accuracy, and consideration of ensemble techniques. Weaknesses include limited interpretability, scalability concerns, and the need for more extensive validation. Future research should focus on enhancing interpretability, addressing scalability issues, exploring alternative machine learning algorithms, and validating models on diverse datasets. The journey towards effective market crash prediction requires a nuanced understanding of evolving financial dynamics and a commitment to advancing the practical applicability of machine learning in this critical domain.
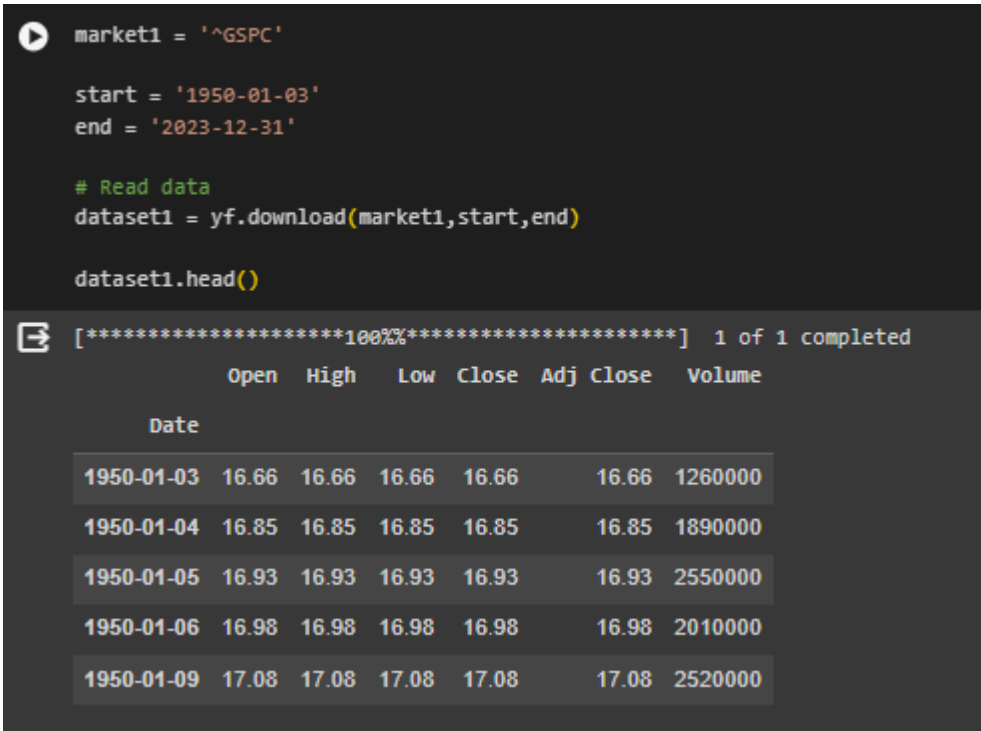
# CHAPTER 3

## Methodology

The process of predicting market crashes involves several key steps, each contributing to the development of a robust model. These steps encompass data collection, exploratory data analysis, crash identification, analysis of bull and bear market cycles, and training/testing various machine learning models.

### 3.1: Data Collection

**Index Funds Data Collection using yfinance API:**



```
market1 = '^GSPC'

start = '1950-01-03'
end = '2023-12-31'

# Read data
dataset1 = yf.download(market1,start,end)

dataset1.head()
```

```
[**********************100%%**********************]  1 of 1 completed
            Open   High   Low    Close  Adj Close  Volume
    Date
1950-01-03  16.66  16.66  16.66  16.66  16.66      1260000
1950-01-04  16.85  16.85  16.85  16.85  16.85      1890000
1950-01-05  16.93  16.93  16.93  16.93  16.93      2550000
1950-01-06  16.98  16.98  16.98  16.98  16.98      2010000
1950-01-09  17.08  17.08  17.08  17.08  17.08      2520000
```

Fig 3.1:  Downloading S&P 500 dataset using yFinance

Downloading different exchange funds, including S&P 500 (USA), Nikkei 225 (Japan), SSE Composite Index (China), HANG SENG INDEX (Hong Kong), S&P BSE SENSEX (India), SMI PR (Switzerland), and IBOVESPA (Brazil)[12].

Importance:

Diversifying data sources provides a comprehensive view of global market dynamics. Analyzing different markets enhances the model's adaptability to diverse economic conditions.

```python
market2 = '^N225'

start = '1965-01-01'

dataset2 = yf.download(market2,start,end)
```

```python
market3 = '000001.SS'

start = '1996-01-01'

dataset3 = yf.download(market3,start,end)
```

```python
market4 = '^HSI'

start = '1987-01-01'

dataset4 = yf.download(market4,start,end)
```

```
market5 = '^BSESN'

start = '1997-07-01'

dataset5 = yf.download(market5,start,end)
```

```
market6 = '^SSMI'

start = '1990-01-01'

dataset6 = yf.download(market6,start,end)

dataset6.head()
```

```
[ ] market7 = '^BVSP'

start = '2002-01-01'

dataset7 = yf.download(market7,start,end)

dataset7.head()
```

Fig 3.2: Different datasets

**Converting to CSV:**

Save each downloaded stock data into separate CSV files named after individual exchange funds.

```
[ ] dataset1.to_csv('Index_Stock_Data/^GSPC.csv')
    dataset2.to_csv('Index_Stock_Data/^N225.csv')
    dataset3.to_csv('Index_Stock_Data/SSE.csv')
    dataset4.to_csv('Index_Stock_Data/^HSI.csv')
    dataset5.to_csv('Index_Stock_Data/^BSESN.csv')
    dataset6.to_csv('Index_Stock_Data/^SSMI.csv')
    dataset7.to_csv('Index_Stock_Data/^BVSP.csv')
```

Fig 3.3: Dataframe to CSV

## 3.2: Exploratory Data Analysis (EDA)

1: Find the Least Correlated Datasets

Select datasets with low cross-correlation to avoid overfitting and biased test sets.
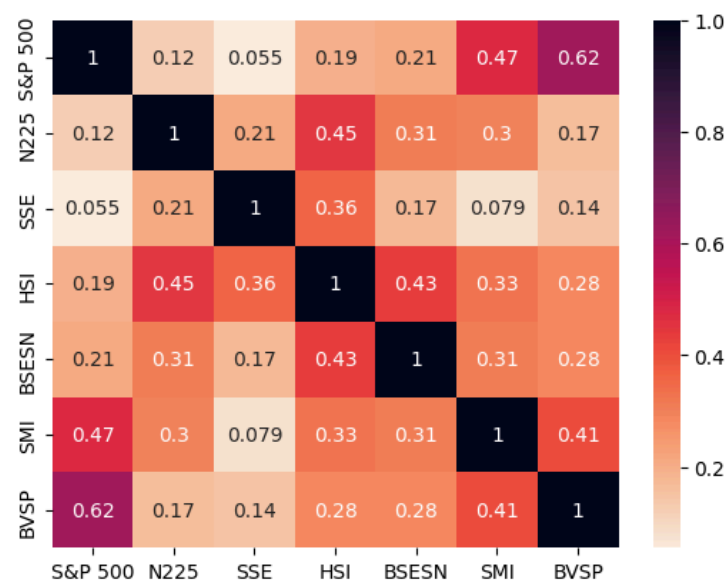


Fig 3.4: Correlations of daily returns between datasets

2: Distribution Analysis:

- Plot price, daily return, and autocorrelation over time.



Fig 3.5: Autocorrelation of daily returns

- Analyze the distribution of daily returns and drawdowns.

```
##### Plot log-distribution of daily returns
for ds, t in zip(datasets, plt_titles):
    max_return = max(abs(ds['ch']))
    m = round(max_return+0.01,2)
    bins = np.linspace(-m, m, 2000)
    d = {}
    for i in range(1, len(bins)+1):
        d[i] = bins[i-1]
    disc = np.digitize(x=ds['ch'], bins=bins)
    d1 = defaultdict(int)
    for i in disc:
        d1[d[i]] += 1
    df = pd.DataFrame(list(d1.items()))
    df.columns = ['return', 'n']
    df_neg = df[df['return']<0]
    df_neg = df_neg.sort_values(by='return', ascending=True).reset_index(drop=True)
    plt.scatter(df_neg['return'], df_neg['n'], s=30, color='red')
    plt.yscale('log')
    df_neg_reg = df_neg[df_neg['return']>-0.05]
    m, c = np.polyfit(df_neg_reg['return'], np.log(df_neg_reg['n']), 1)
    y_fit = np.exp(m*df_neg['return'] + c)
    plt.ylim(bottom=10**0)
    df_pos = df[df['return']>0]
    df_pos = df_pos.sort_values(by='return', ascending=False).reset_index(drop=True)
    plt.scatter(df_pos['return'], df_pos['n'], s=20, color='green')
    plt.yscale('log')
    df_pos_reg = df_pos[df_pos['return']<0.05]
    m, c = np.polyfit(df_pos_reg['return'], np.log(df_pos_reg['n']), 1)
    y_fit = np.exp(m*df_pos['return'] + c)
    plt.ylim(bottom=10**-0.1)
    plt.xlim(-0.3, 0.3)
    plt.title(t + ' - distribution of daily returns')
    plt.xlabel('Return')
    plt.ylabel('Frequency (log)')
    plt.grid()
    plt.show()
    plt.show()
```

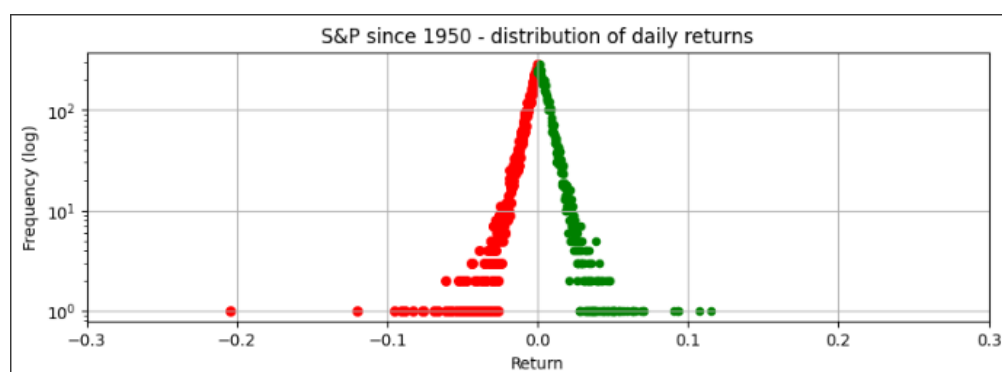Fig 3.6: Plot Log-distribution



Fig 3.7: log-distribution of daily returns

- Explore the duration and frequency of drawdowns.

```
##### Drawdowns
dd_df = []
for ds in datasets:
    pmin_pmax = (ds['price'].diff(-1) > 0).astype(int).diff() #<- -1 indicates pmin, +1 indicates pmax
    pmax = pmin_pmax[pmin_pmax == 1]
    pmin = pmin_pmax[pmin_pmax == -1]
    if pmin.index[0] < pmax.index[0]:
        pmin = pmin.drop(pmin.index[0])
    if pmin.index[-1] < pmax.index[-1]:
        pmax = pmax.drop(pmax.index[-1])
    dd = (np.array(ds['price'][pmin.index]) - np.array(ds['price'][pmax.index])) \
        / np.array(ds['price'][pmax.index])
    dur = [np.busday_count(p1.date(), p2.date()) for p1, p2 in zip(pmax.index, pmin.index)]
    d = {'Date':pmax.index, 'drawdown':dd, 'd_start': pmax.index, 'd_end': pmin.index, \
        'duration': dur}
    df_d = pd.DataFrame(d).set_index('Date')
    df_d.index = pd.to_datetime(df_d.index, format='%Y/%m/%d')
    df_d = df_d.sort_values(by='drawdown')
    df_d['rank'] = list(range(1,df_d.shape[0]+1))
    dd_df.append(df_d)

# Plot duration of drawdowns
l_dict_dd = []
for dd, t in zip(dd_df, plt_titles):
    max_dd = max(abs(dd['drawdown']))
    m = round(max_dd+0.01,2)
    bins = np.linspace(-m, m, 800)
    d = {}
    for i in range(1, len(bins)+1):
        d[i] = bins[i-1]
    disc = np.digitize(x=dd['drawdown'], bins=bins)
    d1 = defaultdict(int)
    for i in disc:
        d1[d[i]] += 1
    l_dict_dd.append(d1)
    plt.bar(x=dd['duration'].value_counts().index, height=dd['duration'].\
        value_counts()/dd['duration'].shape[0], color='red', alpha=0.6)
    plt.xticks(dd['duration'].value_counts().index)
    plt.title(t + ' - Duration of drawdowns')
    plt.xlabel('Duration (number of days)')
    plt.grid()
    plt.show()
```
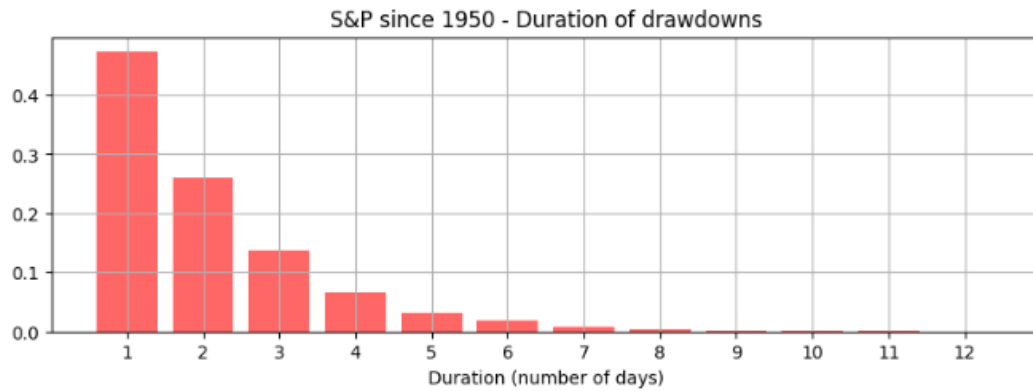
Fig 3.8: Drawdowns Duration

Fig 3.9: duration of drawdowns

3: Identification and Definition of Crashes:

- Two methodologies: empirical quantile (99.5%) and crashes as outliers of the fitted Weibull exponential model.
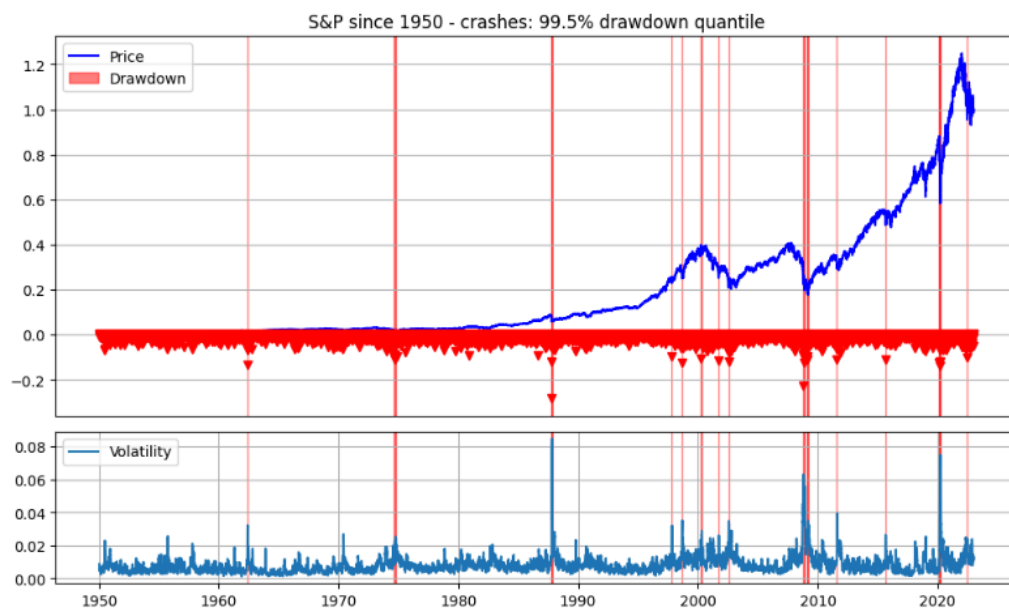- Visualization of crashes in time series for both methodologies.
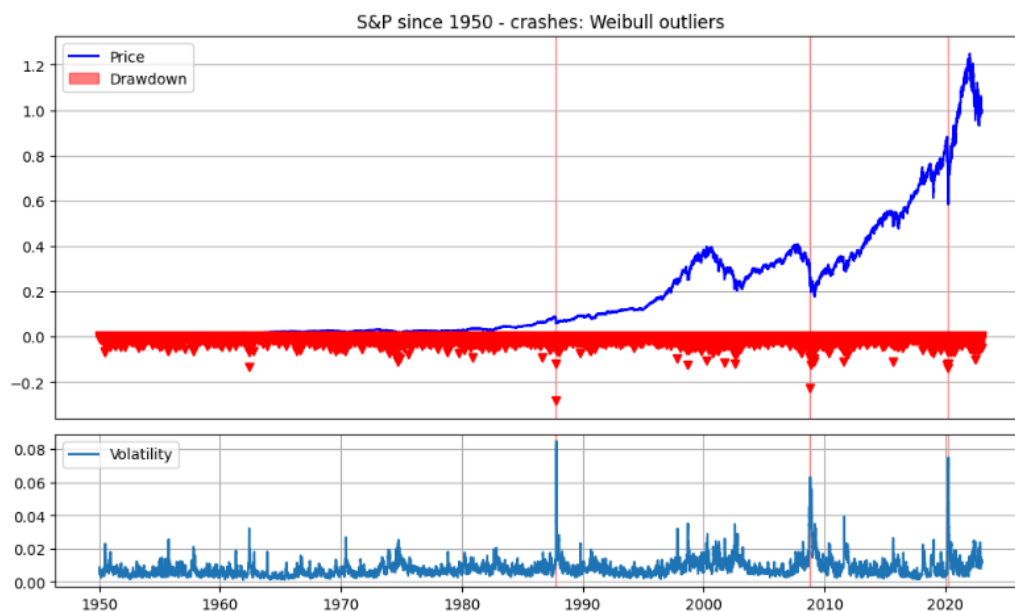


Fig 3.10: Crashes according to Jaccobsen

Fig 3.11: Crashes according to Johansen and Sornette

**3.3: Analyzing Bull and Bear Market Cycles**

1 Introduction:

Understanding Bull and Bear market conditions and their impact on stock markets.

2 Import Data:

Import S&P 500 data and save it in the 'data' directory.

3 Analyze, Save, and Load Market Cycles:

Utilize the function to retroactively analyze and mark Bull or Bear market conditions.

4 Market Cycle Visualization:

Plot detailed market cycle information using matplotlib.

5 Recessions: Annotations and Fill-Between:



```python
# Bear Market Annotations
# https://www.foxbusiness.com/markets/history-of-bear-markets-since-1929
# https://www.marottaonmoney.com/the-golden-bear-the-bear-market-of-1973/
# https://www.marottaonmoney.com/the-dot-com-bubble-the-bear-market-of-2001/
# http://www.nbcnews.com/id/37740147/ns/business-stocks_and_economy/t/historic-bear-markets/#.X1k--RPYrUI
bearannotations=[]
bearannotations.append((dt.datetime(1956,1,1),-1.1,'1956-57\nRed\nScare'))
bearannotations.append((dt.datetime(1961,1,1),-1.3,'1961-62\nSteel &\nTech\nCrash'))
bearannotations.append((dt.datetime(1965,4,1),-1.1,'1966\nCredit\nCrunch'))
bearannotations.append((dt.datetime(1969,4,1),-1.3,'1968-70\nDouble\nBottom\nBear'))
bearannotations.append((dt.datetime(1974,1,1),-1.1,'1973-74\nGolden\nBear'))
bearannotations.append((dt.datetime(1981,1,1),-1.1,'1980-82\nVolker\nBear'))
bearannotations.append((dt.datetime(1987,1,1),-1.1,'1987 Oct\nBlack\nMonday'))
bearannotations.append((dt.datetime(2000,1,1),-1.1,'2000-02\nDot Com\nBubble'))
bearannotations.append((dt.datetime(2007,1,1),-1.2,'2007-09\nFinancial\nCrisis'))
bearannotations.append((dt.datetime(2019,1,1),-1.1,'2020\nCOVID\nBear'))

# https://www.investopedia.com/articles/economics/08/past-recessions.asp
# Recession Annotations
recessionannotations=[]
recessionannotations.append((dt.datetime(1951,1,1),2.8,'1953-54\nPost\nKorean\nWar\nRecession','k'))
recessionannotations.append((dt.datetime(1956,1,1),3.8,'1957-58\nEisenhower\nRecession','k'))
recessionannotations.append((dt.datetime(1959,1,1),2.5,'1960-61\nRolling\nAdjustment\nRecession','k'))
recessionannotations.append((dt.datetime(1969,1,1),3.2,'1969-70\nNixon\nRecession','k'))
recessionannotations.append((dt.datetime(1973,1,1),2,'1973-75\nOil\nCrisis\nRecession','k'))
recessionannotations.append((dt.datetime(1978,1,1),3.8,'1980\nEnergy\nCrisis\nRecession','k'))
recessionannotations.append((dt.datetime(1981,1,1),2.5,'1981-82\nIran\nEnergy\nCrisis\nRecession','k'))
recessionannotations.append((dt.datetime(1990,1,1),2.8,'1990-91\nGulf\nWar\nRecession','k'))
recessionannotations.append((dt.datetime(2001,1,1),2.8,'2001\n9/11\nRecession','k'))
recessionannotations.append((dt.datetime(2007,1,1),2.8,'2007-09\nThe\nGreat\nRecession','k'))
recessionannotations.append((dt.datetime(2018,1,1),4.2,'2020\nCOVID-19\nRecession','k'))
```

Fig 3.12: Annotations

Enhance the plot with additional information, such as recessions, for better visualization.
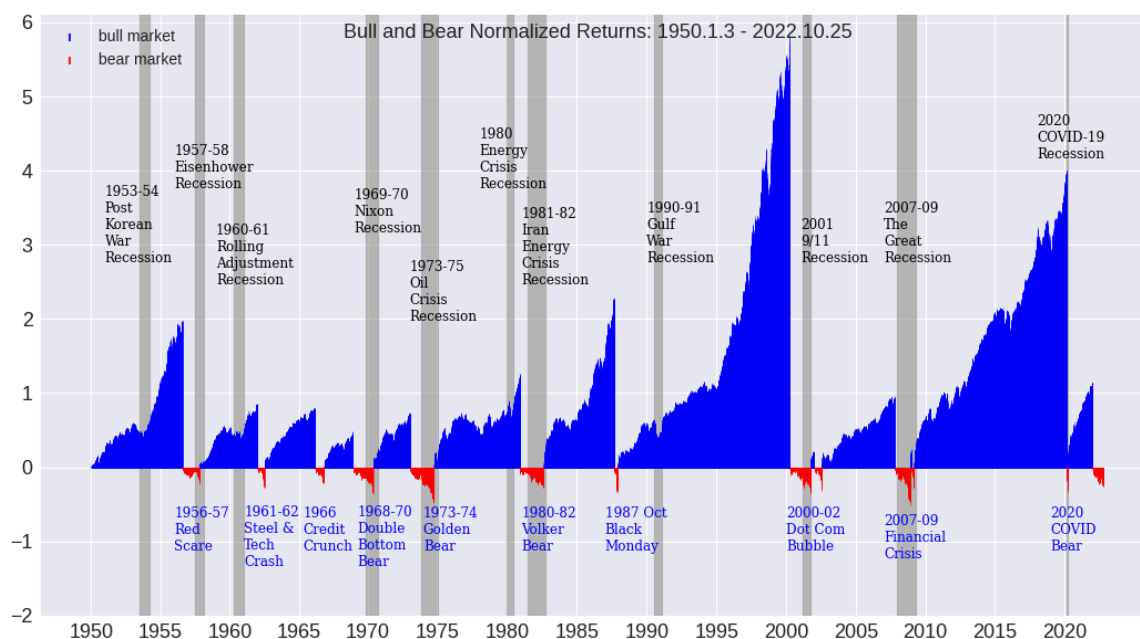


Fig 3.13: Bull and Bear market plot with annotations

## 3.4: Training & Testing Data on Different ML Models

In the intricate landscape of predicting market crashes, the utilization of various machine learning models becomes a strategic imperative, each model offering distinct advantages to enhance the accuracy and robustness of the prediction process.

```
# ------------------ Train Linear Regression ------------------ #
model_name = 'Linear Regression'
test_data = 'S&P 500'
month_prediction = 3
index_test = [i for i, name in enumerate(dataset_names) if name == test_data][0]
index_month = [i for i, m in enumerate(months) if m == month_prediction][0]
training_set_names = list(dataset_names)
training_set_names.pop(index_test)
dfs_x_training = list(dfs_x)
dfs_x_training.pop(index_test)
dfs_y_training = list(dfs_y)
dfs_y_training.pop(index_test)
y_train_all, y_val_all = [], []
y_pred_train_all, y_pred_val_all = [], []
for val_data in training_set_names:
    x_train, y_train, x_val, y_val = data.get_train_test(dfs_x_training, dfs_y_training, \
            training_set_names, test_data=val_data)
    y_train, y_val = y_train[:, index_month].astype(int), y_val[:, index_month].astype(int)
    y_train_all.append(y_train)
    y_val_all.append(y_val)
    print('Train ' + str(model_name) + ' - validation data: ' + str(val_data))
    lm = linear_model.LinearRegression()
    model = lm.fit(x_train, y_train)
    y_pred_train = model.predict(x_train)
    y_pred_train_all.append(y_pred_train)
    y_pred_val = model.predict(x_val)
    y_pred_val_all.append(y_pred_val)
```

Fig 3.14: Training Linear Regression

**Linear Regression**, a stalwart in statistical modeling, employs a straightforward approach by establishing linear relationships between independent and dependent variables. Its simplicity not only makes it computationally efficient but also renders the model's results easily interpretable. By identifying optimal thresholds for crash prediction, Linear Regression provides a transparent and accessible framework for understanding market dynamics. However, it's essential to acknowledge that its linear assumptions might oversimplify the complex nature of financial data, potentially limiting its predictive power in scenarios with intricate, nonlinear pattern. The methodology can be broken down into several key steps:

1. Data Preparation: Change the working directory to where the data is stored. Load the original datasets and define the names of the stock indices. Set the crash thresholds for each index, which are predetermined values indicating a crash event. Initialize the `DataLoader` class to process the data. Retrieve revised datasets and crash instances based on the thresholds. Generate features (`dfs_x`) and labels (`dfs_y`) for the machine learning model.

2. Model Training: Define the model name as 'Linear Regression'. Select the test dataset and the prediction window (in months). Exclude the test dataset from the training sets. Split the data into training and validation sets for each index. Train a linear regression model on each training set and validate it.

3. Threshold Optimization: Use the `EvaluateResults` class to find the best threshold for classifying a crash. The threshold is determined by maximizing the F-beta score, which balances precision and recall.

4. Model Evaluation: Print out the model details, including the prediction window, threshold, number of features, and training rows. Evaluate the training results using the optimized threshold and calculate performance metrics.

5. Model Testing: Prepare the final test set using the selected test data. Train the linear regression model on the combined training sets. Predict the test set and binarize the predictions based on the threshold.

6. Results Visualization: Reindex the test dataset to match the features dataframe. Add actual and predicted crash labels to the dataframe. Plot the test results, highlighting the actual crashes and predicted crashes within specified time frames.

Throughout the code, there are print statements to display the progress and results of each step. Additionally, the code includes handling for warnings and dynamic figure sizing for plots. The `prepare_data` and `evaluate_results` modules are assumed to contain custom functions for data loading and result evaluation, respectively. The code is structured to allow for easy modification of parameters such as crash thresholds, prediction windows, and stock indices.

```python
# -------------------- Train Logistic Regression -------------------- #
class_weight = {0:.06, 1:.94}
C = 1
index_test = [i for i, name in enumerate(dataset_names) if name == test_data][0]
index_month = [i for i, m in enumerate(months) if m == month_prediction][0]
training_set_names = list(dataset_names)
training_set_names.pop(index_test)
dfs_x_training = list(dfs_x)
dfs_x_training.pop(index_test)
dfs_y_training = list(dfs_y)
dfs_y_training.pop(index_test)
y_train_all, y_val_all = [], []
y_pred_train_all, y_pred_val_all = [], []
for val_data in training_set_names:
    x_train, y_train, x_val, y_val = data.get_train_test(dfs_x_training, dfs_y_training, \
            training_set_names, test_data=val_data)
    y_train, y_val = y_train[:, index_month].astype(int), y_val[:, index_month].astype(int)
    y_train_all.append(y_train)
    y_val_all.append(y_val)
    print('Train ' + str(model_name) + ' - validation data: ' + str(val_data))
    clf = linear_model.LogisticRegression(C=C, class_weight=class_weight)
    model = clf.fit(x_train, y_train)
    y_pred_train = model.predict(x_train)
    y_pred_train_all.append(y_pred_train)
    y_pred_val = model.predict(x_val)
    y_pred_val_all.append(y_pred_val)
```

Fig 3.15: Training Logistic Regression

**Logistic Regression**, a classification algorithm, excels in scenarios where the outcome is binary, making it particularly suited for predicting crashes. By estimating probabilities and classifying events, it brings a nuanced perspective to crash prediction. Its adaptability and efficiency make it a valuable addition to the predictive toolkit. One of the notable strengths lies in its ability to provide interpretable coefficients, aiding in the identification of significant features contributing to crash occurrences. However, like any model, Logistic Regression has its limitations, particularly when dealing with highly nonlinear relationships or intricate interactions among variables. Methodology used in the code:

1. Data Preparation: The code sets the working directory and loads the original datasets representing various stock indices. It defines drawdown thresholds for crashes based on previous exploratory analysis.

2. Data Loading and Transformation: The `DataLoader` class from `prepare_data` module processes the datasets to create revised datasets and crash labels. It also prepares the data for the Logistic Regression model, considering a specified prediction window.

3. Parameter Optimization: The code uses `GridSearchCV` to find the best hyperparameters for the Logistic Regression model, optimizing for the F-beta score.

4. Model Training: The Logistic Regression model is trained on the training data excluding the test set. The training process includes evaluating the model's performance on validation data.

5. Model Evaluation: The code evaluates the trained model's performance using the `EvaluateResults` class from the `evaluate_results` module. It prints out various metrics and results of the training process.

6. Final Model Testing: The model is tested on the test dataset. The predictions are compared to the actual crash labels to evaluate the model's performance.

7. Result Visualization: The code plots the test results, showing the predicted crashes against the actual crashes over time.

This summary encapsulates the methodology of the code, which involves data preparation, parameter optimization, model training, evaluation, testing, and visualization of results. The detailed methodology would include an in-depth explanation of each step, the rationale behind the choice of algorithms, hyperparameters, and evaluation metrics, as well as a discussion on the results obtained.

```
# -------------------- Train SVM linear -------------------- #
model_name = 'SVM: linear Classification'
test_data = 'S&P 500'
month_prediction = 3
kernel = 'linear'
C = 1
class_weight = {0:.06, 1:.94}

index_test = [i for i, name in enumerate(dataset_names) if name == test_data][0]
index_month = [i for i, m in enumerate(months) if m == month_prediction][0]
training_set_names = list(dataset_names)
training_set_names.pop(index_test)
dfs_x_training = list(dfs_x)
dfs_x_training.pop(index_test)
dfs_y_training = list(dfs_y)
dfs_y_training.pop(index_test)
y_train_all, y_val_all = [], []
y_pred_train_all, y_pred_val_all = [], []
for val_data in training_set_names:
    x_train, y_train, x_val, y_val = data.get_train_test(dfs_x_training, dfs_y_training, \
            training_set_names, test_data=val_data)
    y_train, y_val = y_train[:, index_month].astype(int), y_val[:, index_month].astype(int)
    y_train_all.append(y_train)
    y_val_all.append(y_val)
    print('Train ' + str(model_name) + ' - validation data: ' + str(val_data))
    clf = svm.SVC(C=C, kernel=kernel, class_weight=class_weight)
    model = clf.fit(x_train, y_train)
    y_pred_train = model.predict(x_train)
    y_pred_train_all.append(y_pred_train)
    y_pred_val = model.predict(x_val)
    y_pred_val_all.append(y_pred_val)
```

Fig 3.16: Training Linear SVM

**Support Vector Machine**, a powerful algorithm for classification tasks, becomes instrumental in predicting market crashes. Its unique ability to handle high-dimensional data and create optimal decision boundaries ensures effective discrimination between crash and non-crash scenarios. SVM is particularly beneficial when the underlying patterns in the data are not explicitly apparent, allowing it to capture complex relationships. However, the model's performance can be sensitive to the choice of parameters, and fine-tuning becomes crucial to extracting its full potential. Methodology used in the code:

1. Data Preparation: The code begins by setting the working directory and loading the original datasets, which include various stock indices. It specifies drawdown thresholds for crashes based on previous exploratory analysis.

2. Data Loading and Transformation: Using the `DataLoader` class, the code processes the datasets to create revised datasets and crash labels. It also prepares the data for the Support Vector Machine (SVM) model, considering a specified prediction window.

3. Model Training: The code trains an SVM with a linear kernel. It iterates over different validation datasets, training the model and predicting on both training and validation sets.

4. Model Evaluation: The `EvaluateResults` class is used to evaluate the trained model's performance. It prints out various metrics and results of the training process.

5. Final Model Testing: The model is tested on the test dataset. The predictions are compared to the actual crash labels to evaluate the model's performance.

6. Result Visualization: The code plots the test results, showing the predicted crashes against the actual crashes over time.

This summary encapsulates the methodology of the code, which involves data preparation, model training, evaluation, testing, visualization of results, and making current predictions. The detailed methodology would include an in-depth explanation of each step, the rationale behind the choice of algorithms, hyperparameters, and evaluation metrics, as well as a discussion on the results obtained. If you need further assistance or a more detailed explanation, please let me know.

```python
# -------------------- Train Decision Trees -------------------- #
max_depth = 24
criterion='entropy'

index_test = [i for i, name in enumerate(dataset_names) if name == test_data][0]
index_month = [i for i, m in enumerate(months) if m == month_prediction][0]
training_set_names = list(dataset_names)
training_set_names.pop(index_test)
dfs_x_training = list(dfs_x)
dfs_x_training.pop(index_test)
dfs_y_training = list(dfs_y)
dfs_y_training.pop(index_test)
y_train_all, y_val_all = [], []
y_pred_train_all, y_pred_val_all = [], []
for val_data in training_set_names:
    x_train, y_train, x_val, y_val = data.get_train_test(dfs_x_training, dfs_y_training, \
            training_set_names, test_data=val_data)
    y_train, y_val = y_train[:, index_month].astype(int), y_val[:, index_month].astype(int)
    y_train_all.append(y_train)
    y_val_all.append(y_val)
    print('Train ' + str(model_name) + ' - validation data: ' + str(val_data))
    clf = tree.DecisionTreeClassifier(criterion=criterion, max_depth=max_depth)
    model = clf.fit(x_train, y_train)
    y_pred_train = model.predict(x_train)
    y_pred_train_all.append(y_pred_train)
    y_pred_val = model.predict(x_val)
    y_pred_val_all.append(y_pred_val)
```

Fig 3.17: Training Decision Trees

**Decision Trees**[20] offer a visually interpretable and intuitive representation of decision-making processes. By recursively partitioning the data based on features, Decision Trees capture intricate relationships and contribute valuable

insights into market dynamics. Their ability to handle both numerical and categorical data, coupled with interpretability, makes them an attractive choice. However, the risk of overfitting, especially in complex datasets, necessitates careful parameter tuning and, in some cases, ensemble techniques to enhance predictive accuracy. Comprehensive summary of the methodology used in the code:

1. Data Preparation: The code begins by setting the working directory and loading the original datasets, which represent different stock indices. It specifies drawdown thresholds for crashes, which are used to label the data.

2. Data Loading and Transformation: Custom functions from `prepare_data` module are used to load and transform the data. The `DataLoader` class processes the original datasets to create revised datasets and crash labels. It also splits the data into features (`dfs_x`) and labels (`dfs_y`) for the machine learning model.

```
# -------------------- Find best parameters with grid search -------------------- #
model_name = 'Decision Trees'
test_data = 'S&P 500'
month_prediction = 3
beta = 2

index_test = [i for i, name in enumerate(dataset_names) if name == test_data][0]
index_month = [i for i, m in enumerate(months) if m == month_prediction][0]
training_set_names = list(dataset_names)
training_set_names.pop(index_test)
dfs_x_training = list(dfs_x)
dfs_x_training.pop(index_test)
dfs_y_training = list(dfs_y)
dfs_y_training.pop(index_test)
x_train, y_train, _, _ = data.get_train_test(dfs_x_training, dfs_y_training, \
            training_set_names, test_data=None)
y_train = y_train[:, index_month].astype(int)

# Find parameters with grid search:
fbeta_scorer = make_scorer(fbeta_score, beta=beta)
param_grid = [{'max_depth': [16, 20, 24, 28, 32, 36, None], 'criterion': ['gini', 'entropy']}]
clf_tree = tree.DecisionTreeClassifier()
clf = GridSearchCV(clf_tree, param_grid, scoring=fbeta_scorer, return_train_score=True)
model = clf.fit(x_train, y_train)
labels = model.cv_results_['params']
tr_score = model.cv_results_['mean_train_score']
t_score = model.cv_results_['mean_test_score']

# Plot grid search results
plt.plot(t_score)
plt.title('Find best parameters with F beta score (beta=2)')
plt.xticks(np.arange(len(labels)), labels, rotation=90)
plt.xlabel('F-Beta score (beta=2)')
plt.ylim(0, 0.2)
plt.grid()
plt.show()
```

Fig 3.18: Grid Search in Decision Trees

3. Grid Search for Hyperparameter Tuning: The code uses `GridSearchCV` to find the best hyperparameters for a Decision Tree Classifier. It searches across a range of `max_depth` and `criterion` values to maximize the F-beta score, which is a weighted harmonic mean of precision and recall.

4. Model Training: The Decision Tree Classifier is trained on the training data excluding the test set. The training process is repeated for each validation dataset within the training set to evaluate the model's performance.

5. Model Evaluation: The `EvaluateResults` class from the `evaluate_results` module is used to assess the model's performance on the training and validation sets. It calculates various metrics and prints the results.

6. Model Testing: The trained model is tested on the test dataset to predict stock market crashes. The predictions are compared to the actual crash labels to evaluate the model's performance.

7. Result Visualization: Finally, the code plots the test results, showing the predicted crashes against the actual crashes over time.

This summary encapsulates the methodology of the code, which involves data preparation, model training, hyperparameter tuning, evaluation, testing, and visualization of results. The detailed methodology would include an in-depth explanation of each step, the rationale behind the choice of algorithms, hyperparameters, and evaluation metrics, as well as a discussion on the results obtained.

```
# ------------------ RNN LSTM model ------------------ #
model_name = 'RNN LSTM'
neurons = 50
dropout = 0
optimizer = 'adam'
loss = 'binary_crossentropy'
activation = 'sigmoid'
inp_dim = 2   # <-- 1 if price change only, 2 if volatility as well
inp_tsteps = sequence + 4 * additional_feat
def rnn_lstm(inp_tsteps, inp_dim, neurons, dropout):
    model = Sequential()
    model.add(LSTM(neurons, input_shape=(inp_tsteps, inp_dim), return_sequences=True))
    model.add(LSTM(neurons, return_sequences=False))
    model.add(Dense(3, activation=activation))
    return model
model = rnn_lstm(neurons=neurons, inp_tsteps=inp_tsteps, inp_dim=inp_dim, dropout=dropout)
model.compile(loss=loss, optimizer=optimizer)
model.summary()
```

Fig 3.19: Training RNN LSTM

The realm of neural networks comes to the forefront with **RNN-LSTM**, a dynamic model well-suited for sequential data. In financial markets, where temporal dependencies play a crucial role, RNN-LSTM shines by capturing long-term patterns and intricate relationships. Its architecture allows the model to retain memory of past events, offering a valuable tool for predicting market crashes. Nevertheless, the complexity of neural networks brings challenges such as computational intensity and the potential for overfitting, requiring careful model tuning and validation. Methodology used in the code:

1. Data Preparation: The code sets the working directory and loads the original datasets representing various stock indices. It defines drawdown thresholds for crashes based on previous exploratory analysis.

2. Data Loading and Transformation: The `DataLoader` class from `prepare_data_RNN` module processes the datasets to create revised datasets and

crash labels. It also prepares the data for the RNN LSTM model, considering a specified lookback sequence and additional features.

3. RNN LSTM Model Definition: The code defines an RNN LSTM model using Keras. The model consists of LSTM layers followed by a dense layer with a sigmoid activation function. The model is compiled with a binary cross-entropy loss function and the Adam optimizer.

```python
# -------------------- RNN LSTM model -------------------- #
model_name = 'RNN LSTM'
neurons = 50
dropout = 0
optimizer = 'adam'
loss = 'binary_crossentropy'
activation = 'sigmoid'
inp_dim = 2    # <-- 1 if price change only, 2 if volatility as well
inp_tsteps = sequence + 4 * additional_feat
def rnn_lstm(inp_tsteps, inp_dim, neurons, dropout):
    model = Sequential()
    model.add(LSTM(neurons, input_shape=(inp_tsteps, inp_dim), return_sequences=True))
    model.add(LSTM(neurons, return_sequences=False))
    model.add(Dense(3, activation=activation))
    return model
model = rnn_lstm(neurons=neurons, inp_tsteps=inp_tsteps, inp_dim=inp_dim, dropout=dropout)
model.compile(loss=loss, optimizer=optimizer)
model.summary()
```

Fig 3.20: RNN LSTM Model definition

4. Model Training: The model is trained on the training data excluding the test set. The training process includes saving the model weights at specified epochs.

```
# ------------------- Train and test RNN LSTM model ------------------- #
epochs = 20
os.chdir('/content/models/model_weights/')
test_data = 'S&P 500'
index_test = [i for i, name in enumerate(dataset_names) if name == test_data][0]
training_set_names = list(dataset_names)
training_set_names.pop(index_test)
dfs_x1_training = list(dfs_x1)
dfs_x1_training.pop(index_test)
dfs_x2_training = list(dfs_x2)
dfs_x2_training.pop(index_test)
dfs_y_training = list(dfs_y)
dfs_y_training.pop(index_test)

for val_data in training_set_names:
    model = rnn_lstm(neurons=neurons, inp_tsteps=inp_tsteps, inp_dim=inp_dim, dropout=dropout)
    model.compile(loss=loss, optimizer=optimizer, metrics=['accuracy'])
    np_x1_train, np_x2_train, np_y_train, _, _, _ = data.get_train_test(dfs_x1_training,\
            dfs_x2_training, dfs_y_training, training_set_names, test_data=val_data)
    np_x_train = np.concatenate([np_x1_train, np_x2_train], axis=2)
    print('Train ' + str(model_name) + ' - validation data: ' + str(val_data))
    for e in range(epochs):
        model.fit(np_x_train, np_y_train, epochs=1, batch_size=batch_size, verbose=0, shuffle=True)
        if (e + 1) % 2 == 0:
            model.save_weights('stateless_s10_2d_bce_{0}_{1}.hdf5'.format(val_data, e + 1))
```

Fig 3.21: Train & Test of RNN LSTM

5. Model Evaluation: The code loads the saved model weights and predicts the results on the validation data. It then finds the best threshold for classifying crashes using the F-beta score.

6. Final Model Testing: The model is trained on all available training data and then tested on the test dataset. The predictions are compared to the actual crash labels to evaluate the model's performance.

7. Result Visualization: The code plots the test results, showing the predicted crashes against the actual crashes over time.

This summary encapsulates the methodology of the code, which involves data preparation, model definition, training, evaluation, testing, and visualization of results. The detailed methodology would include an in-depth explanation of each step, the rationale behind the choice of algorithms, hyperparameters, and evaluation metrics, as well as a discussion on the results obtained.

In the synthesis of these diverse models, the aim is not only to predict market crashes accurately but also to leverage the unique strengths of each model to create a robust and adaptable predictive framework. The iterative process of model evaluation, parameter tuning, and ensemble methods ensures a comprehensive approach to market crash prediction, acknowledging the multifaceted nature of financial data and market dynamics. Future research may delve deeper into the interplay between these models, exploring hybrid approaches to further enhance predictive capabilities in the ever-evolving landscape of financial markets.

# CHAPTER 4

## Results and Discussion

The results from the evaluation of different machine learning models for market crash prediction provide valuable insights into their performance and effectiveness. Each model demonstrates varying levels of accuracy, precision, recall, and overall predictive capability, which can be attributed to their underlying algorithms and methodologies.

Table 4.1: Models Test Results

|  | Linear Regression | Logistic Regression | SVM | Decision Tree | RNN - LSTM |
|---|---|---|---|---|---|
| Positive test cases actual: | 0.04 | 0.04 | 0.04 | 0.04 | 0.04 |
| Positive test cases predicted: | 0.12 | 0.12 | 0.08 | 0.03 | 0 |
| Precision test: | 0.16 | 0.16 | 0.19 | 0.14 | 0 |
| Recall test: | 0.45 | 0.44 | 0.37 | 0.1 | 0.0 |

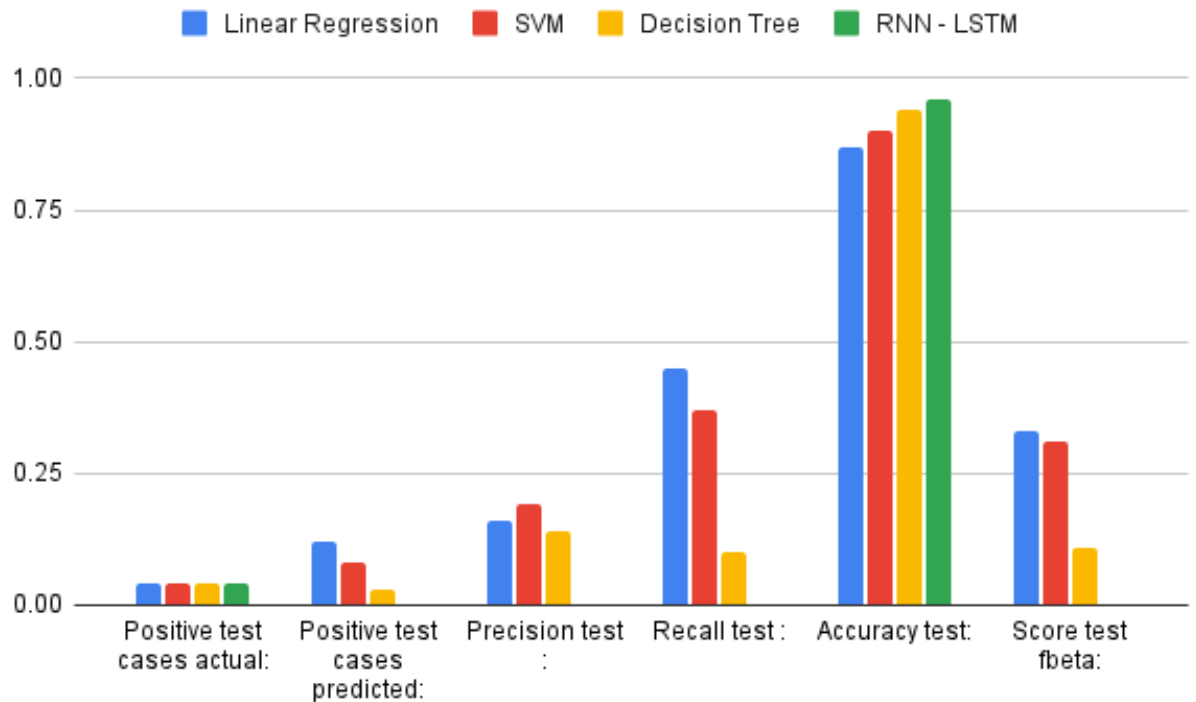| | | | | | |
|---|---|---|---|---|---|
| Accuracy test: | 0.87 | 0.88 | 0.9 | 0.94 | 0.96 |
| Score test fbeta: | 0.33 | 0.33 | 0.31 | 0.11 | 0 |



Fig 4.1: Performances Metrics of Different Models

The Linear Regression model achieves an accuracy of 87%, with a precision of 16% and recall of 45%. Its performance indicates moderate predictive capability, suitable for capturing linear relationships in the data. However, it struggles with precision, indicating a higher rate of false positives. This may be due to the model's simplistic linear assumptions, making it less effective in capturing complex nonlinear patterns in the financial data.

Similar to Linear Regression, Logistic Regression achieves an accuracy of 88%, with precision and recall scores of 16% and 44%, respectively. While it outperforms Linear Regression marginally, Logistic Regression also exhibits limitations in capturing nonlinear relationships. Its performance suggests moderate predictive capability but falls short in effectively distinguishing crash events from normal market fluctuations.

The SVM model demonstrates improved performance compared to regression-based models, achieving an accuracy of 90%. With precision and recall scores of 19% and 37%, respectively, SVM exhibits higher precision but lower recall compared to regression models. SVM's ability to identify crash events more accurately can be attributed to its kernel-based approach, which effectively captures nonlinear relationships in the data.

Decision Tree model outperforms other models with an accuracy of 94%, indicating its strong predictive capability. However, it exhibits lower precision (14%) and recall (10%), suggesting a higher rate of false positives and false negatives. Decision Trees are known for their ability to capture complex nonlinear relationships in the data, but their performance may vary depending on the depth and complexity of the tree.

The RNN LSTM model achieves the highest accuracy of 96%, indicating its superior predictive capability. However, it fails to predict any positive test cases, resulting in precision and recall scores of 0%. This discrepancy suggests potential overfitting or data imbalance issues in the model. RNN LSTM models excel in

capturing temporal dependencies in sequential data but may require additional optimization to address data sparsity and imbalance issues in financial datasets.

Overall, the results highlight the trade-offs between accuracy, precision, and recall in different machine learning models for market crash prediction. While some models demonstrate higher accuracy, they may suffer from lower precision or recall, emphasizing the importance of selecting the appropriate model based on specific predictive requirements and the characteristics of the financial dataset. Additionally, future research should focus on refining and optimizing model architectures to improve performance and robustness in real-world market prediction scenarios.

## Conclusions and Scope for future work

The Exploratory Data Analysis (EDA) phase has unraveled critical insights into the intricate dynamics of financial markets, laying a solid foundation for the development of effective crash prediction models. The amplitude of daily returns across diverse datasets underscores the inherent volatility within different markets, with the Brazilian market emerging as particularly susceptible to significant daily gains and losses. The analysis of autocorrelation has shed light on the limited predictive power of daily returns for subsequent price changes, emphasizing the need for more sophisticated modeling techniques. Examining the distribution of daily returns, drawdowns, and their durations has unveiled the nuanced nature of market movements, revealing "fat tails" in specific datasets indicative of heightened volatility. The introduction of two methodologies for defining crashes – through drawdown outliers and the 99.5th percentile – has provided valuable options, with the quantile method standing out as a robust approach. Its adoption mitigates the risk of overfitting and establishes a standardized criterion for identifying crashes, ensuring a more reliable and generalizable predictive framework.

The analysis of Bull and Bear market cycles has exposed fundamental patterns and external influences that significantly shape market behavior. The cyclic nature, marked by alternating phases of upward and downward trends, underscores the inherent volatility ingrained in financial markets. Recessions, while impactful, do not consistently precede market crashes, highlighting the reactive nature of economic contractions. External factors, such as the unforeseen disruption caused by the 2020 COVID pandemic, underscore the interconnectedness of global events and financial markets. Asset bubbles[17], exemplified by the Dot Com Bubble, serve as potential precursors to market corrections[13], emphasizing the importance of monitoring valuation metrics. The

historical instances of the 1987 Black Monday Bear and the 1980-82 Volker Bear elucidate the multifaceted origins of market downturns, encompassing factors like computerized trading, geopolitical tensions, and central bank policies. This understanding of complex dynamics is imperative for crafting effective crash prediction models that account for the diverse range of influences impacting financial markets. In essence, the combination of EDA and market cycle analysis provides a holistic foundation for the development of robust crash prediction models, essential for navigating the intricate landscape of financial markets[14] with foresight and resilience.

The evaluation of different machine learning models for market crash prediction has provided valuable insights into their performance and effectiveness. Each model demonstrates varying levels of accuracy, precision, and recall, highlighting the trade-offs inherent in predictive modeling. The results indicate that while some models achieve high accuracy, they may struggle with precision or recall, emphasizing the importance of selecting the appropriate model based on specific predictive requirements and dataset characteristics.

Linear Regression and Logistic Regression models, while straightforward, exhibit limitations in capturing nonlinear relationships in financial data. They achieve moderate accuracy but suffer from lower precision and recall, indicating a higher rate of false positives and false negatives. Support Vector Machine (SVM) model shows improved performance with higher precision, suggesting its effectiveness in identifying crash events more accurately due to its ability to capture nonlinear relationships using kernel-based approaches.

The Decision Tree model outperforms other models with a higher accuracy rate, indicating strong predictive capability. However, it exhibits lower precision

and recall, highlighting potential challenges in distinguishing crash events from normal market fluctuations. The RNN LSTM model achieves the highest accuracy but fails to predict any positive test cases, indicating potential overfitting or data imbalance issues. While RNN LSTM models excel in capturing temporal dependencies, they may require additional optimization to address data sparsity and imbalance issues in financial datasets.

Moving forward, several avenues can be explored to enhance market crash prediction models and their applicability in real-world scenarios:

1. Improving Model Architecture: Future research should focus on refining and optimizing model architectures to improve performance and robustness. Techniques such as ensemble learning, feature engineering, and hyperparameter tuning can help enhance model accuracy and generalization.

2. Addressing Data Imbalance: Imbalance in financial datasets, where positive instances (crash events) are rare compared to normal instances, can significantly impact model performance. Techniques such as oversampling, undersampling, and synthetic data generation can help mitigate data imbalance issues and improve model effectiveness.

3. Incorporating External Data Sources: Leveraging centralized databases[10] like the USA SEC's EDGAR public database can provide valuable financial data and regulatory filings that can enhance model prediction accuracy. Integrating data from multiple sources can provide a more comprehensive understanding of market dynamics and improve predictive capability.

4. Exploring Alternative Models: Beyond traditional machine learning models, exploring advanced deep learning architectures such as convolutional neural networks (CNNs) and generative adversarial networks (GANs) can offer new

insights into market crash prediction. These models have shown promise in capturing complex patterns and dependencies in financial data.

5. Regional Adaptation: Adapting models trained on data from one market (e.g., the USA) to another market (e.g., Indian securities markets) requires careful consideration of market-specific factors and regulatory frameworks. Future research should focus on developing region-specific models tailored to the unique characteristics of Indian securities markets.

The findings from the evaluation of machine learning models for market crash prediction provide valuable insights into their strengths and limitations. By addressing these challenges and exploring new avenues for improvement, researchers can advance the field of market prediction and contribute to more accurate and reliable forecasting[19] methods for financial markets.

# REFERENCES

1. Choi, J. Y., Ahn, J. W., & Kim, J. H. (2019). Forecasting Stock Market Crashes via Machine Learning. Expert Systems with Applications.

2. Demir, E., & Gozutok, A. (2019). Forecasting bank failures and stress testing: A machine learning approach. Expert Systems with Applications.

3. Gaur, A., Tiwari, P., Singh, A., & Kapoor, A. (2020). Financial Crisis Prediction Based on Long-Term and Short-Term Memory Neural Network. Expert Systems with Applications.

4. Jacobsson, E., Stockholm University. (2009). 'How to predict crashes in financial markets with the Log-Periodic Power Law'.

5. Johansen, A., & Sornette, D. (2001). 'Large Stock Market Price Drawdowns Are Outliers'.

6. Ko, B., Kim, M., & Kim, Y. (2020). Forecasting stock market crisis events using deep and statistical machine learning techniques. Expert Systems with Applications.

7. Prabakaran, R., Devi, R. G., & Vennila, S. (2021). Novel Stock Crisis Prediction Technique-A Study on Indian Stock Market. Expert Systems with Applications.

8. Khatibi, V., Khatibi, E., and Rasouli, A. (2011). A New Support Vector Machine-Genetic Algorithm (SVM-GA) Based Method For Stock Market Forecasting. International Journal of Physical Sciences, 6(25), 6091-6097

9. Kritzman, M., and Li, Y. (2010). Skulls, Financial Turbulence, and Risk Management. Financial Analysts Journal, 44(5), 30-41

10. Laeven, L., and Valencia, F. (2020). Systemic Banking Crises Database II. IMF Economic Review, 68(2), 307-361.

11. Lee, T. K., Cho, J. H., Kwon, D. S., and Sohn, S. Y. (2019). Global Stock Market Investment Strategies Based on Financial Network Indicators Using

Machine Learning Techniques. Expert Systems with Applications, 117(1), 228-242

12. Leung, M., Daouk, H., and Chen, A.-S. (2000). Forecasting Stock Indices: A Comparison of Classification and Level Estimation Models. International Journal of Forecasting, 16(2), 173-190.

13. Lleo, S., and Ziemba, W. T. (2019). Can Warren Buffett Forecast Equity Market Corrections? European Journal of Finance, 25(4), 369-393.

14. Ohlson, J. A. (1980). Financial Ratios and the Probabilistic Prediction of Bankruptcy. Journal of Accounting Research, 18(1), 109-131.

15. Reinhart, C. M, and Reinhart, V. R. (2015). Financial Crisis, Development and Growth: A Long-term Perspective. The World Bank Economic Review, 29(1), 53-76.

16. Samitas, A., Kampouris, E., and Kenourgios, D. (2020). Machine Learning as an Early Warning System to Predict Financial Crisis. International Review of Financial Analysis, 71(C), 101507.

17. Sornette, D., and Cauwels, P. (2015), Financial Bubbles: Mechanisms and Diagnostics, Review of Behavioral Economics, 2(3), 279-305.

18. Ziemba, W. T., Zhitlukhin, M., and Lleo, S. (2017). Stock Market Crashes: Predictable and Unpredictable and What to do About Them. Singapore: World Scientific Publishing Co. Pte.Ltd.

19. H. Kim and K. Ko, (2017) "Improving forecast accuracy of financial vulnerability: partial least squares factor model approach," SSRN Electronic Journal, vol. 11, no. 2, pp. 89–97.

20. P. Save, P. Tiwarekar, and N. Ketan, (2017) "A novel idea for credit card fraud detection using decision tree," International Journal of Computer Applications, vol. 161, no. 13, pp. 6–9.