

Analisis PySpark menggunakan Algoritma Decision tree dan Random Forest dalam Prediksi Harga Saham PT Bank CIMB Niaga

PySpark analysis using Decision Tree and Random Forest Algorithm in Stock Price

Prediction of PT Bank CIMB Niaga

Tiani Ayu Lestari¹

¹Teknik Informatika, Fakultas Teknik, Universitas Pelita Bangsa

¹tianiayulestari@gmail.com

Abstract

This research explores the utilization of PySpark with Decision Tree and Random Forest algorithms in predicting the stock prices of PT Bank CIMB Niaga based on historical data from Yahoo Finance. In the context of unstable stock price fluctuations, employing Machine Learning models offers a solution to enhance predictions. The results indicate that the Random Forest model outperforms the Decision Tree, demonstrating high accuracy and low Root Mean Square Error (RMSE) values. The use of PySpark as a Big Data Processing framework facilitates efficient handling of large data volumes. The modeling process involves experimental steps with data collection sourced from Yahoo Finance. A comprehensive overview is provided through statistical descriptions and charts depicting stock price movements. Model performance evaluation is demonstrated through graphical analysis, accuracy metrics, and RMSE values. This study applies PySpark and Machine Learning algorithms for stock price prediction, providing valuable insights into financial market dynamics. The Random Forest model emerges as an effective choice for stock price analysis with potential for more accurate predictions. The practical implications extend to investors and financial analysts for making smarter investment decisions. Recommendations for future research include refining the model by considering external factors and employing advanced analytical techniques, establishing a robust foundation for further development.

Keywords: Stock Investment, PySpark, Decision Tree, Random Forest, Stock Price Prediction, Big Data

Abstrak

Penelitian ini mengeksplorasi penggunaan PySpark dengan algoritma Decision Tree dan Random Forest dalam memprediksi harga saham PT Bank CIMB Niaga berdasarkan data historis dari Yahoo Finance. Dalam konteks fluktuasi harga saham yang tidak stabil, model Machine Learning menjadi solusi untuk meningkatkan prediksi. Hasil penelitian menunjukkan bahwa model Random Forest memberikan performa lebih baik dibandingkan Decision Tree, dengan akurasi tinggi dan nilai RMSE yang rendah. Penggunaan PySpark sebagai framework Big Data Processing memungkinkan penanganan volume data besar dengan efisien. Pemodelan melibatkan tahapan eksperimental dengan pengumpulan data dari Yahoo Finance. Deskripsi statistik dan grafik pergerakan harga saham memberikan gambaran yang komprehensif. Evaluasi performa model ditunjukkan melalui analisis grafik, akurasi, dan RMSE. Penelitian ini mengaplikasikan PySpark dan algoritma Machine Learning untuk prediksi harga saham, memberikan wawasan berharga terkait dinamika pasar keuangan. Model Random Forest menonjol sebagai pilihan efektif untuk analisis harga saham dengan potensi prediksi yang lebih akurat. Implikasi praktisnya dapat digunakan oleh investor dan analis keuangan untuk pengambilan keputusan investasi yang lebih cerdas. Rekomendasi untuk penelitian selanjutnya mencakup peningkatan model dengan mempertimbangkan faktor eksternal dan teknik analisis yang lebih mutakhir, menciptakan landasan yang solid untuk pengembangan lebih lanjut.

Kata kunci: Investasi Saham, PySpark, Decision Tree, Random Forest, Prediksi Harga Saham, Big Data

Pendahuluan

Saat ini, investasi dalam saham tengah menjadi favorit di Indonesia. Saham mencerminkan kontribusi finansial seseorang dalam suatu perusahaan atau badan usaha terbatas. Fluktuasi harga saham di Indonesia bersifat tidak stabil atau tidak dapat diprediksi secara pasti[1]. Tidak pastinya pergerakan harga saham ini terjadi baik dalam jangka waktu singkat maupun jangka waktu yang lebih panjang[2]. Penggunaan teknologi dalam menganalisis serta memprediksi pergerakan harga saham telah menjadi suatu kebutuhan. Big Data merupakan terminologi yang sedang trend saat ini. Karakteristik Big Data dapat dianalisis melalui tiga aspek utama, yaitu volume, velocity, dan variety[3]. Big Data mencerminkan pertumbuhan yang besar dalam jumlah dan ragam data yang tidak dapat lagi ditangani oleh sistem basis data konvensional. Dengan bantuan alat analisis, volume data yang terlihat berbeda-beda dan bervariasi ini dapat diolah untuk menemukan pola, yang kemudian dapat menghasilkan pemahaman penting untuk mendukung pengambilan keputusan[4].

Latar Belakang

PT Bank CIMB Niaga, sebagai salah satu institusi keuangan terkemuka di Indonesia, memiliki pergerakan harga saham yang menjadi perhatian bagi banyak investor. Terdapat dua jenis utama analisis harga, yaitu analisis fundamental dan analisis teknikal[5]. Pada penelitian ini melibatkan penggunaan grafik harga historis, volume perdagangan, dan indikator teknikal lainnya untuk mengidentifikasi pola dan tren dalam pergerakan harga saham. Ini membantu para trader dan investor dalam pengambilan keputusan perdagangan.

Tinjauan Literatur Singkat

Investor dalam pengambilan keputusan untuk membeli, menjual, atau menahan saham menggunakan beberapa indikator dalam analisis teknikal[6]. Berbagai penelitian sebelumnya telah menggunakan algoritma Machine Learning, khususnya Decision Tree dan Random Forest[7][8], untuk memprediksi harga saham dengan tingkat keberhasilan yang bervariasi. Penggunaan PySpark sebagai framework Big Data Processing memberikan keunggulan dalam menangani volume data besar dan meningkatkan kinerja analisis. Apache Spark adalah platform komputasi cluster yang memiliki kinerja tinggi dan dirancang khusus untuk perhitungan yang cepat[9]. Dibangun di atas kerangka kerja Hadoop MapReduce, Apache Spark memperluas model MapReduce untuk mendukung berbagai jenis perhitungan dengan lebih efisien, termasuk query interaktif dan pemrosesan data secara real-time (stream processing)[10].

Alasan Penelitian Dilakukan

Penelitian ini bertujuan untuk mengisi kesenjangan (gap) dalam literatur yang belum sepenuhnya mengeksplorasi potensi algoritma Decision Tree dan Random Forest dengan menggunakan framework PySpark dalam konteks prediksi harga saham PT Bank CIMB Niaga. Dengan memanfaatkan teknologi Big Data dan Machine Learning[11], diharapkan penelitian ini dapat memberikan kontribusi signifikan dalam meningkatkan ketepatan prediksi pergerakan harga saham, yang pada gilirannya dapat membantu para investor dan pelaku pasar untuk mengambil keputusan investasi yang lebih cerdas.

Tujuan Penelitian

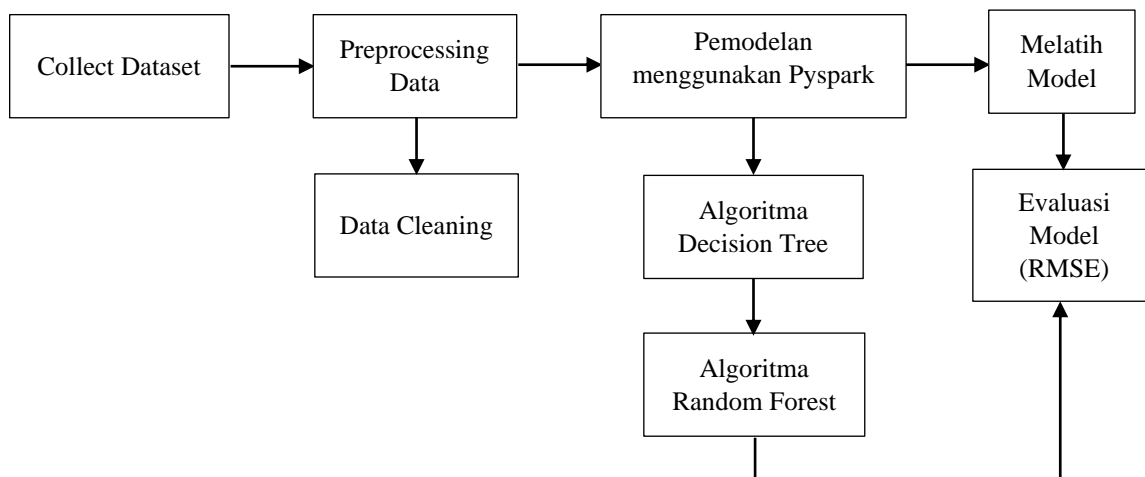
Tujuan utama dari penelitian ini adalah untuk memanfaatkan PySpark dan algoritma Decision Tree serta Random Forest guna membangun model prediksi yang akurat terhadap harga saham PT Bank CIMB Niaga. Melalui pengolahan data yang luas dan pemilihan fitur yang tepat, diharapkan model yang dikembangkan dapat memberikan hasil prediksi yang lebih mendekati realitas. Untuk menilai tingkat keakuratan, digunakan ukuran kesalahan seperti RMSE (Root Mean Square Error) dan MSE (Mean Square Error). Semakin kecil nilai RMSE dan MSE yang dihasilkan, semakin tinggi tingkat keakuratan yang ditunjukkan oleh model. RMSE merupakan nilai akar kuadrat dari MSE dan berguna sebagai metrik untuk mengevaluasi kesalahan suatu model[12].

Metode Penelitian

Metode penelitian ini bertujuan untuk menghasilkan model prediksi harga saham PT Bank CIMB Niaga dengan menggunakan PySpark dan algoritma Decision Tree serta Random Forest. Metode yang digunakan pada teknik pengumpulan data melalui data sekunder. Metode pengumpulan data melibatkan dokumentasi menggunakan sumber data dari Finance Yahoo untuk memperoleh data yang diperlukan[13].

Desain Penelitian

Penelitian ini menggunakan pendekatan eksperimental dengan analisis data historis harga saham PT Bank CIMB Niaga. Desain penelitian akan melibatkan pemrosesan dan analisis data menggunakan teknologi PySpark untuk membangun model prediksi. Langkah-langkah dalam penelitian ini tergambar pada ilustrasi yang disajikan dalam gambar 1.



Gambar 1. Tahapan Penelitian

Ruang Lingkup atau Objek Penelitian

Objek penelitian ini adalah data historis harga saham PT Bank CIMB Niaga dalam rentang waktu yang ditentukan, dengan variabel-variabel yang mencakup harga pembukaan, harga penutupan, harga tertinggi, harga terendah, volume perdagangan, dan faktor-faktor lain yang mempengaruhi pergerakan harga saham.

Tempat Penelitian

Penelitian dilakukan secara virtual menggunakan perangkat lunak komputasi dan pengolahan data yang dibutuhkan. Pengolahan data dan pengujian model dilakukan pada platform PySpark yang berbasis cloud atau komputasi terdistribusi yang memadai.

Teknik Pengumpulan Data

Data akan dikumpulkan dari sumber yang valid dan terpercaya mengenai pergerakan harga saham PT Bank CIMB Niaga pada situs <https://finance.yahoo.com/quote/BNGA.JK/>. Dengan kode saham BNGA.JK meliputi rentang waktu dari 1 Januari 2019 hingga 29 Desember 2023 yang mencakup total 1231 data.

Yahoo Finance adalah platform yang menyediakan berbagai informasi keuangan, termasuk data pasar saham, berita terkini, data historis, analisis, dan alat visualisasi untuk membantu pengguna dalam melacak dan menganalisis pergerakan pasar keuangan global[14]. Platform ini memungkinkan pengguna untuk mengakses data harga saham, indeks pasar, mata uang, serta informasi terkait perusahaan secara gratis atau dengan langganan premium tertentu. Dari kumpulan data harga saham PT. Bank CIMB Niaga, terdapat enam atribut yang terdiri dari nilai pembukaan (*open*), nilai tertinggi (*high*), nilai terendah (*low*), nilai penutupan (*close*), *adj close*, dan *volume*. Semua atribut tersebut merupakan faktor-faktor yang mempengaruhi analisis harga saham.

Teknik Analisis Penelitian

Teknik analisis yang digunakan melibatkan penggunaan PySpark untuk pemrosesan data yang besar dan kompleks. Algoritma Decision Tree dan Random Forest akan diterapkan untuk membangun model prediksi harga saham. Evaluasi model akan dilakukan dengan membandingkan hasil prediksi dengan data aktual menggunakan metrik evaluasi yang relevan. Jika diperlukan, dilakukan penyetelan parameter untuk meningkatkan kinerja model.

Hasil dan Pembahasan

Hasil Pembahasan, dapat disimpulkan bahwa penggunaan PySpark dengan algoritma Decision Tree dan Random Forest dalam memprediksi harga saham PT Bank CIMB Niaga berdasarkan data historis dari Yahoo Finance memberikan wawasan yang berharga terkait dinamika pasar keuangan. Data historis memberikan gambaran yang komprehensif tentang pergerakan harga saham, dengan volatilitas yang tercermin dalam variasi harian. Penelitian ini memberikan kontribusi positif dalam menggali potensi PySpark dan algoritma Machine Learning untuk menganalisis pergerakan harga saham, menciptakan landasan yang solid untuk pengembangan lebih lanjut.

Pengumpulan Data dan Pemrosesan Data Awal

Data awal yang terdiri dari harga pembukaan, harga penutupan, harga tertinggi, harga terendah, dan volume perdagangan telah diambil dari sumber data Yahoo Finance. Proses pengolahan awal dilakukan untuk membersihkan data dari nilai yang hilang atau tidak lengkap.

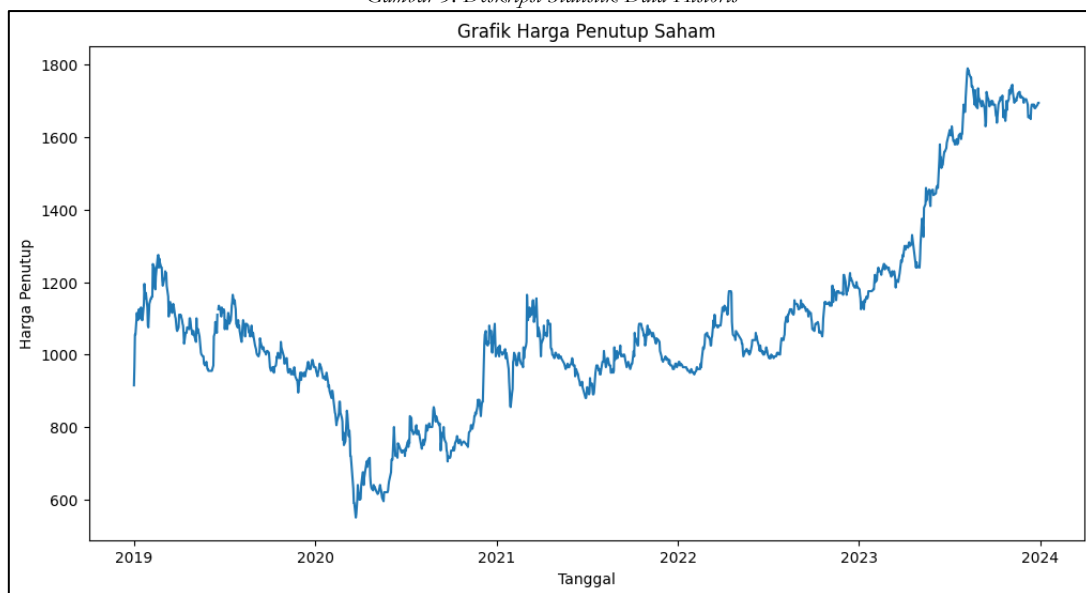
Date	Open	High	Low	Close*	Adj Close**	Volume
Dec 29, 2023	5,750.00	5,750.00	5,675.00	5,725.00	5,641.00	93,126,000
Dec 28, 2023	5,700.00	5,750.00	5,675.00	5,725.00	5,641.00	121,434,600
Dec 27, 2023	5,700.00	5,725.00	5,625.00	5,625.00	5,542.47	122,236,700
Dec 22, 2023	5,650.00	5,700.00	5,600.00	5,675.00	5,591.73	109,411,300
Dec 21, 2023	5,550.00	5,600.00	5,525.00	5,575.00	5,493.20	99,049,600
Dec 20, 2023	5,700.00	5,700.00	5,550.00	5,550.00	5,468.57	138,470,900
Dec 19, 2023	5,450.00	5,550.00	5,450.00	5,550.00	5,468.57	135,207,300
Dec 18, 2023	5,575.00	5,575.00	5,500.00	5,500.00	5,419.30	102,780,900
Dec 15, 2023	5,575.00	5,600.00	5,550.00	5,550.00	5,468.57	252,448,800
Dec 14, 2023	5,450.00	5,550.00	5,425.00	5,550.00	5,468.57	239,261,700
Dec 13, 2023	5,300.00	5,350.00	5,275.00	5,300.00	5,222.24	98,881,600
Dec 12, 2023	5,375.00	5,400.00	5,325.00	5,325.00	5,246.87	134,501,600
Dec 11, 2023	5,325.00	5,375.00	5,300.00	5,300.00	5,222.24	124,468,600
Dec 08, 2023	5,425.00	5,450.00	5,375.00	5,375.00	5,296.14	130,542,600

Gambar 2. Data Harga Saham PT CIMB NLAGA (Sumber: Yahoo Finance)

```
#describe with specific variables
df.describe(['Open', 'High', 'Low', 'Close', 'Adj Close', 'Volume']).show()
```

summary	Open	High	Low	Close	Adj Close	Volume
count	1231	1231	1231	1231	1231	1231
mean	1081.462258326565	1096.4662875710803	1067.3639317627944	1080.6661251015435	926.1658645493533	7233682.615759545
stddev	259.0825244323926	260.1461102136569	258.4851444184299	259.8772439688772	315.42536364234667	1.025783196643582E7
min	540.0	575.0	515.0	550.0	408.88898	0.0
max	1790.0	1815.0	1770.0	1790.0	1790.0	1.511562E8

Gambar 3. Deskripsi Statistik Data Historis



Gambar 4 Grafik Pergerakan Saham

Gambar 2 menunjukkan deskripsi statistik dari data historis harga saham PT Bank CIMB Niaga. Pada Gambar 3 menampilkan data yang terdiri dari total jumlah data, rata-rata, standar deviasi, nilai minimum, kuartil, dan nilai maksimum dari harga pembukaan, penutupan, tertinggi, terendah, dan volume perdagangan. Gambar 4 menampilkan grafik pergerakan harga saham PT Bank CIMB Niaga selama 5 tahun terakhir. Grafik ini menggambarkan tren harga saham harian dari data historis yang diambil dari Yahoo Finance.

Pemodelan dan Evaluasi

```
[12] # Inisialisasi model Decision Tree
      dt = DecisionTreeRegressor(featuresCol='features', labelCol='close')

      # Melatih model Decision Tree
      dt_model = dt.fit(train_data)

[13] # Inisialisasi model Random Forest
      rf = RandomForestRegressor(featuresCol='features', labelCol='close')

      # Melatih model Random Forest
      rf_model = rf.fit(train_data)
```

Gambar 5 Pemodelan menggunakan Decision Tree dan Random Forest

```
[14] # Evaluasi performa model Decision Tree pada data pengujian
dt_predictions = dt_model.transform(test_data)
dt_evaluator = RegressionEvaluator(labelCol='Close', metricName='rmse')
dt_rmse = dt_evaluator.evaluate(dt_predictions)
print("Root Mean Squared Error (Decision Tree):", dt_rmse)

Root Mean Squared Error (Decision Tree): 19.07783220809059

[15] # Evaluasi performa model Random Forest pada data pengujian
rf_predictions = rf_model.transform(test_data)
rf_evaluator = RegressionEvaluator(labelCol='Close', metricName='rmse')
rf_rmse = rf_evaluator.evaluate(rf_predictions)
print("Root Mean Squared Error (Random Forest):", rf_rmse)

Root Mean Squared Error (Random Forest): 16.456338280164854
```

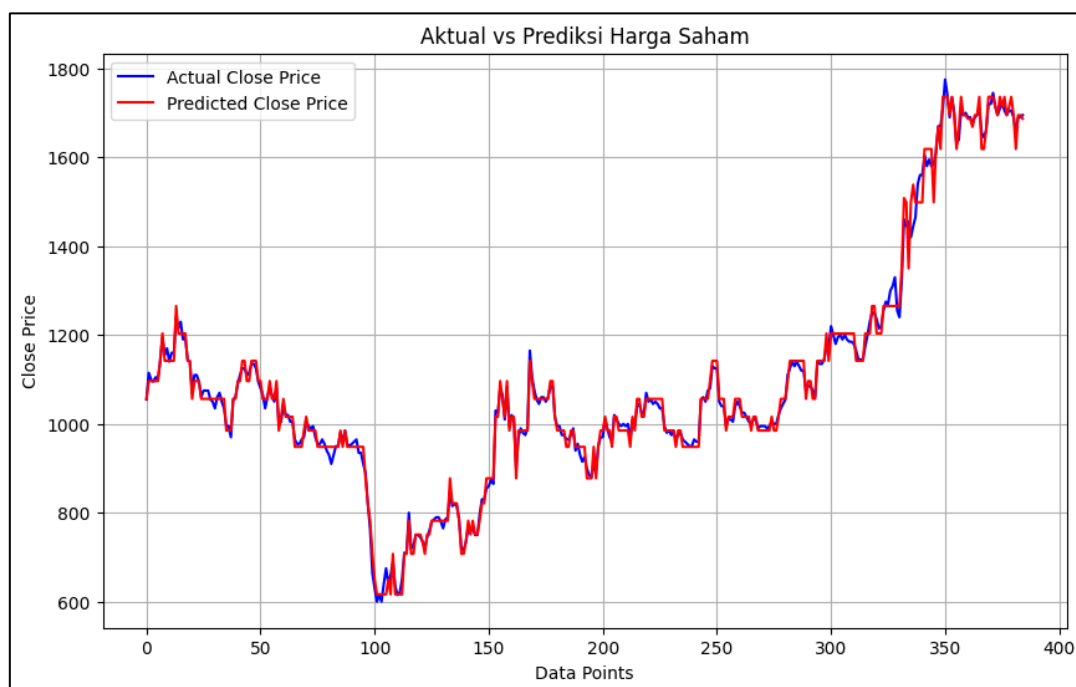
Gambar 6 Akurasi menggunakan RMSE

Pada gambar 5 dan gambar 6 menunjukkan hasil evaluasi performa model Decision Tree dan Random Forest dalam memprediksi harga saham PT Bank CIMB Niaga. Metrik evaluasi menggunakan RMSE dan akurasi R2 atau R-Squared untuk mengevaluasi kinerja kedua model pada Gambar 7.

```
# Menghitung R-squared
r2 = evaluator.evaluate(predictions)
print("R-squared:", r2)

R-squared: 0.9960837448384626
```

Gambar 7 Akurasi R-Squared



Gambar 8 Hasil Analisis

Deskripsi statistik data menunjukkan bahwa grafik pergerakan harga saham menunjukkan tren naik dan turun yang berkaitan dengan faktor-faktor eksternal yang memengaruhi pasar. Dari hasil evaluasi, model Random Forest menghasilkan *Root Mean Squared Error* (RSME) sebesar 16.45 menunjukkan performa yang lebih baik dibandingkan dengan Decision Tree menghasilkan *Root Mean Squared Error* (RSME) sebesar 19,08 dengan akurasi R-Squared sebesar 0,99 atau sebesar 99%. Ini mengindikasikan bahwa model Random Forest mampu memberikan prediksi yang lebih akurat terhadap pergerakan harga saham PT Bank CIMB Niaga berdasarkan data historis dari Yahoo Finance. Dan dapat ditampilkan pada Gambar 8 bagaimana harga saham aktual (garis biru) dengan nilai yang diprediksi oleh model (garis merah) untuk setiap titik data yang ada. Ini memberikan gambaran visual tentang seberapa baik model dapat memprediksi harga saham berdasarkan data yang ada.

Hasil ini menegaskan bahwa penerapan algoritma Random Forest dan Decision Tree pada data historis harga saham menggunakan PySpark dapat memberikan nilai prediksi yang lebih dapat diandalkan, sehingga dapat menjadi landasan bagi pengambilan keputusan investasi yang lebih cerdas dan akurat di pasar keuangan.

Kesimpulan

Pada hasil pembahasan, dapat disimpulkan bahwa penggunaan PySpark dengan algoritma Decision Tree dan Random Forest dalam memprediksi harga saham PT Bank CIMB Niaga berdasarkan data historis dari Yahoo Finance memberikan wawasan yang berharga terkait dinamika pasar keuangan. Data historis memberikan gambaran yang komprehensif tentang pergerakan harga saham, dengan volatilitas yang tercermin dalam variasi harian. Lebih lanjut, evaluasi performa model menunjukkan bahwa Random Forest mampu memberikan prediksi yang lebih akurat dibandingkan Decision Tree. Dengan nilai akurasi yang lebih tinggi, model Random Forest menjadi pilihan yang lebih efektif dalam mengatasi kompleksitas dan variasi data historis harga saham.

Hasil ini memiliki implikasi praktis yang signifikan bagi para pemangku kepentingan pasar keuangan, terutama investor dan analis keuangan. Model prediksi yang diperoleh dapat digunakan sebagai panduan dalam pengambilan keputusan investasi, membantu mengidentifikasi tren potensial, dan meningkatkan pemahaman terhadap risiko yang terkait dengan pergerakan harga saham. Rekomendasi untuk penelitian selanjutnya melibatkan peningkatan model dengan mempertimbangkan faktor-faktor eksternal dan teknik analisis yang lebih mutakhir. Pengembangan ini diharapkan dapat menghadapi dinamika pasar yang terus berubah dan meningkatkan ketepatan prediksi harga saham secara keseluruhan. Kesimpulannya, penelitian ini memberikan kontribusi positif dalam menggali potensi PySpark dan algoritma Machine Learning untuk menganalisis pergerakan harga saham, menciptakan landasan yang solid untuk pengembangan lebih lanjut.

Daftar Rujukan

- [1] L. I. E. S. Setyo Wira Rizki, "Analisis Pengukuran Kinerja Portofolio Optimal Indeks Saham Lq45 Dengan Model Black-Litterman," *Bimaster Bul. Ilm. Mat. Stat. dan Ter.*, vol. 8, no. 3, pp. 555–562, 2019, doi: 10.26418/bbimst.v8i3.33904.
- [2] M. Jannah, F. Mardika, L. H. Hasibuan, and D. M. Putri, "Pemodelan Data Saham Menggunakan Analisis Time Series Dengan Pendekatan Copula Gaussian," *Math Educ. J.*, vol. 5, no. 2, pp. 163–174, 2021, doi: 10.15548/mej.v5i2.3124.
- [3] C. Z. Tumbel, H. Sitepu, and M. Hutagalung, "Analisis Big Data Berbasis Stream Processing Menggunakan Apache Spark," *J. Telemat.*, vol. 11, no. 1, p. 6, 2017, [Online]. Available: <http://journal.ithb.ac.id/telematika/article/view/145>.
- [4] R. Purnomo, Wowon Priatna, and Tri Dharma Putra, "Implementasi Big Data Analytical Untuk Perguruan Tinggi Menggunakan Machine Learning," *J. Inform. Inf. Secur.*, vol. 2, no. 1, pp. 77–88, 2021, doi: 10.31599/jiforty.v2i1.633.

- [5] N. A. 'Izzah *et al.*, "Analisis Teknikal Pergerakan Harga Saham Dengan Menggunakan Indikator Stochastic Oscillator Dan Weighted Moving Average," *Keunis*, vol. 9, no. 1, p. 36, 2021, doi: 10.32497/keunis.v9i1.2307.
- [6] P. Prihatiningsih, E. Duriany, A. Sunindyo, and M. A. Kodir, "Strategi Investasi Saham Di Bursa Efek Indonesia Dengan Analisis Teknikal," *J. Aktual Akunt. Kenang. Bisnis Terap.*, vol. 5, no. 2, p. 198, 2022, doi: 10.32497/akunbisnis.v5i2.4060.
- [7] P. Studi and T. Informatika, "Penerapan Decision Tree untuk memprediksi pergerakan harga saham dengan analisis Foreign Flow," no. 13519196.
- [8] A. Pratomo, R. F. Umbara, and A. A. Rohmawati, "Prediksi Pergerakan Harga Saham Dengan Metode Random Forest Menggunakan Trend Deterministic Data Preparation (Studi Kasus Saham Perusahaan Pt Astra International Tbk, Pt Garuda Indonesia Tbk, Dan Pt Indosat Tbk)," *eProceedings Eng.*, vol. 6, no. 1, pp. 2545–2556, 2019.
- [9] R. A. Fauzi, I. Cholissodin, and B. Rahayudi, "Pemanfaatan Spark untuk Analisis Sentimen Mengenai Netralitas Berita dalam Membahas Pemilu Presiden 2019 Menggunakan Metode Naive Bayes Classifier," *J. Pengemb. Teknol. Inf. dan Ilmu Komput.*, vol. 5, no. 3, pp. 1070–1077, 2021, [Online]. Available: <https://j-ptiik.ub.ac.id/index.php/j-ptiik/article/view/8741>.
- [10] S. Oliviandi, A. B. Osmond, and R. Latuconsina, "Implementasi Apache Spark Pada Big Data Berbasis Hadoop Distributed File System," *e-Proceeding Eng.*, vol. 5, no. 1 Maret, pp. 1005–1012, 2018.
- [11] J. S. Prasetyo, "Stock Price Prediction Using Machine Learning With Long Short Term Memory Method (LSTM)," *Kilat*, vol. 12, no. 1, pp. 64–78, 2023, doi: 10.33322/kilat.v12i1.1723.
- [12] W. Hastomo, A. S. B. Karno, N. Kalbuana, E. Nisfiani, and L. ETP, "Optimasi Deep Learning untuk Prediksi Saham di Masa Pandemi Covid-19," *J. Edukasi dan Penelit. Inform.*, vol. 7, no. 2, p. 133, 2021, doi: 10.26418/jp.v7i2.47411.
- [13] M. K. Abdi, C. N. Farida, D. A. Ramadhanik, L. Anifa, A. Devit, and C. Putra, "Analisis Pengambilan Keputusan Saham Dengan Metode Decision Analysis Pada Investasi Saham PT . Indofood Sukses Makmur Tbk .," *Snhrp*, no. April, pp. 68–76, 2022, [Online]. Available: <https://snhrp.unipasby.ac.id/prosiding/index.php/snhrp/article/view/296%0Ahttps://snhrp.unipasby.ac.id/prosiding/index.php/snhrp/article/download/296/246>.
- [14] Reza Maulana and Devy Kumalasari, "Analisis Dan Perbandingan Algoritma Data Mining Dalam Prediksi Harga Saham Ggrm," *J. Inform. Kaputama*, vol. 3, no. 1, pp. 22–28, 2019, [Online]. Available: <https://finance.yahoo.com/quote/GGRM.J>.