

Prediksi Harga Saham Menggunakan Improved Multiple Linear Regression Untuk Pencegahan Data Outlier

Abidatul Izzah^{*1}, Ratna Widyastuti²

^{1,2}Politeknik Kediri

abidatul.izzah90@poltek-kediri.ac.id^{*1}, nana89widya@gmail.com²

Abstrak

Seiring berkembangnya bisnis modern, prediksi harga saham selalu mendapat perhatian khusus oleh pakar ekonomi. Prediksi nilai saham menggunakan model Multiple Linear Regression (MLR) telah terbukti memberikan nilai prediksi yang presisi dan cukup baik. Namun di sisi lain, regresi linear memiliki beberapa kelemahan terhadap data outlier. Oleh karena itu pada penelitian ini dilakukan prediksi nilai harga saham menggunakan MLR yang dibantu dengan teknik K-Means dan Moving Average (MA) untuk mengatasi pengaruh data outlier. Pengujian diawali dengan pengumpulan data dan pra-proses data. Data harga saham yang akan digunakan dalam pengujian diperoleh dari laman *finance.yahoo.com* dengan kategori "Jakarta Composite Index (^JKSE)". Selanjutnya proses prediksi dilakukan dengan hybrid MLR dengan K-Means dan MA untuk mengatasi titik-titik saham yang outlier. Dari hasil yang diperoleh, dapat dilihat bahwa pendekatan paling baik ditunjukkan oleh metode MLR dan MA yakni dengan nilai MSE sebesar 15087.465, RMSE sebesar 122.831, dan MAPE sebesar 3.255.

Kata kunci: K-Means, Moving Average, Outlier, Prediksi, Regresi, Saham

Abstract

Nowadays, stock price prediction got special attention by economist or investor. Besides that, stock prediction based on Multiple Linear Regression (MLR) shows a good prediction. However, linear regression has a weakness to outlier data. Therefore, in this study, the stock price prediction using MLR is aided by K-Means and Moving Average (MA) to help MLR in outlier case. In this paper, stock data is obtained from the *finance.yahoo.com* on category "Jakarta Composite Index (^JKSE)". Improved MLR with K-Means and MA are used to overcome the outlier stock. From the results obtained that the best approach is shown by MLR and MA with the value of MSE is 15087.465, RMSE is 122.831, and MAPE is 3.255.

Keywords: K-Means, Moving Average, Outlier, Prediction, Regression, Stock

1. Pendahuluan

Peramalan atau prediksi merupakan sebuah proses untuk mengetahui nilai pada waktu yang akan datang. Kebutuhan prediksi pun banyak ditemui, antara lain prediksi nilai tukar rupiah, prediksi curah hujan, atau prediksi Indeks Harga Saham Gabungan (IHSG). Seiring berkembangnya bisnis modern, prediksi harga saham selalu mendapat perhatian khusus oleh pakar ekonomi. Dengan adanya kemampuan prediksi atau menebak nilai harga saham, sebuah perusahaan/investor dapat menganalisis dan memprediksi langkah kebijakan yang optimal untuk membuat keputusan pembelian/penjualan saham yang sesuai. Sehingga, prediksi atau peramalan menjadi alat bantu penting bagi perencanaan yang lebih efektif dan efisien. Prediksi dilakukan dengan mencatat riwayat harga saham harian, mingguan, bulanan, bahkan tahunan. Dalam melakukan prediksi harga saham dibutuhkan data secara berkala dengan satuan waktu yang dapat digolongkan sebagai data *time series*. Data *time series* dapat diartikan sebagai serangkaian pengamatan yang berasal dari suatu sumber tetap yang terjadi berdasarkan waktu *t* secara berurutan dengan interval waktu yang tetap [1]. Beberapa penelitian tentang pendekatan prediksi harga saham telah dilakukan menggunakan berbagai metode antara lain Jaringan Syaraf Tiruan [2], Hybrid regresi dengan Algoritma Genetika [3], dan Support Vector Regression [4]. Salah satu metode yang umum digunakan untuk memprediksi data adalah metode Regresi karena memiliki keunggulan, yakni perhitungannya yang mudah. Dalam penelitian Rahmi dkk [3]

telah dilakukan prediksi nilai saham menggunakan model MLR dan memberikan nilai prediksi yang presisi dan cukup baik. Namun disisi lain, regresi linear memiliki beberapa kelemahan, yakni lemah terhadap data outlier. Jika data outlier digunakan untuk perhitungan nilai koefisien regresi maka akan menimbulkan bias pada hasil prediksi. Pada penelitian sebelumnya [5], Sari dkk telah melakukan uji coba pada *simple linear regression* (SLR) atau Regresi Linear Sederhana dengan improvisasi pemilihan titik data saat penentuan koefisien regresi untuk mengoreksi hasil citra kamera ponsel. Pemilihan titik ini dimaksudkan agar titik-titik outlier tidak diperhitungkan. Pemilihan data ini dilakukan menggunakan teknik *clustering*, yakni K-Means. Dari hasil yang diperoleh menunjukkan bahwa dengan pemilihan titik, hasil warna yang diperoleh lebih baik. Hal ini memungkinkan teknik pemilihan data juga bisa diterapkan pada metode Multiple Linear Regression (MLR).

Dari uraian tersebut, pada penelitian ini akan dilakukan prediksi nilai harga saham menggunakan MLR dimana data yang diperhitungkan adalah data yang dipilih menggunakan teknik K-Means. Diharapkan dengan demikian hasil prediksi harga saham dapat lebih presisi karena tidak mengandung outlier. Selain itu, metode yang juga sering digunakan untuk penghalusan data hasil prediksi adalah Moving Average (MA) yang memperhitungkan nilai rata-rata bergerak. Oleh karena itu, pada penelitian ini akan dilakukan prediksi nilai harga saham menggunakan MLR yang diimprovisasi menggunakan K-Means dan MA.

2. Metode Penelitian

2.1 Data Mining dan *Forecasting*

Data Mining (Pengolahan Data) merupakan disiplin ilmu yang mempelajari teknik pencarian informasi dari data yang jumlahnya sangat besar. Proses pencarian informasi tersebut dapat dibedakan menjadi beberapa komponen sebagai berikut [6]:

1. Klasifikasi (*Classification*)

Teknik klasifikasi digunakan untuk mengidentifikasi data sesuai kelompok data yang telah tersedia. Proses identifikasi didasari dengan proses pembelajaran yang disebut dengan *supervised learning*. Metode yang kerap digunakan untuk klasifikasi data adalah Neural Networks dan Decision Trees.

2. Pengelompokan (*Clustering*)

Teknik pengelompokan data digunakan untuk menentukan posisi sebuah data yang dikelompokkan berdasarkan kesamaan atribut yang dimiliki. Metode yang kerap digunakan untuk pengelompokan data adalah K-Means, K-Medoids, dan Fuzzy C Means.

3. Asosiasi (*Association*)

Asosiasi digunakan untuk menganalisis sebab terjadinya peristiwa sehingga mengakibatkan sebuah peristiwa. Teknik asosiasi dilakukan dengan membangkitkan *rule* yang dievaluasi dengan perhitungan *support* dan *confident*. Metode yang umum digunakan adalah Apriori Algorithm.

4. Prediksi (*Prediction*)

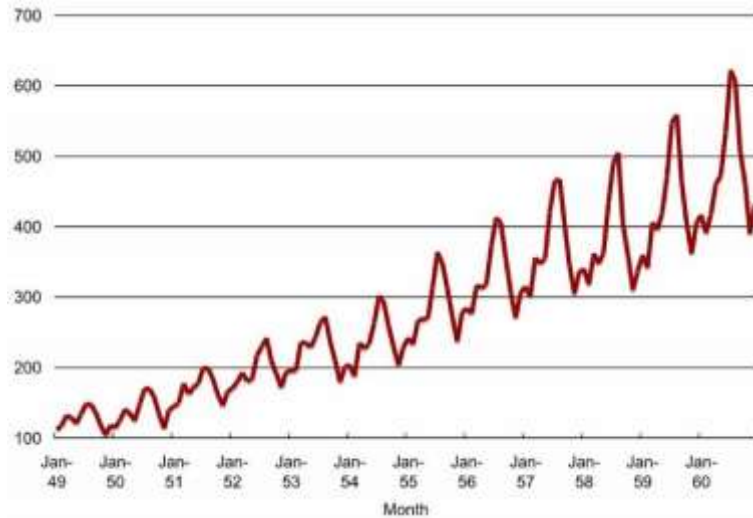
Prediksi dilakukan untuk meramalkan (*forecasting*) kejadian apa yang terjadi di masa depan dengan memperhitungkan kejadian yang sekarang. Metode yang sering digunakan untuk prediksi adalah Regression atau Neural Network.

Lebih lanjut, [7] mengartikan prediksi sebagai seni dan ilmu untuk memperkirakan kejadian di masa depan. Hal ini dapat dilakukan dengan melibatkan pengambilan data historis dan memproyeksikannya ke masa mendatang dengan suatu bentuk model matematis atau prediksi intuisi bersifat subyektif, atau menggunakan kombinasi model matematis. Sedangkan menurut Prasetya [8], prediksi merupakan suatu usaha untuk meramalkan keadaan masa mendatang melalui pengujian keadaan pada masa lalu. *Forecasting* berkaitan dengan upaya memperkirakan apa yang terjadi di masa depan berbasis pada metode ilmiah (ilmu dan teknologi) serta dilakukan secara matematis.

2.2 Data *Time Series*

Data *time series* adalah susunan observasi berurut menurut waktu. Data *time series* merupakan serangkaian data pengamatan yang berasal dari satu sumber tetap yang terjadinya berdasarkan indeks waktu secara berurutan dengan interval waktu yang tetap. Digunakan notasi untuk menyatakan nilai numerik dari pengamatan, dengan menunjukkan periode waktu terjadinya pengamatan. Proses prediksi data *time series* tidak melibatkan variabel independen lain selain

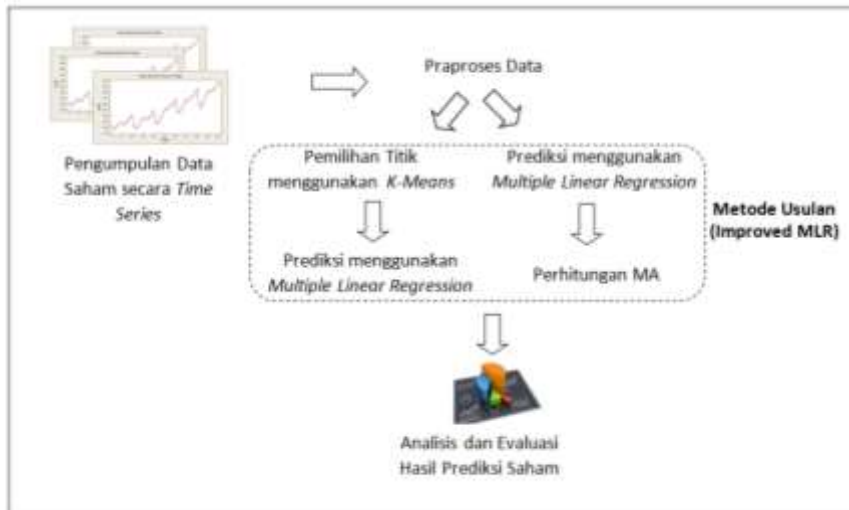
indeks waktu itu sendiri [1]. Data tersebut didasarkan pada urutan dari titik-titik data yang berjarak sama dalam waktu seperti mingguan, bulanan, kuartalan, dan lainnya. Metode Peramalan Time Series terdiri dari Pendekatan Naif (Naïve Approach), Rataan Bergerak (Moving Average), Penghalusan Eksponensial (Exponential Smoothing), dan Proyeksi Tren (Trend Projection) [7]. Gambar 1 adalah grafik yang menunjukkan plot data *time series*.



Gambar 1. Contoh Grafik Time Series

2.3 Metode

Adapun langkah-langkah yang akan dilakukan pada untuk memprediksi nilai harga saham diawali dengan pengumpulan data dan pra-proses data. Selanjutnya dilakukan proses prediksi menggunakan MLR. Proses prediksi dilakukan menggunakan *hybrid* MLR dengan K-Means dan MA untuk mengatasi titik-titik saham yang outlier. Setelah hasil prediksi diperoleh, maka hasil tersebut dibandingkan, dianalisis, dan dievaluasi. Tahapan ini dapat ditunjukkan pada Gambar 2.



Gambar 2. Alur Metode Penelitian

2.3.1 Data

Data harga saham yang akan digunakan dalam pengujian diperoleh dari laman finance.yahoo.com dengan kategori “Jakarta Composite Index (^JKSE)”. Harga saham ini bernilai rupiah atau IDR. Data harga saham yang akan digunakan adalah data harian selama kurang lebih tiga tahun yang diambil pada bulan 2 Januari 2013 sampai dengan 21 Desember 2015 yang berbentuk *time series* yang disertai variabel *open*, *high*, *low*, *close*, *volume*, dan *adj close* seperti yang ditunjukkan pada Tabel 1.

Tabel 1. Contoh Keterangan Tabel

Tanggal	Open	High	Low	Close	Volume	Adj Close
02-01-13	4322.582	4364.665	4316.424	4346.475	3.28E+09	4346.475
03-01-13	4356.411	4401.334	4356.411	4399.258	4.33E+09	4399.258
...
...
21-12-15	4452.648	4490.68	4452.017	4490.68	2.1E+09	4490.68

Dalam penelitian ini, jumlah parameter masukan (*input*) yang digunakan sebanyak 3 titik, dengan asumsi bahwa jumlah tersebut mewakili data selama 3 hari. Sedangkan target *output* yang akan diprediksi adalah data pada hari ke-4. Sehingga data harga saham hari ke 1 sampai dengan hari ke 3 digunakan untuk perhitungan prediksi nilai harga saham pada hari ke 4 dan seterusnya. Dalam penelitian ini, prediksi dilakukan untuk *variable* harga penutupan saham (*close price*) karena *close price* merupakan nilai terpenting dalam mencerminkan posisi harga dimana investor bisa mengambil keputusan membeli/menjual saham.

Langkah selanjutnya adalah pra-proses data yang dilakukan agar data siap diolah. Pada data *time series*, terdapat beberapa langkah pra-proses data, yaitu Penanganan *missing value* dan outlier data. Penanganan *missing value* adalah dengan menghilangkan satu titik atau satu hari yang mengandung *missing value*. Sedangkan penanganan outlier data diasumsikan hal ini terjadi pada bulan-bulan tertentu atau karena ada kejadian tak terduga, sehingga nilai harga saham tidak normal. Untuk kasus ini titik tersebut tetap diperhitungkan karena ingin dibuktikan keunggulan metode yang diusulkan. Setelah data melakukan pra-proses, data dirancang dalam bentuk parameter *input* dan *output* seperti pada Tabel 2.

Tabel 2. Pemodelan Data Saham

Pola ke-	Input (X_1, X_2, X_3)			Output	Target
	X_1	X_2	X_3		
1	4346.475	4399.257	4410.020	?	4392.378
2	4399.257	4410.020	4392.378	?	4362.928
...

2.3.2 Multiple Linear Regression

Multiple Linear Regression atau Regresi Linear Berganda dapat digunakan dalam prediksi atau peramalan yang disusun atas dasar pola hubungan data yang relevan dimasa lalu. Pada metode regresi umumnya variabel yang diprediksi seperti penjualan atau permintaan suatu produk, dinyatakan sebagai variabel yang dicari (*dependent variable*), variabel ini dipengaruhi besarnya oleh variabel bebas (*independent variable*). Pada dasarnya terdapat dua macam analisa hubungan dalam penyusunan peramalan [9], yaitu analisa *cross section* atau model sebab akibat (*causal model*) dan analisa deret waktu (*time series*) yang akan dibahas dalam penelitian ini.

Suatu langkah yang penting dalam peramalan *time series* adalah mempertimbangkan jenis pola yang terdapat dari data observasi, sehingga metode tersebut dapat diuji kembali. Pola yang ditunjukkan dengan analisa regresi yang sederhana mengasumsikan bahwa hubungan diantara >2 variabel dapat dinyatakan dengan suatu garis lurus. Regresi linear berganda dapat dihitung berdasarkan Persamaan 1.

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n \quad (1)$$

Dimana $\beta_0, \beta_1, \beta_2, \dots, \beta_n$ merupakan koefisien MLR yang dihitung berdasarkan Persamaan 2.

$$\beta = (X^T X)^{-1} X^T Y \quad (2)$$

2.3.3 Pengelompokan Data K-Means

K-Means merupakan salah satu teknik pengelompokan data. Tujuan dari algoritma ini adalah membagi data menjadi beberapa kelompok yang memiliki kesamaan atribut. Algoritma K-Means *clustering* [10] diawali dengan menentukan bilangan bulat k , dimana k merepresentasikan

jumlah kelompok. Kemudian k buah titik pusat (*centroid*) dipilih secara acak. Kemudian kelompokkan masing-masing data kepada titik pusat terdekat berdasarkan persamaan Euclid yang dihitung pada Persamaan 3 sehingga terbentuk k buah kelompok (*cluster*). Dalam setiap kelompok tersebut, nilai titik pusat diperbarui dengan Persamaan 4. Proses ini diulangi sampai nilai dari titik *centroid* tidak lagi berubah.

$$d = \sum_{j=1}^k \sum_{i=1}^n \|x_i^j - c_i\|^2 \quad (3)$$

$$c = \frac{\sum_{i=1}^n x_i}{n} \quad (4)$$

2.3.4 Moving Avarage

Peramalan rataan bergerak dilakukan dengan menggunakan sejumlah data aktual masa lalu untuk menghasilkan peramalan ke- n . Secara matematis, untuk menghitung rataan bergerak sederhana dinyatakan pada persamaan 5.

$$\bar{X}_n = \frac{\sum_{i=1}^n X_i}{n} \quad (5)$$

2.3.5 Improved Multiple Linear Regression

Metode Improved Multiple Linear Regression yang dimaksud dalam penelitian ini adalah metode prediksi MLR yang dibantu dengan K-Means dan Moving Average. Teknik K-Means diimplementasikan pada pemilihan titik riwayat data saham untuk membentuk pola dan perhitungan koefisien regresi. Metode ini disusun seperti pada Gambar 3.

Algoritma 1. Multiple Linear Regression Improved by K-Means()
<pre> Start input history data as X, k, n init N as number of data object X, c, beta, y'; for i=i:N get X(1:n,:) calculate c = kmeans(X,k) Calculate beta using (2) Calculate y' using (1) endfor end </pre>

Gambar 3. Metode Multiple Linear Regression Improved by K-Means

Metode prediksi Multiple Linear Regression selanjutnya merupakan implementasi MA dalam memperhalus hasil data prediksi. Metode ini disusun seperti pada Gambar 4.

Pseudocode 2. Multiple Linear Regression Improved by MA()
<pre> Start input history data as X, n init N as number of data object X, beta, y'; for i=i:N get X(1:n,:) Calculate beta using (2) Calculate y' using (1) Smooth y' using (5) endfor end </pre>

Gambar 4. Metode Multiple Linear Regression Improved by MA

3. Hasil Penelitian dan Pembahasan

3.1 Hasil Uji Coba

Pada tahap ini akan dilakukan evaluasi dengan cara membandingkan hasil prediksi *time series* dengan metode Multiple Linear Regression dan Multiple Linear Regression dengan pemilihan titik Improved Multiple Linear Regression. Ukuran kesalahan pola hasil prediksi adalah kesalahan yang terjadi antara data prediksi dan data aktual. Kesalahan tersebut direpresentasikan menggunakan *mean square error* (MSE), yakni merupakan rata-rata selisih kuadrat antara nilai yang diprediksikan dengan diamati, *root mean square error* (RMSE) merupakan akar dari MSE, dan *mean absolute percentage error* (MAPE), yaitu rata-rata diferensiasi absolut antara nilai yang diprediksi dan aktual, dinyatakan sebagai persentase nilai aktual yang dihitung berdasarkan Persamaan 6, Persamaan 7, dan Persamaan 8.

$$MSE = \frac{1}{n} \sum (Y_t - \hat{Y}_t)^2 \quad (6)$$

$$RMSE = \sqrt{\frac{1}{n} \sum (Y_t - \hat{Y}_t)^2} \quad (7)$$

$$MAPE = \frac{1}{n} \sum_{t=1}^n \frac{|Y_t - \hat{Y}_t|}{Y_t} \quad (8)$$

Pengujian pertama diawali dengan skenario prediksi harga saham menggunakan metode MLR dan K-Means (KM). Data yang digunakan adalah data riwayat saham selama tiga hari yang telah dimodelkan seperti pada Tabel 2. Untuk pembentukan pola diambil 5 *record* data untuk kemudian dipilih $k=3$ titik *centroid* diantaranya. Hasil pengujian menunjukkan nilai ukuran evaluasi MSE, RMSE, dan MAPE yang diperoleh MLR-KM lebih tinggi (Tabel 3). Hal ini menggambarkan keadaan bahwa adanya K-Means tidak membantu MLR dalam memprediksi harga saham. Hasil perbandingan prediksi menggunakan MLR dan MLR-KM dapat dilihat pada Gambar 5.

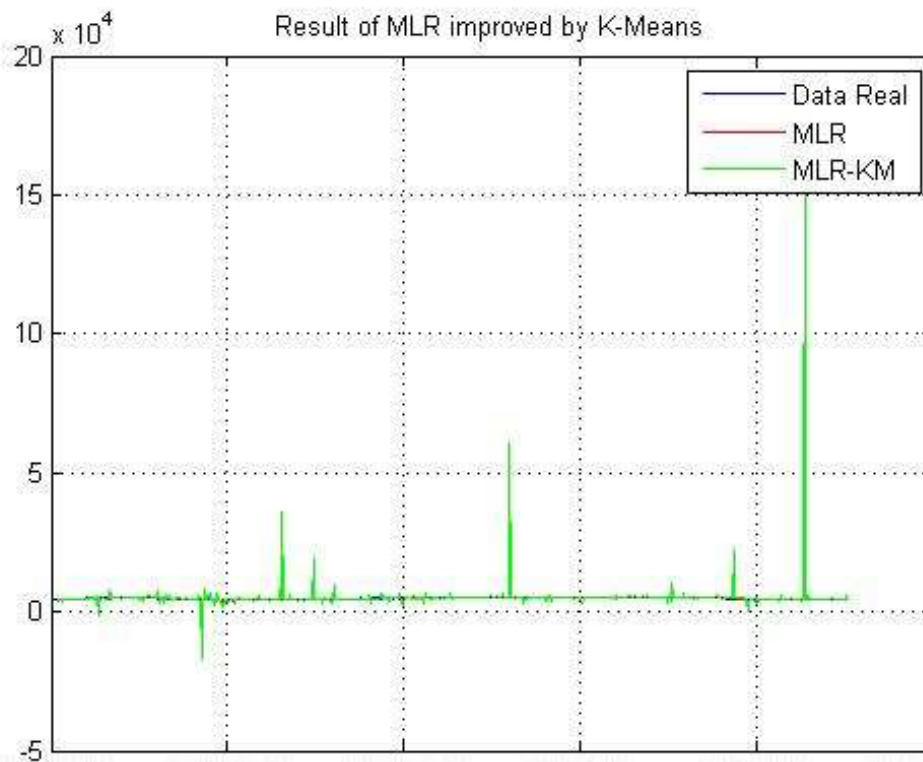
Skenario selanjutnya prediksi harga saham menggunakan metode MLR dan MA. Data yang digunakan adalah data riwayat saham selama tiga hari dan untuk pembentukan pola diambil 3 *record* data. Kemudian hasil prediksi diperhalus menggunakan teknik MA. Hasil yang diperoleh menunjukkan nilai ukuran evaluasi MSE, RMSE, dan MAPE yang diperoleh MLR-MA lebih rendah (Tabel 3). Hal ini menggambarkan keadaan bahwa MA membantu MLR dalam memprediksi harga saham dengan memperhalus hasil prediksi MLR dengan rata-rata hasil prediksi sebelumnya sehingga data outlier dapat ditangani. Hasil perbandingan prediksi menggunakan MLR dan MLR-MA dapat dilihat pada Gambar 6.

Tabel 3. Hasil Uji Coba

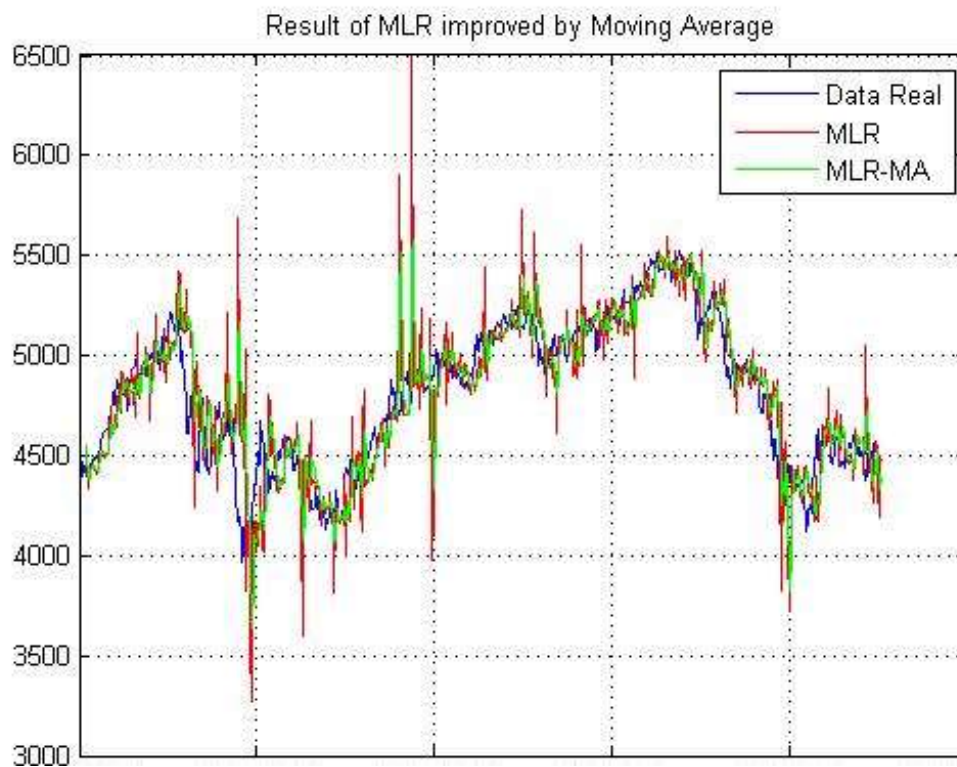
Ukuran Evaluasi	Metode		
	MLR	MLR-KM	MLR-MA
MSE	26849.807	55143754.119	15087.465
RMSE	163.859	7425.884	122.831
MAPE	5.754	12199.802	3.255

3.2 Pembahasan

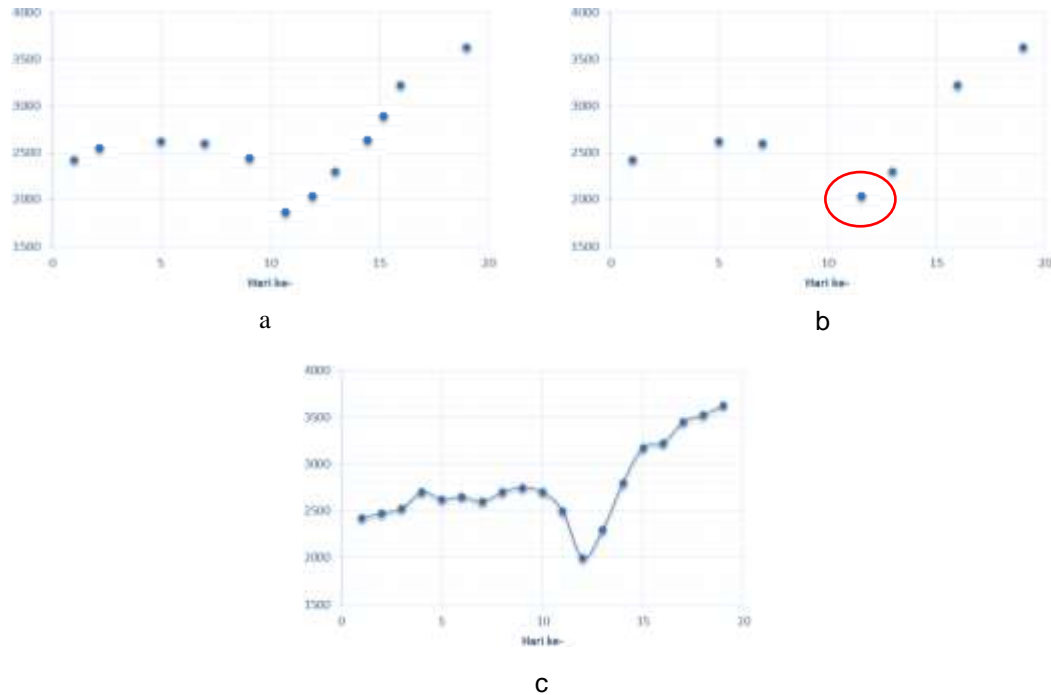
Berdasarkan hasil uji coba, dapat dilihat bahwa metode MLR yang diperbaiki menggunakan teknik MA menunjukkan *error* yang paling rendah. Sebaliknya, hasil yang buruk ditunjukkan oleh MLR-KM. Hal ini disebabkan K-Means tidak berfungsi baik dalam hal menangani data outlier pada data saham. *Centroid* yang merupakan hasil pemilihan K-Means pada titik-titik saham nyatanya tidak menghilangkan titik-titik outlier sesungguhnya. Sebagai contoh, Gambar 7A menunjukkan data saham dalam runtun waktu tertentu. Kemudian pada data tersebut dilakukan pengelompokan dan pemilihan titik. Namun, pada kelompok data outlier tetap diwakili dengan salah satu *centroid* seperti terlihat pada Gambar 7B, selanjutnya hasil pengelompokan data dapat dilihat pada ilustrasi Gambar 7C.



Gambar 5. Perbandingan Hasil MLR dan MLR-KM



Gambar 6. Perbandingan Hasil MLR dan MLR-MA



Gambar 7. Hasil Pengelompokan Tetap Mewakili Data Outlier
 (a) Data Saham Dalam Runtun Waktu Tertentu (b) Centroid yang Mewakili Kelompok Data Outlier (c) Data Saham Hasil Pengelompokan

4. Kesimpulan

Pada penelitian ini telah dilakukan prediksi harga saham menggunakan *Multiple Linear Regression* dengan K-Means dan Moving Average. Dari hasil yang diperoleh, dapat dilihat bahwa pendekatan paling baik ditunjukkan oleh metode MLR dan MA, yakni dengan nilai MSE sebesar 15087.465, RMSE sebesar 122.831, dan MAPE sebesar 3.255. Jika pada penelitian sebelumnya K-Means menunjukkan hasil yang baik dalam memperbaiki kinerja regresi linear untuk mengoreksi warna citra, namun hal ini tidak ditunjukkan untuk kasus prediksi. Hal ini disebabkan K-Means tidak dapat menangani data *outlier* yang ada pada data saham. Lebih lanjut, metode prediksi *Improved Multiple Linear Regression* ini dapat digunakan untuk memprediksi nilai saham mingguan atau bulanan. Metode ini juga dapat digunakan dalam pembuatan aplikasi prediksi berbasis desktop, web, maupun *mobile*.

5. Daftar Notasi

- Y : Nilai target saham
- \hat{Y} : Hasil prediksi saham
- X_i : Data saham pada periode ke- i
- β : Koefisien regresi MLR
- X^T : Matriks transpose dari Matriks X
- X^{-1} : Matriks invers dari Matriks X
- d : Jarak euclidian
- x_i^j : Titik data ke- i
- c : Titik *centroid*.
- n : Jumlah titik data
- \bar{X}_n : Rataan bergerak pada periode ke- n

Referensi

- [1] D. Hatidja, "Penerapan Model Arima Untuk Memprediksi Harga Saham PT Telkom Tbk.," *Jurnal Ilmiah Sains*, Vol. 11, No. 1, Pp. 116–123, 2011.
- [2] S. C. Jaya, M.L. Khodra, "Model Prediksi Harga Saham Dengan Jaringan Syaraf Tiruan (Studi Kasus: Saham Tikm Di Bursa Efek Indonesia)," *Prosiding Konferensi Nasional Informatika (KNIF)*, 2015, Pp. 94–99.
- [3] A. Rahmi, W. Mahmudy, "Prediksi Harga Saham Berdasarkan Data Historis Menggunakan Model Regresi Yang Dibangun Dengan Algoritma Genetika," *JTIK Univ. Brawijaya*, Vol. 5, No. 12, Pp. 1–9, 2014.
- [4] L. Kurniawati, H. Tjandrasa, I. Arieshanti, "Prediksi Pergerakan Harga Saham Menggunakan Support Vector Regression," *Jurnal Scan*, Vol. 8, No. 2, Pp. 11–21, 2013.
- [5] Y. Sari, R. Ginardi, N. Suciati, "Color Correction Using Improved Linear Regression Algorithm," In *International Conference On Information, Communication Technology And System (ICTS) Proceeding*, 2015.
- [6] E. Han, A. Srivastava, V. Kumar, "Parallel Formulation Of Inductive Classification Learning Algorithm," 1996.
- [7] J. Heizer, B. Render, "Manajemen Operasi, Edisi 9 ", Terjemahan Chriswan Sungkono, In *Salemba Empat*, Jakarta, 2006.
- [8] F. L. H. Prasetya, "Manajemen Operasi," *Cetakan Pertama*, In *PT Buku Kita*, Jakarta, 2009.
- [9] R. Zunaidhi, W. Saputra, N. Sari, "Aplikasi Peramalan Penjualan Menggunakan Metode Regresi Linier," *Jurnal Scan*, Vol. 2, No. 3, Pp. 41–45, 2008.
- [10] T. Kanungo, D. Mount, N. Netanyahu, "An Efficient K-Means Clustering Algorithm: Analysis And Implementation," *IEEE Transactions on Pattern Analysis And Machine Intelligence*, Vol. 24, No. 7, 2002.