

Reduksi Data Latih pada K-Support Vector Nearest Neighbor Menggunakan Entropy

Eko Prasetyo¹, R. Dimas Adityo²

Program Studi Teknik Informatika, Fakultas Teknik,
Universitas Bhayangkara Surabaya
Jl. Ahmad Yani 114, Surabaya

¹eko@ubhara.ac.id, ²dimas@ubhara.ac.id

Nanik Suciati³, Chastine Fatichah⁴

Departemen Informatika, Fakultas Teknologi Informasi dan
Komunikasi

Institut Teknologi Sepuluh Nopember
Kampus ITS, Jl. Raya ITS, Sukolilo, Surabaya

³nanik@if.its.ac.id, ⁴chastine@if.its.ac.id

Abstrak—Pemilihan sebagian data latih atau reduksi data latih yang mempunyai pengaruh pada garis keputusan klasifikasi penting dilakukan. Tujuannya untuk mengurangi beban sistem pada tahap pelatihan. Sebagai metode reduksi data, K-Support Vector Nearest Neighbour (K-SVNN) mendapatkan hasil berdasarkan ketinggian nilai Significant Degree (SD) masing-masing data. Nilai SD dihitung menggunakan variabel LVRV (Left Value dan Right Value). Sayangnya, LVRV hanya dapat digunakan pada kasus klasifikasi biner. Penelitian ini melakukan uji coba penggunaan Entropy untuk menghitung SD. Secara konseptual, Entropy memberikan nilai kemurnian distribusi kelas data sehingga dimungkinkan penggunaan Entropy untuk menghitung SD pada kasus multi kelas. Pada makalah ini, disajikan analisis perbandingan perilaku nilai SD antara menggunakan LVRV dan Entropy. Hasil reduksi data menggunakan threshold (T) > 0 , didapatkan akurasi yang sama pada kedua metode, sedangkan klasifikasi dengan reduksi data latih memberikan nilai akurasi lebih tinggi daripada tanpa reduksi. Hal ini membuktikan bahwa entropy dapat digunakan untuk menggantikan LVRV untuk menghitung SD.

Kata kunci—reduksi data, entropy, left value right value, Nearest Neighbor, jenis pohon mangga

I. PENDAHULUAN

Pemilihan sebagian data latih yang digunakan pada proses pelatihan merupakan salah satu pekerjaan dalam pra-pemrosesan data. Data terpilih tersebut biasanya merupakan data yang mempunyai pengaruh pada garis batas keputusan (*decision boundary*) klasifikasi. Alasan dilakukan reduksi data ini adalah memberikan keuntungan pada kebanyakan metode daripada hanya didasarkan pada tetangga terdekat [1]. Alasan lainnya adalah mengurangi beban sistem pada tahap pelatihan karena ada sebagian data yang sebenarnya tidak mempunyai pengaruh pada garis batas keputusan klasifikasi. Jika data tersebut telah disisihkan, maka hanya data yang mempunyai pengaruh pada garis batas keputusan saja yang dibaca pada tahap pelatihan. Hasilnya, proses pelatihan menjadi lebih ringan dan cepat.

Reduksi data latih bertujuan untuk memilih sebagian sampel set data latih semula, mempunyai kemampuan memilih sampel yang relevan dan membuang data *noise* atau redundan tanpa membangkitkan data buatan baru yang sering disebut

metode Prototype Generation atau abstraksi [2]. Banyak metode reduksi data yang diusulkan oleh peneliti, diantaranya Condensed Nearest Neighbour Rule (CNN) [3] merupakan metode paling tua yang diusulkan, CNN melakukan reduksi data dengan mencari sebagian data dari data latih menjadi anggota hasil reduksi dimana jarak masing-masing anggota pada kelas yang sama lebih pendek daripada jarak pada anggota kelas lainnya, Ullmann's CNN [4] berusaha meningkatkan CNN dengan konsep similaritas dan matrik biner, Template Reduction KNN (TRKNN) [5] mengenalkan konsep rantai tetangga terdekat pada pencarian data hasil reduksi, dan K-Support Vector Nearest Neighbour (K-SVNN) [6] menggunakan konsep pemanggilan sebagai tetangga terdekat. Konsep K-SVNN menggunakan variabel LV (Left Value) dan RV (Right Value). LV digunakan untuk menampung jumlah pemanggilan data sebagai tetangga terdekat dari kelas yang sama, sedangkan RV untuk menampung jumlah pemanggilan dari kelas berbeda. Jika jumlah kelas set data lebih dari tiga, maka K-SVNN tidak dapat digunakan. Untuk menyelesaikan masalah tersebut, penelitian sebelumnya [7] mengusulkan Entropy pada perhitungan SD. Alasannya Entropy dapat mengukur kemurnian distribusi kelas data, sehingga pemilihan SD dapat didasarkan pada entropy tinggi. Tetapi penelitian tersebut belum membuktikan penggunaan Entropy untuk menghitung SD pada kasus klasifikasi biner.

Makalah ini memaparkan analisis perbandingan perilaku nilai SD antara penggunaan LVRV (Left Value dan Right Value) dan Entropy pada kasus klasifikasi biner. Analisis ini dilakukan untuk mengetahui perbedaan perilaku nilai SD yang dihasilkan. Pada kasus klasifikasi jenis pohon mangga, penulis menggunakan jenis mangga Madu dan Kepodang, masing-masing 100 daun. Masing-masing data direpresentasikan dengan 260 fitur.

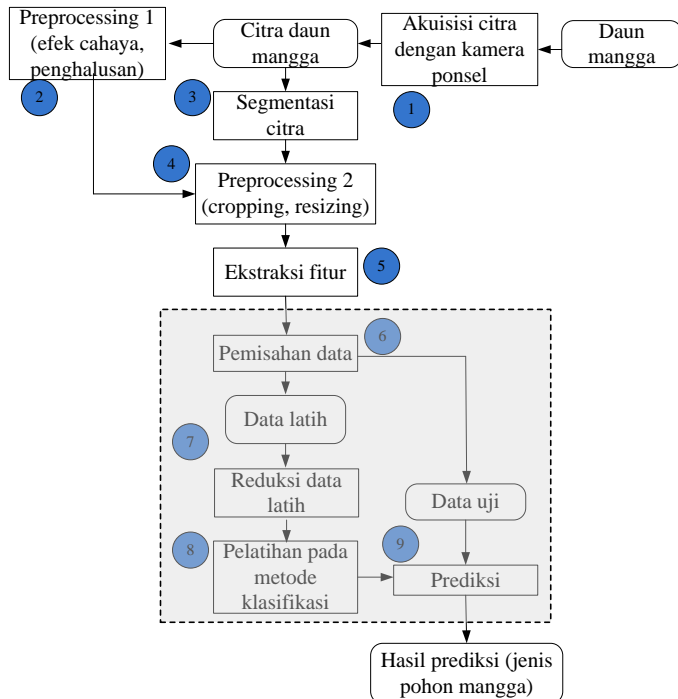
Makalah ini terbagi menjadi empat bagian. Bagian pertama berisi pendahuluan, bagian ini menyajikan latar belakang masalah dan penelitian lain yang terkait. Bagian kedua adalah metodologi penelitian, bagian ini menyajikan langkah-langkah penelitian deteksi jenis pohon mangga dan penelitian terkait permasalahan reduksi data. Bagian tiga adalah hasil dan

pembahasan, bagian ini menyajikan data-data dan analisis hasil pengujian. Bagian empat menyajikan simpulan dan saran.

II. METODOLOGI PENELITIAN

A. Kerangka Kerja Penelitian Klasifikasi Jenis Pohon Mangga

Kerangka kerja pada penelitian klasifikasi jenis pohon mangga terbagi menjadi beberapa tahap sebagai berikut: (1) akuisisi citra, (2) pemrosesan awal 1, (3) segmentasi, (4) pemrosesan awal 2, (5) ekstraksi fitur, (6) pemisahan data, (7) reduksi data latih, (8) pelatihan menggunakan metode klasifikasi, dan (9) prediksi. Diagram kerangka kerja disajikan pada Gambar 1.



Gambar 1. Kerangka kerja penelitian klasifikasi jenis pohon mangga

Tahap akuisisi citra merupakan tahap untuk mendapatkan citra daun mangga menggunakan kamera ponsel. Pada penelitian ini, citra dihasilkan dari akuisisi satu daun mangga dengan latar belakang lantai halaman. Resolusi citra yang digunakan adalah 2592 x 1944 piksel dengan efek normal. Pemrosesan awal digunakan untuk memperbaiki kualitas citra dan mengurangi bagian-bagian citra yang tidak dibutuhkan. Pemrosesan awal yang digunakan diantaranya adalah: penghalusan, pemrosesan morfologi, penghapusan bagian daun yang mendapat pengaruh cahaya tinggi [8], pemotongan bagian daun citra dan pengubahan ukuran citra agar sesuai dengan kebutuhan dalam komputasi. Segmentasi digunakan sebagai tahap untuk memisahkan obyek daun mangga dari latar belakang. Segmentasi menggunakan metode Otsu pada komponen Cr dari ruang warna YCbCr [9]. Tahap ekstraksi fitur merupakan tahap untuk membangkitkan fitur klasifikasi. Ada 260 fitur yang dibangkitkan, terdiri dari 256 fitur tekstur Weighted Rotation- and Scale invariant Local Binary Pattern Average (WRSI-LBP-avg) [10], 2 fitur warna yaitu rata-rata

dan standard deviasi [11], dan 2 fitur bentuk (compactness dan circularity). Tahap berikutnya adalah pemisahan data hingga prediksi. Tahap ini menjadi fokus yang yang dipaparkan dalam makalah ini, seperti yang disajikan pada Gambar 1, dibatasi oleh garis putus-putus. Didalamnya ada tahap reduksi data latih, pada tahap ini dilakukan pengurangan pada data latih yang kemungkinan besar tidak dibutuhkan pada proses pelatihan. Hal ini sebabkan data tersebut lokasinya jauh dari garis batas keputusan (*decision boundary*) klasifikasi. Pada makalah ini dipaparkan hasil perbandingan perhitungan SD antara penggunaan LVRV (Left Value dan Right Value) dan Entropy. Perbandingan diterapkan pada kasus klasifikasi biner, untuk kasus penelitian klasifikasi jenis mangga penulis menggunakan mangga jenis Madu dan Kepodang.

B. Reduksi Data dengan K-Support Vector Nearest Neighbour

Hasil penelitian yang disajikan pada makalah ini menggunakan 200 daun mangga yang direpresentasikan 260 fitur. Ada 2 jenis mangga yang digunakan dalam penelitian ini, yaitu: Madu, dan Kepodang. Perbandingan LVRV dan Entropy pada KSVNN dilakukan untuk kasus klasifikasi biner. Karena ukuran jumlah data dan fitur besar maka dibutuhkan tahap reduksi data latih. Tujuannya untuk mengurangi waktu komputasi dan kesederhanaan sistem. Banyak metode reduksi data yang diusulkan oleh peneliti sebelumnya, diantaranya adalah Condensed Nearest Neighbour Rule (CNN) [3], Template Reduction KNN (TRKNN) [5], dan K-Support Vector Nearest Neighbour (K-SVNN) [6]. Metode CNN melakukan reduksi data dengan mencari sebagian data dari data latih menjadi anggota hasil reduksi dimana jarak masing-masing anggota pada kelas yang sama lebih pendek daripada jarak pada anggota kelas lainnya [3]. Metode TRKNN menggunakan konsep rantai tetangga terdekat untuk mengurangi data latih [5]. Sedangkan K-SVNN mengeluarkan data yang tidak punya pengaruh pada garis batas keputusan (*decision boundary*) klasifikasi [5].

Data hasil reduksi K-SVNN didapatkan dari seleksi nilai *Significant Degree* (SD) setiap data. Data dengan nilai SD nol berarti data tersebut tidak mempunyai pengaruh sama sekali pada garis batas keputusan klasifikasi. Semakin besar nilai SD maka semakin tinggi pengaruh data latih tersebut pada garis batas keputusan. Pada penelitian sebelumnya [7, 12], seleksi data latih menggunakan nilai $SD > 0$, artinya sekecil apapun pengaruh data latih tetap digunakan sebagai hasil reduksi. Nilai SD didapatkan dari Left Value (LV) dan Right Value (RV). LV merupakan jumlah tetangga yang memanggil data tersebut sebagai tetangga terdekat pada kelas yang sama, sedangkan RV pada kelas berbeda [12]. Persamaan yang digunakan untuk menghitung SD pada data ke-*i* seperti disajikan pada Persamaan (1).

$$SD_i = \begin{cases} 0 & , LV_i = RV_i = 0 \\ \frac{LV_i}{RV_i} & , LV_i < RV_i \\ \frac{RV_i}{LV_i} & , LV_i > RV_i \\ 1 & , LV_i = RV_i \end{cases} \quad (1)$$

Konsep LV dan RV tersebut berakibat K-SVNN hanya dapat bekerja pada klasifikasi dua kelas saja. Alasannya adalah tidak ada tempat untuk menyimpan hitungan sebagai tetangga terdekat dari kelas ketiga, keempat, dan seterusnya. Maka, dalam makalah hasil penelitian ini dipaparkan uji coba penggunaan Entropy untuk menghitung SD. Entropy dapat mengukur kemurnian distribusi kelas data, sehingga pemilihan SD dapat dilakukan berdasarkan nilai Entropy tinggi. Dengan menggunakan Entropy, maka LV dan RV setiap data digantikan oleh nilai $V_i(k)$, $k=1, \dots, C$, dimana $V_i(k)$ adalah jumlah pemanggilan sebagai tetangga terdekat kelas k pada data ke- i . Sedangkan C adalah jumlah kelas. Untuk menentukan nilai $V_i(k)$ digunakan (2).

$$V_i(k) = \sum_{j=1}^N I(i, k, j) \quad (2)$$

Dimana $I(i, k, j)$ adalah hasil pemeriksaan ketika data ke- i terpanggil sebagai K tetangga terdekat data ke- j , sedangkan data ke- i mempunyai kelas k . Hal tersebut dilakukan ketika data ke- j mencari tetangga terdekat dan data ke- i terpanggil sebagai salah satu dari K tetangga terdekat. Jika k merupakan kelas data ke- i ternyata sama dengan kelas data ke- j maka $V_i(k)$ dinaikkan 1, jika tidak sama maka $V_i(k)$ pada semua kelas selain k dinaikkan 1. Seperti disajikan pada (3) dan (4).

$$I(i, k, j) = \begin{cases} 1, & C_i(k) = C(j) \\ 0, & \text{lainnya} \end{cases} \quad (3)$$

$$I(i, \sim k, j) = \begin{cases} 1, & C_i(\sim k) = C(j) \\ 0, & \text{lainnya} \end{cases} \quad (4)$$

Dimana $C(j)$ adalah label kelas data ke- j yang sedang diproses. $C_i(k)$ adalah label kelas data ke- i . $I(i, \sim k, j)$ adalah hasil pemeriksaan data ke- i untuk kelas selain k ketika terpilih sebagai tetangga terdekat oleh data ke- j . Sebelum menghitung SD, nilai V pada setiap data dinormalisasikan menggunakan (5).

$$V_i^{norm}(k) = \frac{V_i(k)}{\sum_{k=1}^C V_i(k)} \quad (5)$$

Nilai SD dihitung menggunakan nilai Entropy ternormalisasi seperti pada (6).

$$SD_i = entropy_i = - \sum_{k=1}^C \left(V_i^{norm}(k) \times \log_2 \left(V_i^{norm}(k) \right) \right) \quad (6)$$

III. HASIL DAN PEMBAHASAN

A. Set Data dan Skenario Pengujian

Penulis melakukan perbandingan antara LVRV dan Entropy ketika digunakan untuk menghitung SD. Perhitungan LVRV menggunakan persamaan dalam penelitian sebelumnya [12], sedangkan Entropy menggunakan (2) hingga (6). Set data yang digunakan dalam penelitian ini dibangkitkan dari 200 citra daun mangga yang terbagi menjadi dua jenis yaitu Madu dan Kepodang. Masing-masing jenis daun terdapat 100 citra. Masing-masing citra direpresentasikan oleh 260 fitur, terdiri dari 256 fitur tekstur Weighted Rotation- and Scale Invariant Local Binary Pattern Average (WRSI-LBP-avg) [10], 2 fitur warna yaitu rata-rata dan standard deviasi [11], dan 2 fitur bentuk yaitu compactness dan circularity. Penulis menggunakan metode K-Fold Cross Validation dengan proporsi 50:50, yaitu 50% sebagai data latih yang akan dilakukan reduksi dan 50% sebagai data uji untuk prediksi.

Data yang disajikan dari hasil perbandingan ini adalah:

1. Pola nilai SD yang terbentuk

Rentang nilai SD yang dihasilkan oleh LVRV adalah 0 hingga 1, sedangkan Entropy mulai dari 0. Dengan mengurutkan nilai SD semua data kemudian digambar menjadi grafik maka dapat diketahui bentuk pola nilai SD. Hal ini penting untuk mengetahui perkiraan nilai batas SD yang lolos seleksi sebagai hasil reduksi. K tetangga terdekat yang digunakan dalam pengujian ini adalah 5, 9, dan 11.

2. Prosentase pengurangan data latih

Metrik ini dapat digunakan untuk mengukur tingkat reduksi yang dihasilkan berdasarkan parameter K dan nilai batas SD. Semakin tinggi reduksi maka semakin sedikit data yang lolos sebagai hasil reduksi. Selanjutnya, proses pelatihan juga semakin ringan karena jumlah data yang diproses semakin sedikit. Metrik ini dipengaruhi oleh pemilihan nilai batas SD, semakin tinggi nilai batas maka jumlah data yang lolos sebagai hasil reduksi semakin sedikit, begitu pula sebaliknya.

3. Akurasi prediksi

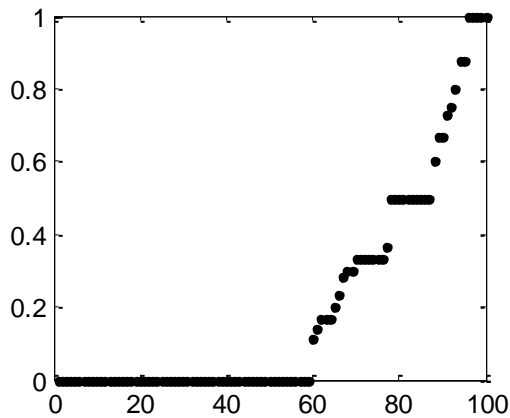
Prediksi merupakan tujuan akhir tahap klasifikasi. Semakin tinggi akurasi yang diberikan maka kinerja *classifier* semakin baik. Maka metrik akurasi perlu digunakan dalam pengujian untuk mengetahui pengaruh hasil reduksi pada hasil kerja *classifier*.

B. Hasil Pengujian

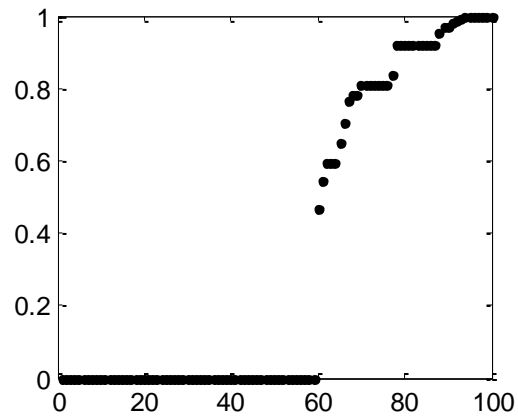
Penulis melakukan pengujian dengan membandingkan nilai SD antara menggunakan LVRV dan Entropy. Dengan menggunakan 2-Fold Cross Validation pada 200 data maka ada 100 data latih yang diproses. Nilai SD yang didapatkan kemudian diurutkan, hasilnya disajikan pada grafik seperti pada Gambar 2. Sumbu x pada Gambar 2 menyatakan nomor data yang sudah diurutkan nilai SD dari kecil ke besar. Sumbu y menyatakan nilai SD pada setiap data. Pada grafik di Gambar 2 (a), (c), dan (e) tersebut dapat diamati bahwa pola nilai SD yang dihitung menggunakan LVRV secara umum membentuk garis berbanding lurus dan proporsional antara urutan data dan nilai SD. Hal ini disebabkan persamaan untuk menghitung SD adalah persamaan linear. Sedangkan pola nilai

SD yang dihitung menggunakan Entropy pada Gambar 2 (b), (d), dan (f) membentuk garis lengkung. Bentuk garis ini merupakan pola grafik nilai yang didapatkan dengan persamaan non linear.

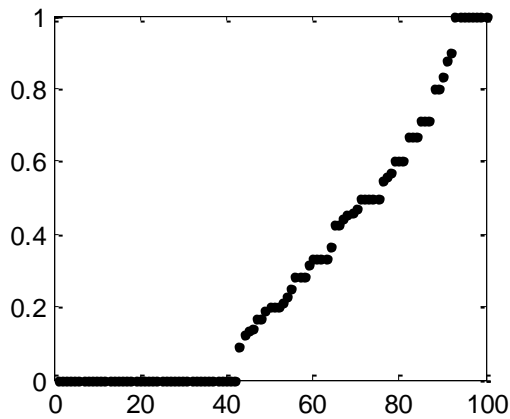
Efek yang diakibatkan dari bentuk pola nilai SD pada kedua metode ini adalah nilai batas SD (T) yang digunakan untuk seleksi data latih menjadi support vector. Penelitian sebelumnya menggunakan nilai $T > 0$. Pada gambar tersebut dapat diamati bahwa nilai SD pada metode LVRV beragam mulai dari 0 hingga 1, sedangkan pada metode Entropy secara umum mulai 0.2 hingga 1. Berdasarkan bentuk grafik maka nilai T yang diberikan pada kedua metode bisa berbeda jika menginginkan hasil yang sama.



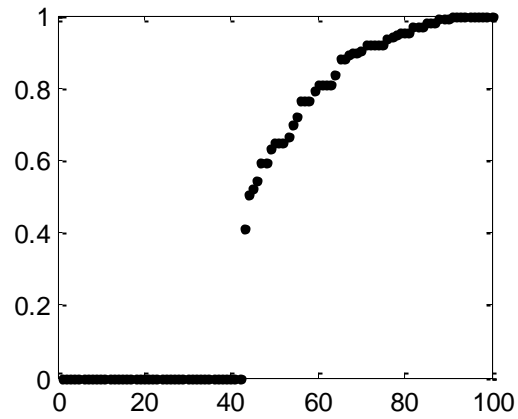
a. LVRV, K=5



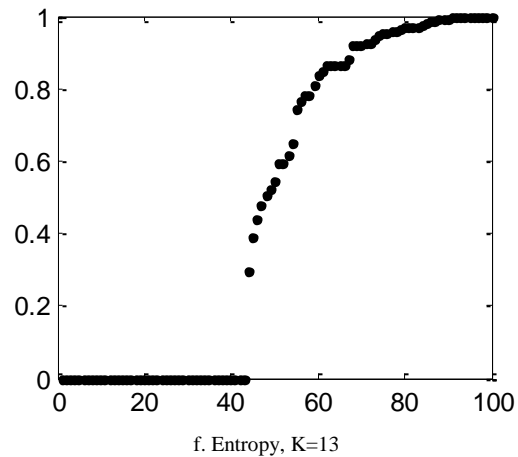
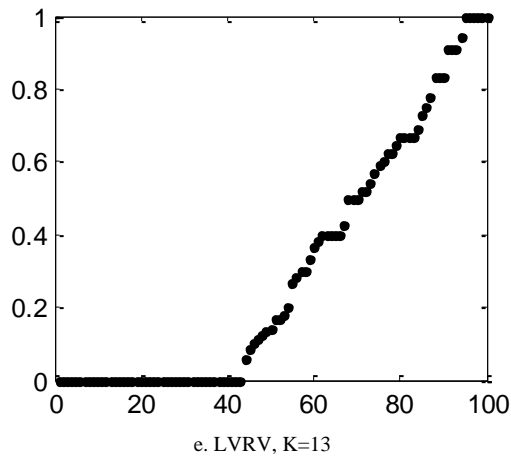
b. Entropy, K=5



c. LVRV, K=9



d. Entropy, K=9



Gambar 2. Pola nilai SD 100 data latih (ditampilkan secara urut naik); (a) (c) (e) Hasil SD menggunakan LVRV; (b) (d) (f) Hasil SD menggunakan Entropy

Data hasil reduksi seperti pada Gambar 2 juga diuji menggunakan metode klasifikasi Support Vector Machine (SVM) dengan kernel Linear. Data uji yang digunakan adalah 100 data lainnya. Pengujian menggunakan 2-Fold Cross Validation. Penulis juga membandingkannya dengan klasifikasi SVM tanpa reduksi data latih. Hasil yang dicapai disajikan pada Tabel 1. Dari nilai yang disajikan pada tabel tersebut, secara umum nilai pada semua parameter untuk LVRV dan Entropy selalu sama. Hal ini bahwa pada parameter yang sama yang digunakan dalam penelitian ini, kedua metode memberikan hasil identik, tidak ada hasil berbeda. Ini sejalan dengan tujuan memperluas penggunaan K-SVNN pada kasus multi kelas yaitu hasil yang sama bisa dicapai ketika LVRV digantikan oleh Entropy. LVRV terbatas hanya untuk dua kelas sedangkan Entropy dapat digunakan untuk multi kelas. Reduksi yang didapatkan mempunyai pola turun sejalan dengan semakin tingginya K yang digunakan. Ini adalah sifat yang dimiliki oleh K-SVNN seperti disampaikan pada penelitian sebelumnya [6]. Sedangkan akurasi yang didapatkan oleh SVM dengan reduksi lebih tinggi daripada tanpa reduksi. Akurasi tertinggi didapatkan dengan reduksi 0.5 sedangkan tanpa reduksi 0.48.

TABEL I. PERBANDINGAN KINERJA

K	Keterangan	LVRV	Entropy
5	Reduksi	0.7875	0.7875
	Akurasi dengan reduksi	0.5	0.5
	Akurasi tanpa reduksi	0.48	0.48
9	Reduksi	0.7025	0.7025
	Akurasi dengan reduksi	0.5	0.5
	Akurasi tanpa reduksi	0.43	0.43
13	Reduksi	0.685	0.685
	Akurasi dengan reduksi	0.5	0.5
	Akurasi tanpa reduksi	0.48	0.48

IV. SIMPULAN DAN SARAN

Simpulan yang didapatkan dari hasil penelitian ini adalah hasil reduksi data menggunakan threshold (T) > 0 , didapatkan akurasi yang sama pada kedua metode, sedangkan klasifikasi dengan reduksi data latih memberikan nilai akurasi lebih tinggi daripada tanpa reduksi. Hal ini membuktikan bahwa entropy dapat digunakan untuk menggantikan LVRV untuk menghitung SD.

Saran yang dapat diberikan dari penelitian ini adalah fitur yang digunakan belum dilakukan seleksi untuk mendapatkan fitur yang informative sekaligus mengurangi kompleksitas komputasi karena jumlah fitur terlalu besar.

UCAPAN TERIMA KASIH

Penulis mengucapkan terima kasih pada Direktorat Riset dan Pengabdian Masyarakat (DRPM) DIKTI yang memberikan pendanaan penelitian pada skim Penelitian Kerjasama Antar Perguruan Tinggi (PKPT) tahun 2018 antara Universitas Bhayangkara Surabaya dan Institut Teknologi Sepuluh Nopember dengan nomor kontrak penelitian: 009/SP2H/LT/K7/KM/2018 tanggal 26 Pebruari 2018.

REFERENSI

- [1] A. Arnaiz-González, Díez-Pastor, J.F., Rodríguez, J.J., García-Osorio, C., "Instance selection of linear complexity for big data," *Knowledge-Based Systems*, vol. 107, pp. 83-95, 2016.
- [2] G. Ditzler, Roveri, M., Alippi, C., "Learning in Nonstationary Environments: A Survey," *IEEE Computational Intelligence Magazine* vol. 10, pp. 12-25, 2015.
- [3] P. E. Hart, "The condensed nearest neighbour rule," *IEEE Transactions on Information Theory*, vol. 18, pp. 515-516, 1968.
- [4] G. Gates, "The reduced nearest neighbor rule (Corresp.)," *IEEE Transactions on Information Theory*, vol. 18, 1972.
- [5] H. Fayed, Atiya, A., "A novel template reduction approach for the k-nearest neighbor method," *IEEE Transactions on Neural Networks*, vol. 20, pp. 890-896, 2009.
- [6] E. Prasetyo, "K-Support Vector Nearest Neighbor Untuk Klasifikasi Berbasis K-NN," in *Seminar Nasional Sistem Informasi Indonesia*, Surabaya, 2012, pp. 245-250.

- [7] E. Prasetyo, Adityo, R.D., Suciati, N., and Faticah, C., "Multi-class K-Support Vector Nearest Neighbor for Mango Leaf Classification," *Telkomnika*, vol. Accepted Paper, 2018.
- [8] E. Prasetyo, Adityo, R.D., Suciati, N., and Faticah, C., "Deteksi Wilayah Cahaya Intensitas Tinggi Citra Daun Mangga Untuk Ekstraksi Fitur Warna dan Tekstur Pada Klasifikasi Jenis Pohon Mangga," presented at the Seminar Nasional Teknologi Informasi, Komunikasi dan Industri, UIN Sultan Syarif Kasim Riau, 2017.
- [9] E. Prasetyo, Adityo, R.D., Suciati, N., and Faticah, C., "Mango Leaf Image Segmentation on HSV and YCbCr Color Spaces Using Otsu Thresholding," in *International Conference on Science and Technology*, Yogyakarta, 2017.
- [10] E. Prasetyo, Adityo, R.D., Suciati, N., and Faticah, C., "Average and Maximum Weights in Weighted Rotation- and Scale-invariant LBP for Classification of Mango Leaves," *Journal of Theoretical and Applied Information Technology*, 2017.
- [11] Fazal-E-Malik, "Mean and Standard Deviation Features of Color Histogram using Laplacian Filter for Content-based Image Retrieval," *Journal of Theoretical and Applied Information Technology*, vol. 34, pp. 1-7, 2011.
- [12] E. Prasetyo, "Reduksi Data Latih Dengan K-SVNN Sebagai Pemrosesan Awal Pada ANN Back-Propagation Untuk Pengurangan Waktu Pelatihan," *SIMETRIS*, vol. 6, pp. 223-230, 2015.