

## CLASSIFICATION OF LIDAR DATA OVER BUILDING ROOFS USING K-MEANS AND PRINCIPAL COMPONENT ANALYSIS

### *Classificação de dados LiDAR sobre telhados de edificações usando k-médias e análise de componentes principais*

Renato César dos Santos<sup>1</sup> - ORCID: 0000-0003-0263-312X

Mauricio Galo<sup>1, 2</sup> - ORCID: 0000-0002-0104-9960

Vilma Mayumi Tachibana<sup>1, 3</sup> - ORCID: 0000-0002-8804-6163

<sup>1</sup> Universidade Estadual Paulista Júlio de Mesquita Filho - UNESP, Programa de Pós-Graduação em Ciências Cartográficas - PPGCC, Presidente Prudente - SP, Brasil.

E-mail: renato\_cstos@hotmail.com

<sup>2</sup> Universidade Estadual Paulista Júlio de Mesquita Filho – UNESP, Departamento de Cartografia, Presidente Prudente - SP, Brasil.

E-mail: galo@fct.unesp.br

<sup>3</sup> Universidade Estadual Paulista Júlio de Mesquita Filho – UNESP, Departamento de Estatística, Presidente Prudente - SP, Brasil.

E-mail: vilma@fct.unesp

Received in March 13<sup>th</sup>, 2017.

Accepted in January 16<sup>th</sup>, 2018.

#### **Abstract:**

The classification is an important step in the extraction of geometric primitives from LiDAR data. Normally, it is applied for the identification of points sampled on geometric primitives of interest. In the literature there are several studies that have explored the use of eigenvalues to classify LiDAR points into different classes or structures, such as corner, edge, and plane. However, in some works the classes are defined considering an ideal geometry, which can be affected by the inadequate sampling and/or by the presence of noise when using real data. To overcome this limitation, in this paper is proposed the use of metrics based on eigenvalues and the k-means method to carry out the classification. So, the concept of principal component analysis is used to obtain the eigenvalues and the derived metrics, while the k-means is applied to cluster the roof points in two classes: edge and non-edge. To evaluate the proposed method four test areas with different levels of complexity were selected. From the qualitative and quantitative analyses, it could be concluded that the proposed classification procedure gave satisfactory results, resulting in completeness and correctness above 92% for the non-edge class, and between 61% to 98% for the edge class.

**Keywords:** Classification of LiDAR points; Edge points; K-means method; Principal component analysis; Eigenvalues.

**How to cite this article:** dos Santos, R.C.; et al Classification of Lidar data over building roofs using k-means and principal component analysis. *Bulletin of Geodetic Sciences*, Vol. 24, issue 1, 69-84, Jan-Mar, 2018.



This content is licensed under a Creative Commons Attribution 4.0 International License.

**Resumo:**

A classificação é uma importante etapa na extração de primitivas geométricas sobre dados LiDAR. Normalmente, a classificação é utilizada para identificar os pontos amostrados sobre primitivas de interesse. Na literatura são encontrados vários trabalhos que exploram o uso dos autovalores para classificar os pontos LiDAR em diferentes estruturas ou classes, tais como: quina, borda, e plano. Entretanto, alguns trabalhos desenvolvidos se baseiam em parâmetros obtidos a partir de uma geometria ideal, que pode fornecer resultados não adequados quando a amostragem for insuficiente ou quando da presença de ruídos. Para contornar esta limitação, é proposto o uso de métricas estimadas a partir de autovalores e do uso do método k-médias. O conceito de análise de componentes principais é utilizado para determinar os autovalores e algumas métricas derivadas, enquanto que o método k-médias é aplicado para agrupar os pontos de telhados em duas classes: borda e não borda. Para avaliar a metodologia foram selecionadas quatro áreas teste com diferentes níveis de complexidade. A partir dos resultados, foi possível concluir que o procedimento de classificação apresentou resultados satisfatórios, obtendo-se nível de acerto e completude acima de 92% para os pontos da classe não borda e entre 61% e 98% para a classe borda.

**Palavras-chave:** Classificação de pontos LiDAR; Pontos de borda; Método k-médias; Análise de componentes principais; Autovalores.

## 1. Introduction

The LASER scanning system installed in an airplane is basically composed of four technologies (El-Sheimy et al. 2005): a LASER emitter and receiver, a scanning system, a satellite positioning system, and an inertial system. These technologies, which must operate in a synchronized manner, generate integrated measurements enabling the determination of the 3D position of points sampled on the scanned surface. According to Wehr and Lohr (1999) this measuring system can be also called LiDAR (Light Detection And Ranging). In recent decades, LiDAR technology has been used in several sciences, for example in the geodesic sciences, where this principle has been applied in many works related to the extraction of objects, such as buildings (Machado and Mitishita 2006, Galvanin and Dal Poz 2012), trees (Xiao 2012), roads (Wang et al. 2011), power lines (Cruz and Silveira 2011), planar surfaces (Lari et al. 2017), among others.

In the context of automatic building extraction, obtaining straight line segments and planes are important tasks. The line segments can be related to ridge lines or building contours. According to Bretar (2009), straight line segments in 3D space can be extracted in two ways: the intersection between adjacent planes (Habib et al. 2005), or directly on LiDAR points (Gross and Thoennessen 2006). In this second case, the extraction is preceded by a classification process, in which the points sampled on edges and/or ridge lines are identified first. Many works in the literature explore the use of eigenvalues to identify LiDAR points sampled on different classes or structures (corner, edge, plane). In Gross and Thoennessen (2006), Jutzi and Gross (2009), Santos (2015), and Santos and Galo (2014, 2016) the classification is performed by comparing the estimated eigenvalues for each point with theoretical eigenvalues from each class, which are determined analytically, as suggested by Gross and Thoennessen (2006). The restriction, in this case, is related to the use of theoretical eigenvalues to describe the classes, which are obtained considering an

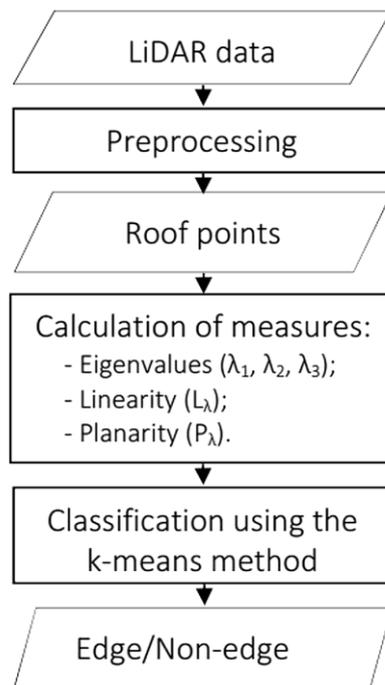
ideal situation and may not reflect the characteristics of real data set, especially in the case of edges. It may happen since it is not possible to guarantee that the edge points are sampled, depending on the scanning rate and point cloud density during the LiDAR scanning. In Sampath and Shan (2008, 2010) and He et al. (2012), the authors estimate a normalized eigenvalue for each point. The classification is then performed by comparing the normalized eigenvalue with a planarity threshold, which is determined empirically. The points are labeled in two classes: planar and non-planar. The limitation of this process is the empirical determination of the threshold, which can vary according to the characteristics of the data set.

To overcome the limitations mentioned, the use of the k-means method and exploration of different measures is proposed. K-means is a technique of cluster analysis, which automatically assigns objects to different classes by analyzing the similarities or dissimilarities between them. In this case, it is not necessary to know *a priori* the behavior of the classes or perform experiments to determine the threshold used to separate the classes. Another advantage is that the classes can be modeled according to object features that can be quantified by different measures. The only aspect that can be considered as disadvantageous is related to *a priori* definition of the number of clusters. The k-means method has been used by some authors to classify LiDAR points, but in different contexts. Chehata et al. (2008) employ the method to identify ground and non-ground points. In Sampath and Shan (2010), the method is applied to group the sampled points on different plane roofs. With regard to the measures used in the classification, the eigenvalues and measures derived from the eigenvalues (West et al. 2004, Pauly et al. 2003) as linearity and planarity were explored.

In this paper, a procedure is proposed for classification of LiDAR points sampled on roofs considering some techniques of multivariate analysis, such as principal component analysis and the k-means clustering method. Furthermore, the use of different sets of measures is explored in the context of classification. In general, the classification method can be divided into three main steps. In the first, a preprocessing of original set points using the *LAStools* software is performed to select LiDAR points sampled on building roofs. In the second, the concept of principal component analysis is used to determine the eigenvalues, and the linearity and planarity measures, related to each LiDAR point and its neighborhood. In the last step, the points are grouped into two classes: edge and non-edge points, using the k-means method.

## 2. Method

Figure 1 shows a flowchart of the proposed classification procedure divided into three main steps: preprocessing of the original point cloud, calculation of measures (eigenvalues, linearity and planarity) related to each LiDAR point, and classification of points into two classes: edge and non-edge points.



**Figure 1:** Flowchart of the steps involved in the classification of LiDAR points over building roofs.

## 2.1 Preprocessing

In this paper, the preprocessing step is performed to select LiDAR points sampled on building roofs. The selection is performed using some tools from *LAStools* software (<http://rapidlasso.com/lastools/>) and can be divided into two sub-steps.

In the first sub-step, the LiDAR points are separated into two groups: ground and non-ground points (buildings, trees, cars, etc). This separation can be conducted by implementing methods found in the literature, for example, the method based on the concept of mathematical morphology (Zhang et al. 2003, Carrilho 2016), and on the concept of slope (Vosselman 2000) or any other filtering approach. Since the focus of this paper is the classification of roof points, *LAStools* software was used, considering the default parameters. In this case, the filtering is run using the *lasground* tool, which is based on the progressive TIN densification filter (Zhang and Lin 2013).

In the second sub-step the sampled points on building roofs are selected using the *lasclassify* tool. For this purpose, the normalized digital surface model (nDSM) is generated from DTM (Digital Terrain Model) and DSM (Digital Surface Model) by the following operation:  $nDSM = DSM - DTM$ . This operation is performed using *lasheight* tool. From nDSM, the sampled points on roofs are selected by analyzing parameters such as height, planarity and roughness.

## 2.2 Selection of measures using principal component analysis

To perform the classification of a set of objects, it is necessary to define the measures that enable the separation of the classes of interest. The measures are used to identify standards related to

different classes. In the case of sampled points in 3D space, these standards can be identified in the variance-covariance matrix or in the correlation matrix, for example. In addition, it is also possible to use some statistical techniques that describe characteristics from the variance-covariance matrix, as in the case of factor analysis and principal components analysis, to determine others types of measures, for example, eigenvalues and/or eigenvectors (Johnson and Wichern 2007).

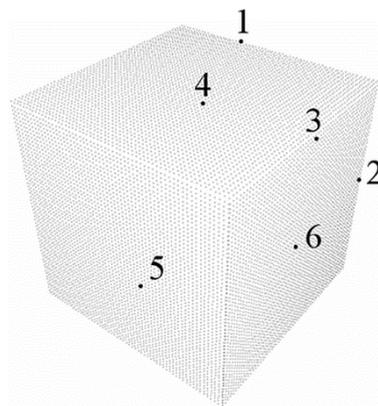
According to West et al. (2004) and Pauly et al. (2003) other measures can be derived from the eigenvalues ( $\lambda_1, \lambda_2, \lambda_3$ ). As an example, we can mention linearity ( $L_\lambda$ ) and planarity ( $P_\lambda$ ) measures. Assuming that the eigenvalues are available and ordered (as  $\lambda_1 \geq \lambda_2 \geq \lambda_3$ ), the linearity and planarity measures can be calculated as given by Equations 1 and 2, respectively.

$$L_\lambda = \frac{\lambda_1 - \lambda_2}{\lambda_1} \quad (1)$$

$$P_\lambda = \frac{\lambda_2 - \lambda_3}{\lambda_1} \quad (2)$$

In this paper, three different groups of measures were explored. The first is composed of eigenvalues ( $\lambda_1, \lambda_2, \lambda_3$ ), as performed by Gross and Thoennessen (2006), Yang and Dong (2013), and Santos and Galo (2016). The second is formed by linearity and planarity measures ( $L_\lambda, P_\lambda$ ), and the last is composed only by linearity measure ( $L_\lambda$ ).

Figure 2 and Table 1 are shown in order to visualize the effect of using these measures ( $\lambda_1, \lambda_2, \lambda_3, L_\lambda, P_\lambda$ ) on points belonging to the edge and non-edge classes. In Figure 2, thousands of points sampled on a cube surface are shown, where six of them are highlighted: three sampled on edge (1, 2 and 3) and three on non-edge regions (4, 5 and 6). The values of these measures, calculated considering a neighborhood around each point of interest, are shown in Table 1.



**Figure 2:** Sampled points on the cube surface. Points 1, 2 and 3 are sampled on edge, and points 4, 5 and 6 are located on non-edge regions.

The values given in Table 1 reinforce the idea of associating the relation between the eigenvalues discussed by Tong et al. (2004), Demantke et al. (2011), and Yang and Dong (2013), and the theoretical eigenvalues determined by Gross and Thoennessen (2006). The difference between the sampled points on edge and non-edge can also be observed analyzing the linearity and planarity measures.

**Table 1:** The values of the measures  $\lambda_1, \lambda_2, \lambda_3, L_\lambda, P_\lambda$  estimated for the highlighted points in Figure 2, considering the edge and non-edge classes.

		Measures				
	Points	$\lambda_1$	$\lambda_2$	$\lambda_3$	$L_\lambda$	$P_\lambda$
Edge	1	0.26	0.12	0.05	0.54	0.27
	2	0.23	0.11	0.00	0.52	0.43
	3	0.25	0.13	0.02	0.48	0.44
Non-Edge	4	0.24	0.24	0.00	0.00	1.00
	5	0.26	0.25	0.02	0.04	0.88
	6	0.23	0.22	0.01	0.04	0.91

The eigenvalues related to each LiDAR point can be obtained through the variance-covariance matrix. In this case, it is necessary to define a neighborhood around a point of interest. In this work, the neighborhood is determined considering a sphere with a radius R centered on a point of interest. The radius value is obtained automatically using the concept of entropy, as performed by Demantke et al. (2011), Yang and Dong (2013), and Santos and Galo (2016).

Assuming that at a certain LiDAR point in 3D space has N neighboring points; the correspondent variance-covariance matrix is calculated by the following equation:

$$S_x = \frac{1}{NR^2} \sum_{j=1}^N X_j X_j^T - m_x m_x^T \tag{3}$$

$$m_x = \frac{1}{N} \sum_{j=1}^N X_j \tag{4}$$

where:

$m_x$  is the vector of mean values (center of gravity);

R is the radius of the sphere;

X is the vector of the 3D coordinates (X, Y, Z).

One way of determining the eigenvalues is by using the singular value decomposition (SVD) method, as seen in Press et al. (1992). Considering the SVD concept, the variance-covariance matrix can be obtained through multiplication of the following matrices: eigenvectors matrix (composed by elements  $e_{i,j}$ , with i and j  $\in \{1, \dots, n\}$ ), diagonal matrix of eigenvalues ( $\lambda_i$ ), with i  $\in \{1, \dots, n\}$  and transpose matrix of eigenvectors, as can be seen in Equation 5.

$$\begin{bmatrix} s_{1,1}^2 & s_{1,2} & \dots & s_{1,n} \\ s_{2,1} & s_{2,2}^2 & \dots & s_{2,n} \\ \dots & \dots & \dots & \dots \\ s_{n,1} & s_{n,2} & \dots & s_{n,n}^2 \end{bmatrix} = \begin{bmatrix} e_{1,1} & e_{1,2} & \dots & e_{1,n} \\ e_{2,1} & e_{2,2} & \dots & e_{2,n} \\ \dots & \dots & \dots & \dots \\ e_{n,1} & e_{n,2} & \dots & e_{n,n} \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix} \begin{bmatrix} e_{1,1} & e_{2,1} & \dots & e_{n,1} \\ e_{1,2} & e_{2,2} & \dots & e_{n,2} \\ \dots & \dots & \dots & \dots \\ e_{1,n} & e_{2,n} & \dots & e_{n,n} \end{bmatrix} \tag{5}$$

### 2.3 Classification of LiDAR points using the k-means method

K-means is a non-hierarchical clustering method that aims to divide the objects into k clusters. Each object has a feature vector, which is compared with mean feature vectors (centroid) related

to each  $k$  clusters. Then the object is assigned to the cluster having the nearest centroid or distance (Johnson and Wichern 2007). Mathematically, the  $k$ -means method can be described as a clustering scheme that divides the input data into  $k$  clusters, in such a way that the following function is minimized:

$$\phi = \sum_{j=1}^k \sum_{i=1}^m d(\mathbf{F}_i, \bar{\mathbf{F}}_j) \quad (6)$$

where:

$\mathbf{F}_i$  is  $i$ <sup>th</sup> feature vector related to  $i$ <sup>th</sup> object;

$\bar{\mathbf{F}}_j$  is the mean feature vector of cluster  $j$  (cluster center  $j$ );

$d$  is the distance function;

$k$  is the number of clusters;

$m$  is the total number of data points.

In Equation 6 different metrics can be considered (Shan and Sampath 2009): Euclidean distance, Manhattan distance and Mahalanobis distance. In Equation 7 the formulation of Euclidean distance between two feature vectors,  $\mathbf{F}_i$  and  $\mathbf{F}_j$ , as considered in this paper, is shown:

$$d(\mathbf{F}_i, \bar{\mathbf{F}}_j) = \sqrt{(\mathbf{F}_i - \bar{\mathbf{F}}_j)^T (\mathbf{F}_i - \bar{\mathbf{F}}_j)} \quad (7)$$

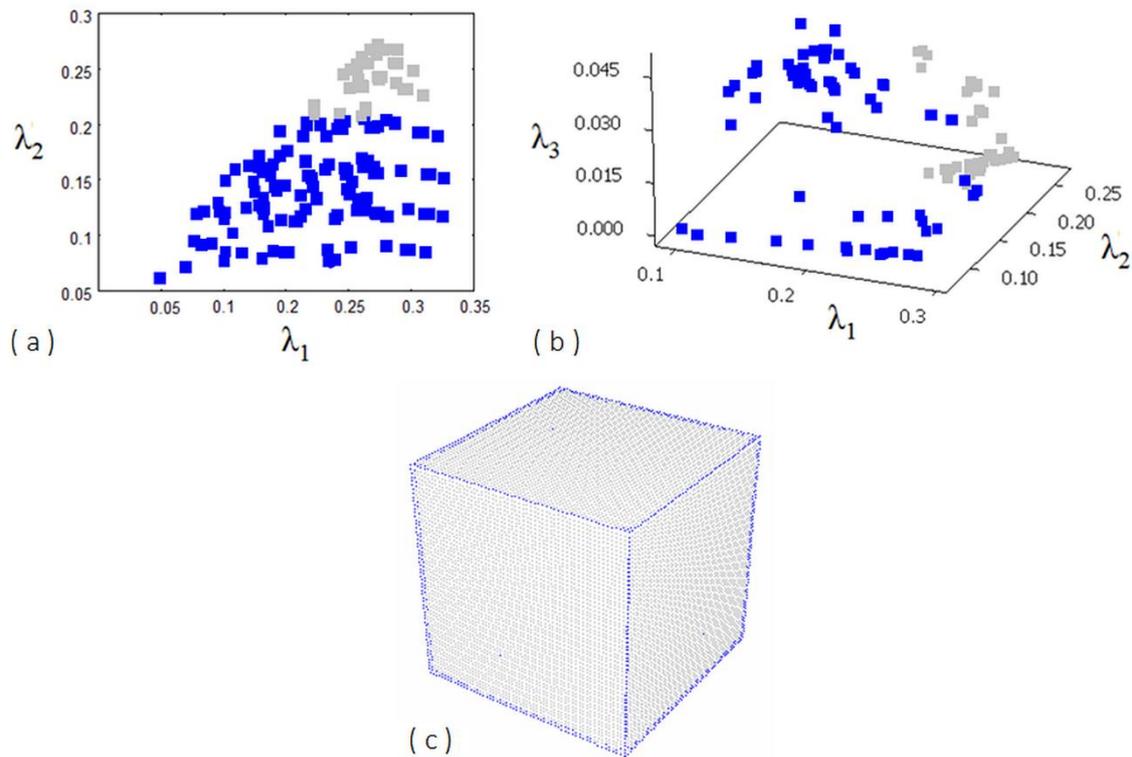
In this paper three different groups of measures were explored, as discussed in Subsection 2.2. Table 2 shows the composition of the attribute vector ( $\mathbf{F}_i$ ), and centroid ( $\bar{\mathbf{F}}_j$ ) related to each set of measures, where the index  $i$  represents a generic LiDAR point  $i$ , and  $j$  is a generic class.

**Table 2:** Feature vector and centroid related to each group of measures.

Measures	Feature vector ( $\mathbf{F}_i$ )	Centroid ( $\bar{\mathbf{F}}_j$ )
$\lambda_1, \lambda_2, \lambda_3$	$\mathbf{F}_i = [\lambda_1^{(i)}, \lambda_2^{(i)}, \lambda_3^{(i)}]^T$	$\bar{\mathbf{F}}_j = [\bar{\lambda}_1^{(j)}, \bar{\lambda}_2^{(j)}, \bar{\lambda}_3^{(j)}]^T$
$L_\lambda, P_\lambda$	$\mathbf{F}_i = [L_\lambda^{(i)}, P_\lambda^{(i)}]^T$	$\bar{\mathbf{F}}_j = [\bar{L}_\lambda^{(j)}, \bar{P}_\lambda^{(j)}]^T$
$L_\lambda$	$\mathbf{F}_i = [L_\lambda^{(i)}]^T$	$\bar{\mathbf{F}}_j = [\bar{L}_\lambda^{(j)}]^T$

Since the points sampled on roofs of buildings are predominantly classified into two classes: edge and non-edge, with most of the points labeled in the non-edge class, as can be seen in Santos and Galo (2014, 2016), the number of classes was defined as two.

Figure 3 shows the result of the  $k$ -means method applied to the cube shown in Figure 2. In this case, only the eigenvalues were used as measures. Figures 3a and 3b show the two generated clusters in 2D and 3D space, respectively, while Figure 3c shows the set of points sampled on the cube surface, after clustering.

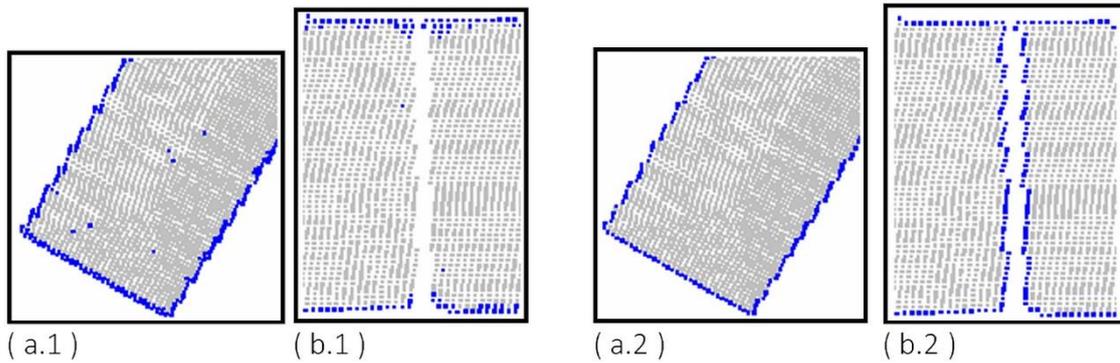


**Figure 3:** Result of the classification using the k-means method. Representation of the clusters in spaces 2D (a) and 3D (b), and set of points sampled on cube after classification (c).

## 2.4 Quality evaluation

The results derived from classification can be evaluated through qualitative and quantitative analysis. Qualitative evaluation was performed through visual analysis, basically by analyzing the consistency of the classification results. Aerial images of the same region were used to help this analysis. In the case of the numerical analysis, various quality parameters were calculated: completeness, correctness, and F-score (Wiedemann et al. 1998, Sokolova et al. 2006). These parameters are determined by comparing the results obtained with the reference data. However, this is no minor task as it is complicated to generate the reference data.

In this paper, the alpha-shape algorithm in 2D space (Edelsbrunner et al. 1983) was explored to determine the reference data. This algorithm enables the identification of edge and non-edge points, and the extraction of building contours. However, the results derived from the alpha-shape algorithm depends on certain factors: the process used to segment the points sampled on different roofs, and the value of the radius adopted ( $\alpha$  parameter). To obtain the reference data the following procedure was undertaken: first, the points related to each building roof were manually segmented. The alpha-shape algorithm was then performed with the set of points related to each building roof. The value of  $\alpha = 0,60$  m, which was empirically obtained, was used in this calculation. Figure 4 shows some clippings made on LiDAR points. Figures 4a.1 and 4b.1 show the LiDAR points labeled in edge and non-edge classes after a classification process, while Figures 4a.2 and 4b.2 show the reference data generated using the alpha-shape algorithm.



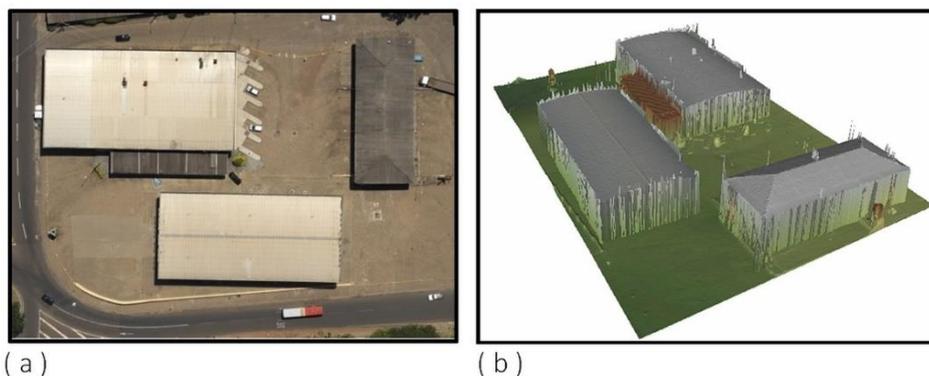
**Figure 4:** LiDAR points classified in the edge and non-edge classes (Figures a.1 and b.1). Reference data obtained using the alpha-shape algorithm (Figures a.2 and b.2).

As can be observed in Figures 4(a.2) and 4(b.2), the edges obtained through the alpha-shape algorithm are thinner and consistent compared with the real edges of the buildings sampled from LiDAR data.

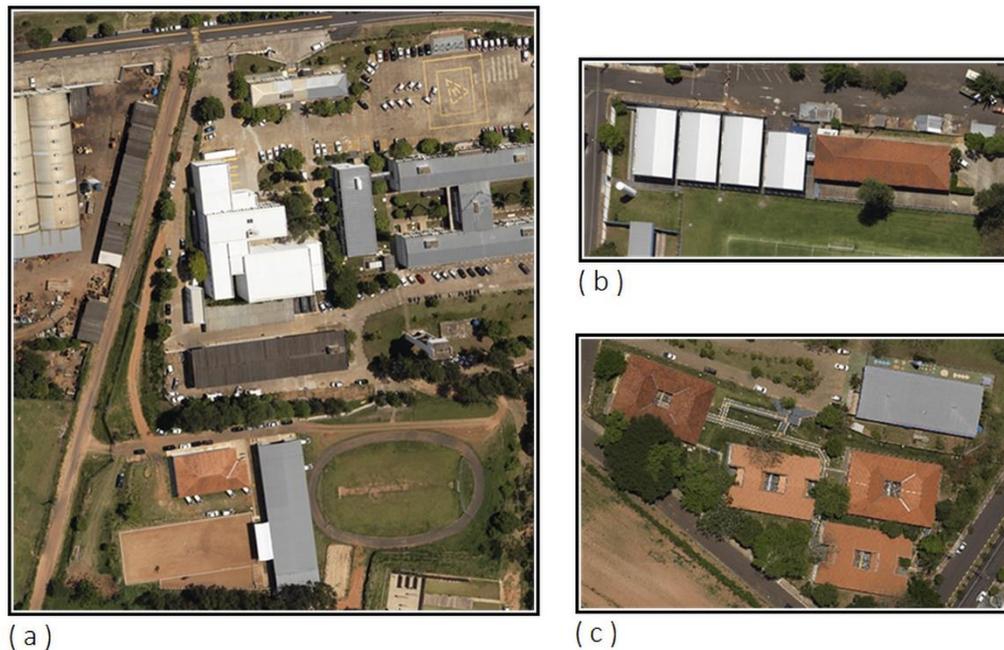
### 3. Results

This section presents the results obtained using the proposed classification method. The input data is a set of irregularly distributed LiDAR points saved in “las” format. For each point in this file is stored the Universe Transverse Mercator coordinate (E, N), an ellipsoid height (h), a pulse LASER return intensity (I), and others information. For details on “las” format specification the following reference is suggested: ASPRS (2013). The LiDAR data set used was acquired over the city of Presidente Prudente/Brazil by the company *Sensormap Geotecnologia*, which belongs to the *Engemap Group*. The average density of the cloud point is around 5 points/m<sup>2</sup> and the airborne LASER scanning system used to perform the recording was a RIEGL LMS-Q680i ([www.riegl.com/com](http://www.riegl.com/com)).

To verify the performance of the proposed method, four different areas were selected. Area 1 (Figure 5a) has large buildings, without the presence of trees in the proximity. The buildings are composed of curved roofs and planes, as can be seen in the digital surface model (Figure 5b). Area 2 (Figure 6a) has a greater quantity of buildings distributed across a larger area, where some have a complex form. Area 3 (Figure 6b) has closely juxtaposed buildings of similar height. Finally, Area 4 (Figure 6c) has complex buildings surrounded by vegetation.



**Figure 5:** Aerial image showing Area 1 (a), and digital surface model (b).

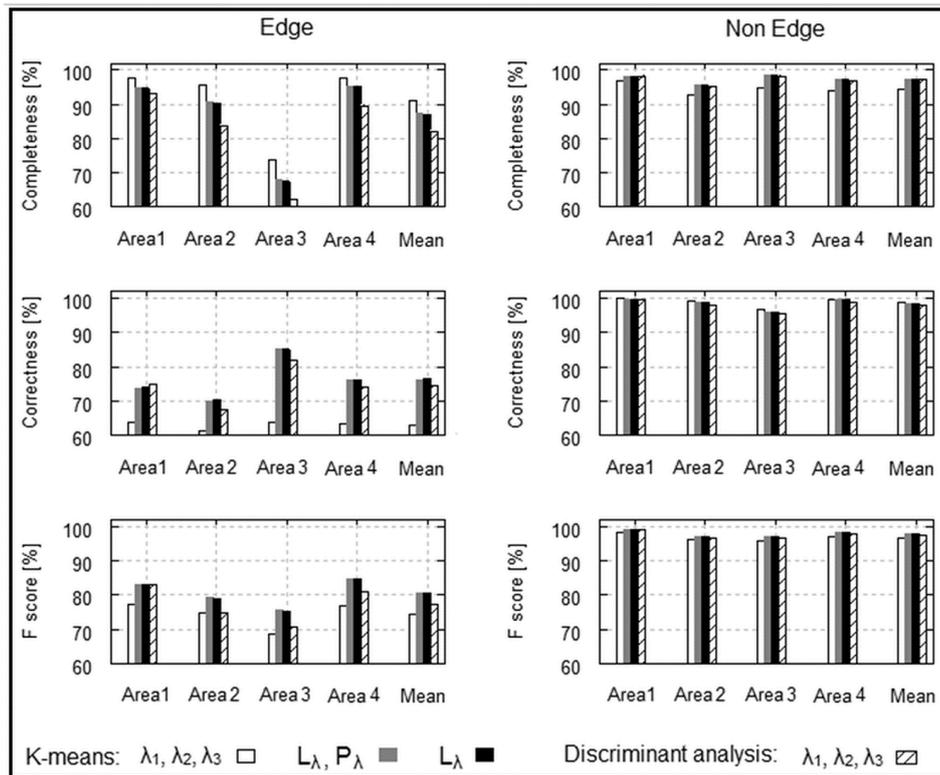


**Figure 6:** Aerial images related to Areas 2 (a), 3 (b) and 4 (c).

The filtering of the points sampled on roofs was performed using tools from *LAStools* software. The tools used were *lasground*, *lasheight* and *lasclassify*; and for each tool, the default parameters were considered. To perform the next steps, as shown in Figure 1, a program in C programming language was developed through *Code::Blocks*. The kd-tree storage structure and search functions, both implemented in the FLANN library (Muja and Lowe 2016), were used to perform neighborhood search operations. The visualization of results was performed using the FugroViewer and Gnuplot software.

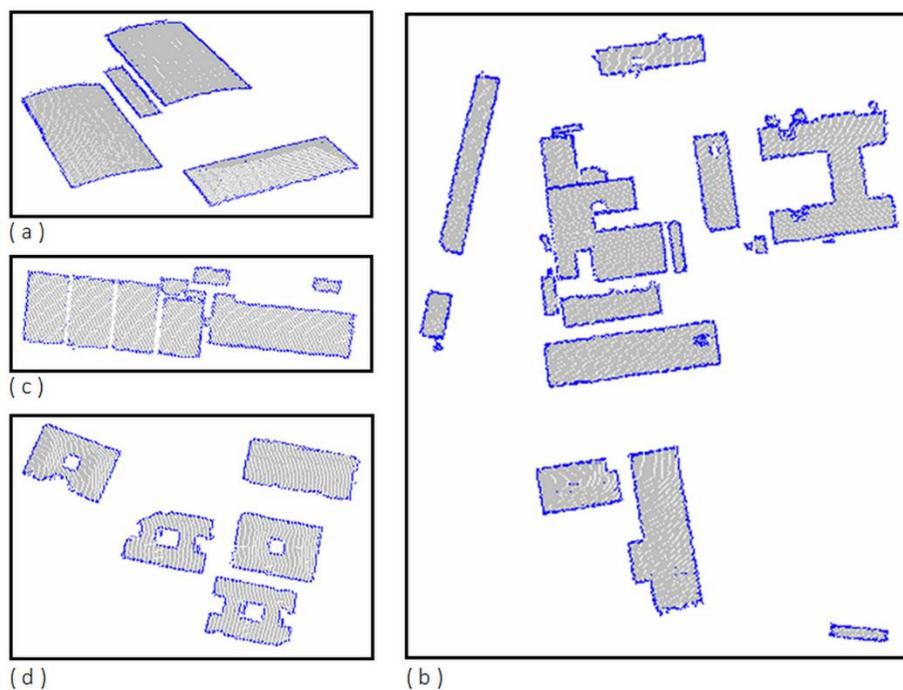
In order to perform the quantitative analysis, quality parameters were estimated for four different cases. In the first three cases, the LiDAR points were classified by the k-means method, considering three different groups of measures, as described in Subsection 2.3 and in Table 2. In the fourth case, the classification step was performed through discriminant analysis, similar to that used by Santos and Galo (2016), using theoretical eigenvalues related to eight different classes, as presented in Gross and Thoennessen (2006). For purposes of comparison, the eight classes were subdivided into two groups: edge and non-edge points. Points classified in the classes line, half plane, quarter plane, and two planes were marked as edge whereas points classified in the other four classes were marked as non-edge.

In Figure 7 quality parameters (completeness, correctness, and F-score) are shown related to the edge and non-edge classes. These parameters were calculated for four test areas (Figures 5 and 6). In addition, the mean parameters for the four areas were also calculated.



**Figure 7:** Quality parameters related to edge class (left) and non-edge class (right).

Figure 8 presents the results obtained from the proposed method using the linearity measure as a feature. They were generated considering as input data the LiDAR point clouds related to Areas 1, 2, 3 and 4. As it can be observed, the method identifies the points sampled on building roofs and classifies them into two classes: edge (blue) and non-edge (gray).



**Figure 8:** Classification results obtained by the proposed method using linearity measure as feature for Areas 1 (a), 2 (b), 3 (c) and 4 (d).

## 4. Discussion of Results

It is possible to see in Figure 7 that the values of quality parameters related to the non-edge class were higher than 92% for all areas and classification procedures, showing a mean value around 97%. The values of quality parameters for the edge class ranged from 61% to 98%, showing a mean around 79%. It can therefore be concluded that the edge class is most affected by classification errors. This result was expected due to the discrete nature of LiDAR points, which does not guarantee the correct sampling of building edges since the success of the sampling process on this kind of points is strongly dependent on the scanning rate and point density.

By comparing the results obtained by the k-means method for the edge class (Figure 7) and considering different groups of measures, it can be seen that the use of eigenvalues as measures generated better results in terms of completeness. However, in terms of correctness, the sets formed by linearity and planarity, and by linearity alone, showed better results. The selection of the set used as measure is directly related to the application. For example, if the goal is the extraction of features of a particular class of object, it is appropriate to use the set of measures that yield higher correctness for this class.

Comparing the values of the F-score related to the k-means method (Figure 7), it can be seen that the use of eigenvalues as variables did not produce the best results. This characteristic was observed in all test areas, being more evident for the class edge. The best results were obtained considering the set of measures composed by linearity and planarity, and the set composed only by linearity. As can be observed in Figure 7, the generated results using these two set of measures were similar. The F-score is a general parameter, which is calculated considering both completeness and correctness. It can therefore be concluded that the classification through the k-means method presents the best results when the linearity and planarity measures are considered together, or when linearity is considered alone.

Performing a comparative analysis between the results generated by the k-means method (Figure 7), considering as measure the linearity and planarity, or only linearity, and also the results generated by the discriminant analysis method, it is possible to observe that the k-means method produced the best results in almost all cases. In particular, the results are similar for Area 1, in which the buildings are separated and have no trees around them.

When comparing the results from the k-means method using eigenvalues as measures with the results derived from discriminant analysis, it is possible to see that an alternation occurred between the processes. Therefore, it can be affirmed that the k-means method presents satisfactory and consistent results in the context of the classification of LiDAR points, mainly when linearity and planarity, or only linearity were considered as measures.

Figure 8 shows graphically the results obtained through the k-means method and using the linearity measure as a feature, considering four test areas with different characteristics. From visual analysis, it can be observed that most of the LiDAR points were correctly classified. However, it is possible to identify some situations in which the classification generated incorrect results.

Figure 8a shows the results of Area 1. This is an area with little complexity, composed of four isolated buildings, two formed with curved roofs. Performing the visual analysis, it is noted that the points were correctly classified without major inconsistencies. This can also be observed when the quality parameters (Figure 7) are analyzed, thus underlining the efficiency of the method in the classification of areas composed of isolated buildings.

Figure 8b shows the results related to Area 2. This area is composed of a wide diversity of objects: trees, cars, buildings, etc, as can be seen in Figure 6a. With regard to buildings, this area is formed by isolated buildings, located close to each other, of different dimensions and heights, and with different types of roofs (curved, inclined, and formed by various planes). By carrying out the visual analysis, it can be noted that, in general, the points were correctly classified. Analyzing the correctness of the edge class (Figure 7), it is noted that the lowest value is related to Area 2. Despite this, the value of correctness was around 70%, indicating the efficiency of the method also in areas with a higher level of complexity.

Figure 8c shows the results of Area 3, which is composed of rectangular buildings, where four of them are located quite close to each other. Due to this factor, several edge points were incorrectly labeled in the non-edge class. This classification error is reflected in the value of the completeness parameter, as can be observed in the edge class related to Area 3 in Figure 7, which presented the lowest value. Thus, the proposed method has a limitation with respect to classification of edges points located very close to other buildings.

Area 4 is composed of isolated buildings, some of which are partially covered by trees, as can be seen in Figure 6c. In regions where occlusions occurred, due to the presence of dense vegetation around the buildings, it is possible to note that the roofs are incomplete, as expected. Despite occlusion of these areas, in general the points were correctly classified (Figure 8d). This can also be observed by analyzing the quality parameters in Figure 7. The method can therefore be seen to be efficient in the classification of points sampled on roofs partially covered by vegetation.

In summary, quantitative and qualitative analyses indicate that the proposed method can be used to classify the LiDAR points sampled on roofs in the edge and non-edge classes, obtaining satisfactory results.

## 5. Conclusion

A procedure for the classification of LiDAR points sampled on roofs in two classes, edge and non-edge points, has been proposed in this paper. This method combines the concept of principal component analysis and the k-means clustering. The main advantage of the proposed approach is the automatic class definition, without the need to establish thresholds.

The results, obtained for real data, show that even in the face of the complexity related to buildings in urban areas, which are composed of different types of roofs (curved, inclined and formed by various planes), and, in some cases, roofs partially covered by vegetation, the results obtained by the proposed method were satisfactory. The efficiency of the classification method was verified through visual and quantitative analysis. In summary, for the non-edge class, both the completeness and correctness were above 92%, and between 61% to 98% for the edge class. In addition, it was possible to verify that the set of measures influences the classification result.

The results generated can be applied in the selection and extraction of various geometric primitives, such as planes, line segments, curved segments, corner points, among others. These primitives are essential in the geometric modeling of buildings.

In future research, the use of this classification procedure is suggested to extract building contours, followed by building contours regularization. Besides that, it is also necessary to improve the building contours extraction from the situations in which edge points from neighboring

building are very close. Finally, it is suggested that the proposed method can be applied to the point cloud obtained by TLS (Terrestrial LASER Scanning) systems.

## ACKNOWLEDGEMENT

The authors would like to thank *Sensormap Geotecnologia* for providing the LiDAR data used in the experiments, and CNPq (grant nº 304189/2016-2) and FAPESP (grant nº 2016/12167-5) for their financial support.

## REFERENCES

- ASPRS *LAS Specification – Version 1.4 – R13*, July, 2013. URL: [www.asprs.com](http://www.asprs.com). Acess: June/2017.
- Bretar, F. Feature extraction from LiDAR data in urban areas. SHAN, J., TOTH, C. K. *Topographic laser ranging and scanning: principles and processing*. CRC Press, Taylor & Francis Group, 2009, 590p.
- Carrilho, A. C. *Aplicação de técnicas de processamento e análise de imagens para detecção de edificações e vegetação a partir de dados LiDAR*. Dissertação (Mestrado em Ciências Cartográficas) PPGCC - Programa de Pós-Graduação em Ciências Cartográficas, Faculdade de Ciências e Tecnologia, Universidade Estadual Paulista. Presidente Prudente/SP, 2016
- Chehata, N.; David, N.; Bretar, F. LiDAR data classification using hierarchical k-means clustering. *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, Beijing, p. 325-330, 2008.
- Cruz, Z. Q.; Silveira, J. C. Aplicação do sistema laser scanner aerotransportado para identificação de linhas de transmissão e atualização de plantas de perfil topográfico. *Simpósio Brasileiro de Sensoriamento Remoto (SBSR)*, Curitiba, 2011.
- Demantke, J.; Mallet, C.; David, N., Vallet, B. Dimensionality based scale selection in 3D lidar point clouds. *ISPRS Workshop on Laser Scanning 2011*, Calgary - Canadá, 2011.
- El-Sheimy, N.; Valeo, C.; Habib, A. *Digital terrain modeling: acquisition, manipulation and applications*. London: Artech House, 2005, 257 p.
- Edelsbrunner, H.; Kirkpatrick, D. G.; Seidel, S. On the shape of set of points in the plane. *IEEE Transactions on Information Theory*, v. IT-29, n. 4, p. 551-559, 1983.
- Galvanin, E. A. S.; Dal Poz, A. P. Extraction of building roof contours from LiDAR data using a Markov-Random-Field-Based approach. *IEEE Transactions on Geoscience and Remote Sensing*, v. 50, n. 3, 2012.
- Gross, H.; Thoennessen, U. Extraction of lines from laser point clouds. *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, Istanbul – Turkey, v. 36, p. 86-91, part 3, 2006.
- Habib, A.; Ghanma, M.; Morgan, M.; Al-Ruzouq, R. Photogrammetric and LiDAR data registration using linear features. *Photogrammetric Engineering and Remote Sensing*, p. 699-707, 2005.

- He, Y.; Zhang, C.; Awrangjeb, M.; Fraser, C. S. Automated reconstruction of walls from airborne LiDAR data for complete 3D building modelling. *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, XXII ISPRS Congress, Melbourne – Australia, 2012.
- Johnson, R. A.; Wichern, D. W. *Applied multivariate statistical analysis*. Upper Saddle River, NJ: Pearson Prentice Hall, 2007. 773p.
- Jutzi, B.; Gross, H. Nearest neighbor classification on LASER point clouds to gain object structures from buildings. *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*. 2009.
- Lari, Z.; El-Sheimy, N.; Habib, A. A new approach for realistic 3D reconstruction of planar surfaces from Laser scanning data and imagery collected onboard modern low-cost aerial mapping systems. *Remote Sensing*, 9 (3), 212, 2017.
- Machado, A. M. L.; Mitishita, E. A. Detecção automática de contornos de edificações utilizando imagem gerada por câmara digital de pequeno formato e dados LiDAR. *Boletim de Ciências Geodésicas*, v. 12, n. 2, p. 215-233, 2006.
- Muja, M.; Lowe, D. *FLANN - Fast Library for Approximate Nearest Neighbors*. Available in: < <http://www.cs.ubc.ca/research/flann/>>. Access in: June 1, 2016.
- Pauly, M.; Keiser, R.; Gross, M. Multi-scale feature extraction on point-sampled surfaces. *Computer Graphics Forum*, p. 81–89, 2003.
- Press, W. H.; Teukolsky, S. A.; Vetterling, W. T.; Flannery, B. P. *Numerical recipes in C: the art of scientific computing*. 2nd ed. Cambridge: Cambridge University Press, 1992. p. 994.
- Sampath, A.; Shan, J. Building roof segmentation and reconstruction from LiDAR point clouds using clustering techniques. *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, Beijing, p. 279-284, 2008.
- Sampath, A.; Shan, J. Segmentation and reconstruction of polyhedral building roofs from aerial Lidar point clouds. *IEEE Transactionson Geoscience and Remote Sensing*, v. 48, n. 3, 2010.
- Santos, R. C.; Galo, M. Classificação de pontos 3D utilizando o conceito de análise de componentes principais. *Simpósio Brasileiro de Ciências Geodésicas e Tecnologias da Geoinformação*, Recife/PE, 2014.
- Santos, R. C. *Extração de feições retas e cálculo de entidades pontuais a partir de dados LASER para o ajustamento relativo de faixas*. Dissertação (Mestrado em Ciências Cartográficas) PPGCC - Programa de Pós-Graduação em Ciências Cartográficas, Faculdade de Ciências e Tecnologia, Universidade Estadual Paulista. Presidente Prudente/SP, 2015.
- Santos, R. C.; Galo, M. Classificação de nuvem de pontos LASER utilizando o conceito de análise de componentes principais e o fator de não ambiguidade. *Boletim de Ciências Geodésicas*, v. 22, n. 2, p. 196-216, 2016.
- Shan, J.; Sampath, A. Building extraction from LiDAR point clouds based on clustering techniques. In: SHAN, J.; TOTH, C. K. (Editors). *Topographic Laser Ranging and Scanning: principles and processing*. Chapter 15. Boca Raton: CRC Press, Taylor & Francis Group, 2009. p. 421-444.
- Sokolova, M.; Japkowicz, N.; Szpakowicz, S. Beyond accuracy F-score and ROC: a family of discriminant measures for performance evaluation. *AI 2006: Advances in Artificial Intelligence*, v. 4304. Springer, Berlin, Heidelberg, p. 1015–1021, 2006.

- Tong, W. S.; Tang, C.K.; Mordohai, P. First order augmentation to tensor voting for boundary inference and multiscale analysis in 3D. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, v. 26, n. 5, p. 594-611, 2004.
- Vosselman, G. Slope based of LASER altimetry data. *International Archives of Photogrammetry, Remote Sensing*, v. 33 (B3), pp 935- 942, Amsterdam, 2000.
- Xiao, W. *Detecting changes in trees using multi-temporal airborne LIDAR point clouds*. Thesis submitted to the Faculty of Geo-Information Science and Earth Observation of the University of Twente in partial fulfilment of the requirements for degree of Master. Netherlands, 2012.
- Wang, G.; Zhang, Y.; Li, J.; Song, P. 3D road information extraction from LiDAR data fused with aerial-images. *Spatial IEEE Transactions on Geoscience and Remote Sensing*, p. 362-366, 2011.
- Wehr, A.; Lohr, U. Airborne laserscanning - an introduction and overview. *ISPRS Journal of Photogrammetry and Remote Sensing*, p. 68-82, 1999.
- West, K. F.; Webb, B. N.; Lersch, J. R.; Pothier, S.; Triscari, J. M.; Iverson, A. E. Context-driven automated target detection in 3-D data. *Proceedings of SPIE*, v. 5426, p. 133-143, 2004.
- Wiedemann, C.; Heipke, C.; Mayer, H.; Jamet, O. Empirical evaluation of automatically extracted road axes. *Empirical Evaluation Methods in Computer Vision*, Ed. Bowyer. IEEE Computer Society Press, p. 172-187. 1998.
- Yang, B.; Dong, Z. A shape-based segmentation method for mobile laser scanning point clouds. *ISPRS Journal of Photogrammetry and Remote Sensing*, p. 19-30, 2013.
- Zhang, J.; Lin, X. Filtering airborne LiDAR data by embedding smoothness-constrained segmentation in progressive TIN densification. *ISPRS Journal of Photogrammetry and Remote Sensing*, p. 44-59, 2013.
- Zhang, K. Q.; Chen, S. C.; Whitman, D.; Shyu, M. L.; Yan, J.; Zhang, C. A progressive morphological filter for removing nonground measurements from airborne LIDAR data. *IEEE Transactions on Geoscience and Remote Sensing*, 2003.