Manuscript

# SPATIAL CLASSIFICATION OF BUILDINGS FOR URBAN ENERGY PLANNING

**Federica Zagarella**[1,]

[1] Euro-Mediterranean Institute of Science and Technology, IEMEST, Palermo, Italy[2]

**CORRESPONDENCE:** Federica Zagarella
e-mail: federicazagarella@iemest.eu
Phone number: +393475828667

## Abstract

Estimating single buildings' energy demand of a whole city is something costly in terms of time and computational load. Hence, defining accurate but smart methods to assess the aggregated buildings energy demand is relevant.

In this study, a statistical model was implemented to classify buildings. Particularly, 150 buildings in Milan were manually labelled by use category and vintage and a conventional neural network was implemented, based on features related to real buildings characteristics and targets from statistical data. The ANNs' accuracy was validated by means of the Confusion Matrix and Receiver Operating Characteristic Curves. Finally, as an application, the ANN was retrained with a new input, including over 50,000 buildings in Milan, in order to assess the urban building stock which was also mapped in GIS environment for visualization purpose.

This study provides a useful preliminary assessment tool for urban planning, able to fill the gap of missing detailed data.

## Keywords

Urban planning; Buildings classification; Artificial Neural Networks; ROC and confusion analysis

## Introduction

Built environment is responsible for the 40% of energy consumptions [1] and consequently, plays a fundamental role in worldwide policies targeted on the decrease of energy consumption and related greenhouse gas (GHG) emissions.

Determining the energy use of a particular built environment is essential for accomplishing the mentioned goals at the local level. However, this issue is still challenging due to an heterogenous availability of data among Countries. Thus,

several studies regard the definition of methodologies for building typologies and related energy consumption determination dealing with lacking input data [2]. In this context, adopted approaches are widely diverse, varying from bottom-up to top-down based methods, from deterministic methods for accurate buildings modelling [3],[4],[5] to statistical methods surveying the possible correlations among buildings characteristics. In the field of black-box methods, the Artificial Neural Networks (ANNs) are increasingly used for energy use prediction [6],[7]. For instance, in [8] the authors used ANN for validating predicted energy performance certificates of residential buildings in Northern Italy. The model is robust however its replicability is affected by the energy certificate databases usage, which in Italy are available only in some regions and whose data refer only to buildings that underwent energy labelling process. In [9], the authors used ANN to forecast the primary energy consumption for space heating and cooling and the occupants' thermal comfort in Southern Italy, since buildings simulated characteristics (geometry, envelope, operation and heating, ventilation, and HVAC systems) and results.

However, going backward, energy consumptions historical or measured data are not always available and it's necessary

estimating them based on available data. Several research programs[1] and studies [3],[4],[5] have dealt with the issue of estimating the buildings energy use based on buildings typologies, defined by archetyping common characteristics (use category, age of construction, shape, etc.). Considering that some urban contexts are affected by lack of data, although the topic is well-known, even gathering data on each building's typology could be still challenging. Therefore, investigating a method for defining building types in a reliable way is an important task.

In Italy, this issue is particularly evident considering that data on buildings are provided in diverse consistency among the national territory or, if provided uniformly, are often available at larger scale. For dealing with this gap, this study proposes a methodology for spatial classification of buildings characteristics, which adopts an ANN for classifying buildings' typology in terms of period of construction and use category AND Geographic Information System (GIS) for mapping them.

## Methods

The research has been carried out at the Department of Nanotechnologies and

---

[1] The Tabula program, later converted in Episcope program (http://episcope.eu/index.php?id=97), is a relevant research program on buildings classification financed by European Union covering all countries building stocks.

Innovative Materials at the IEMEST research institute, in Palermo, Italy. As the NIMA dept. activities focus on buildings materials characterization, the accomplished research aimed at developing an opensource tool for large scale building characteristics mapping. In detail, the developed methodology consists of preliminary data gathering from available data sources and subsequent mapping by means of GIS. Then, the approach for choosing the samples, in terms of number and criteria was defined and some buildings were manually labelled. Then, a set of key variables was defined as features of the ANN, as derived by processing original data: three different and progressively extended sets of features were used and compared. Then, the most typical buildings' classes were defined as targets of the ANN and three different label methods were compared. After this pre-processing phase, the ANN model was implemented in MATLAB; 18 combinations were tested to find the most performing. Finally, as an application, the model was used with a larger dataset including all buildings in Milan city, Italy.

All outlined phases are following described with more detail.

<u>Data gathering and mapping.</u>

In Italy, one of the main issues for accomplishing effective energy policies, is the availability and quality of data. Two data sources, publicly available, were used for this study purpose.

The buildings base map comes from the Regional Topographic Data Base (DbTR)[2] from Lombardy Region, as one of the most complete in Italy. Among the Themes the DbTR is made, the *Built Environment* one was considered. It is made of four Classes: *Volume Unit*, *Roof Element*, *Architectural Detail*, *Building*, which in turn includes the buildings *Footprint* and the *MaxFloor*. The *Footprint* is a polygonal vector layer providing information (*attributes*) on building state, perimeter, and area. There are over 80 thousand samples in the original dataset that decrease to 52423 buildings, by removing minor, damaged and under construction ones. The *Volume Unit* is defined as each elementary unit of a building, having a uniform roof surface and a constant height. Both were used to model the buildings geometry, calculating additionally the volumes and the centroids.

The second buildings data source is the dataset form the National Institution of Statistics (ISTAT)[3]. ISTAT used to accomplish the "General Census of Population and Houses" throughout the Italian territory every 10 years until 2011.

ISTAT disseminates geographic data in shapefile format trough the set of so-called *Territory Bases*, with progressively increasing dimension. Among the available *Territory Bases*, the so-called *Census section base* was used, as it is the lowest unit for census relief (it usually corresponds to a block or a small group of buildings), which is provided as a polygonal vector layer, including information for the identification, perimeter, and area. Additionally, for characterizing the built environment, it was necessary acquiring the information contained in the so-called *Census Sections Variables* spreadsheet, provided for each census section of each Italian Region, and including aggregated data on population, houses, and buildings.

Definition of samples.

For the city of Milan, Italy, the ISTAT database contains over 6,000 census sections, while the DbTR base map is made of over 50,000 buildings. To develop the defined method, 15 census sections and 150 buildings were selected throughout the city (Figure 1). The selection of census sections was accomplished with the aim to consider same shares of buildings in central/semi-central/suburban areas, residential and tertiary, and by vintage. Then, selected buildings have been visually checked one by

one by means of Google Maps Street View[4], to assign more accurately the use category and vintage (as shown in the example in Figure 2). This elaboration was accomplished through comparison between the ISTAT and the DbTR databases to assure coherence.
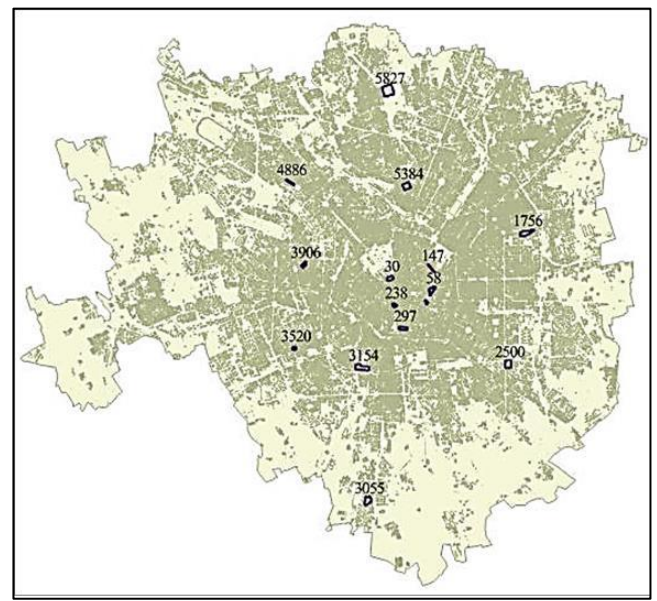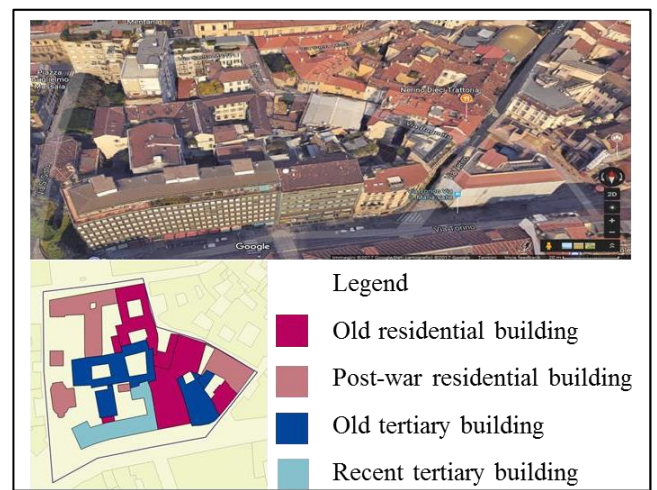


*Figure 1. GIS map of Milan with selected samples.*



*Figure 2. Example of samples labelling for one selected census section.*

[4] Google Maps Street View
https://www.google.com/maps/@42.2222411,14.3916736,3a,75y,167.25h,89.83t/data=!3m6!1e1!3m4!1s1gyyvKCzf-_aHCEtQatI-Q!2e0!7i13312!8i6656?hl=it-IT

<u>Definition of key variables.</u>

The subsequent step regarded the calculations of some key variables from the data included in the two used data sources, accomplished in QGIS environment[5], i.e.:

- **Buildings Perimeter to Floor Area Ratio (Perim/Surf):** the ratio of the buildings perimeter over their floor area was calculated to appreciate the buildings shape and compactness;

- **Buildings Volume to Floor Area Ratio (Vol/Surf):** the ratio of the buildings volume over floor area was calculated to appreciate the buildings shape;

- **Buildings Relative Radius (Radius%):** the distance of a building from the city centre has been taken into account because in radial cities, like Milan, the urban texture gradually becomes more recent towards the edges. In QGIS the centroids and related coordinates were determined for each building. Among them, the urban main church, the so-called "Duomo", was used as reference of Milan city centre to that end. Thus, the difference of coordinates values between each building centroid and Duomo centroid was calculated. Then,

the relative radius was calculated, as the ratio between the distance of a building from the city centre and the maximum distance between a building and the city centre according to the following equation:

$$Radius = abs(514978.175 - \text{x\_coord}) + \newline + abs\,(5034534.681 - \text{y\_coord});$$

- **Share of Tertiary buildings (Tertiary_Blds%):** the use category of buildings is an essential data to appreciate different operation and energy consumption profiles. From ISTAT, the percentage of not residential buildings (field E4) on the total number of used buildings (field E2) was calculated for each census section;

- **Share of residential buildings with construction different from masonry (NotMasonry_Blds%):** ISTAT technicians empirically survey the construction material of residential buildings split by masonry, concrete and other material (steel, mixed, etc.). Since buildings in masonry represent the greatest share, the percentage of buildings built in concrete or other material (fields E6 and E7) on the total number of used residential buildings (field E3) for each census section was used.

---

[5] QGIS https://qgis.org/it/site/

- **Share of residential buildings built before 1945 (Old_Blds%):** ISTAT technicians empirically survey the vintage of residential buildings by assigning intervals. As, different thermal transmittance and energy behavior could be associated based on building vintage, the first two periods (<1919 and 1945-60) were combined and the percentage of old buildings (fields E8 and E9) on the total number of used residential buildings (field E3) were determined for each census section;

- **Share of residential buildings built before 1990 (Postwar_Blds%):** similarly, the periods relative to construction post World Wars (i.e., from 1961-70 to 1981-90) were combined and the percentage of post-war buildings (fields from E10 to E13) on the total number of used residential buildings (field E3) were determined for each census section;

- **Share of residential buildings with more than 8 flats (Large_Blds%):** ISTAT technicians survey the number of flats in residential buildings dividing them by intervals. For simplicity, the ones relative to biggest buildings (9-15 flats interval and more than 16 flats interval) were combined and the percentage of larger buildings (fields E25 and E26) on the total number of used residential buildings (field E3) was determined for each census section.

Implementation of an Artificial Neural Network

As described in [10],[11],[12], the ANNs have been developed as generalizations of mathematical models of biological nervous systems. The basic unit of an ANN is the neuron i.e., a processing element for computing a non-linear function of the input. Every neuron is connected to other neurons though an assigned weight that defines its impact on the output. The network consists of at least three layers of neurons: an input layer, a hidden layer, and an output layer. The output is obtained as a non-linear function of the total weighted input minus a bias term. Optimizing the number of hidden layers is a challenging task as well as connecting a proper number of nodes and setting appropriate weights. The structure of nodes, connections and weights determines the final behavior of the ANN.

For analyzing the ANN model performances, a Receiver Operating Characteristic (ROC) Curve and a Confusion Matrix were used [13]. The ROC curve is basically used for assessing the performance of classification methods. Given a set of values to be

predicted according to n classes, it involves the count of true versus false (positive and negative) associations. Results are given on a chart, where y axis represents the sensitivity of the test, while the x axis its specificity. A threshold between true and false associations is represented by the diagonal: roughly, the more the curve follows the left and the top sides, the more it is accurate. Moreover, as a measure of the accuracy, the calculated Area Under the Curve (AUC) should be closed to 1.

Regarding the implemented model, **input data (i.e., features)**, were defined as follows. The used input data correspond to a matrix of 150 rows (buildings from the selected section census sections) and n columns. From the 8 features described before, three different combinations of buildings features were tested to appreciate their different accuracy:

A. 3 geometry features: Perim_Surf, Vol_Surf, and Radius%;
B. 5 typological features: Perim_Surf. Vol_Surf, Radius%, Tertiary_Blds%, and Old_Blds%;
C. 8 combined features: Perim_Surf, Vol_Surf, Radius%, Tertiary_Blds%, Old_Blds%, NotMasonry_Blds%, Postwar_Blds%, and Large_Blds%.

Regarding the definition of the **target data**, a matrix of 150 rows * 6 classes was defined, corresponding to these building classes:

(1) Residential Old building;
(2) Residential Post-war building;
(3) Tertiary Old building;
(4) Tertiary Post-war building;
(5) Residential Recent Building;
(6) Tertiary Recent building.

Once the features and classes were defined, another essential task was defining the labelling procedure, since it can affect the quality of outcomes and the possibility of replicate this procedure in other contexts. Considering this, 3 approaches for labelling the 150 building samples were proposed:

α. "**visual individual assignment**": each building' class was visually checked by means of Google Street View and manually labelled;

β. "**analytical rough assignment**": based on ISTAT aggregated data, the determined census section's main properties were assigned to all enclosed buildings. The census section main use was assigned based on the highest number of buildings with residential or tertiary use; the census section main period of construction was assigned based on the highest number of residential buildings built in a certain

period (the main period of residential buildings was assigned also to the tertiary ones as a building block often dates back the same construction period);

γ. "**analytical crossed assignment**": based on ISTAT aggregated data, the determined census section main properties were assigned to all enclosed buildings. The census section main use was assigned as in approach "β" while the main period of construction was determined by correlating the data on the age with the residential and tertiary buildings size.

By means of Neural Pattern Recognition, in MATLAB software,[6] an ANN was implemented using pre-defined codes for assessing the mentioned 9 cases. An additional variation concerned the number of hidden neurons, set proportionally to features number: 3 and 6 for target α, 5 and 10 for target β, 8 and 16 for target γ. In all 18 tested cases (Figure 3), the training set was made of 104 samples (75%), while the validation and the testing sets by 23 samples (15%).
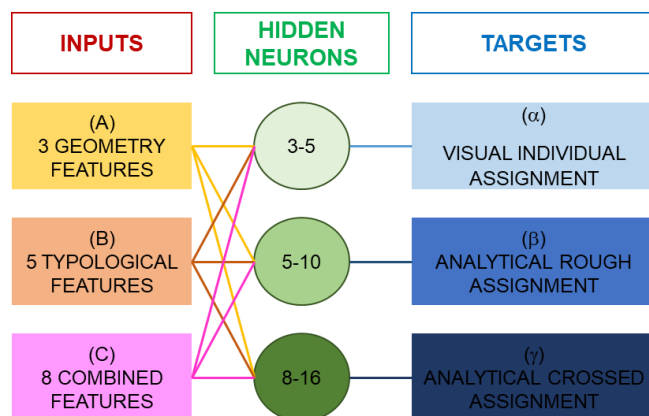


*Figure 3. Scheme of performed 18 ANN models.*

## Results

In the following sections, results of running the ANN 18 cases are presented with reference to the Error Histogram, the Confusion Matrixes, and the ROC Curves. For clarity, results are separately discussed by each set of input data. Then, the results of an application of the best performing ANN at urban level are presented and discussed, too.

Analysis with input A.

In tests regarding the association with target α (Figure 4a), the Error Histogram shows very random paths with irregular tails, even if the most of instances has small error. The Confusion Matrixes shows several null true associations. Consistently, the ROC Curves widely cross the diagonal.

In tests regarding the association with target β (Figure 4b), the Error Histogram has a quite irregular path; in the Confusion Matrixes, classes 2 and 3 have no true values; the ROC Curves still have not an acceptable

behavior. Generally, performances decrease increasing hidden neurons number.

In tests about the association with target $\gamma$ (Figure 4c), most of instance have a small error and small tails; in Confusion Matrixes, class 2 lacks true values; the ROC Curves are closed to left and top sides. Better performances than previous cases can be noted, although they cannot be considered sufficient yet.

Summarizing, results clearly show that this model, based on only three geometry features (input A), is not able to carry out good predictions of building typologies, even if an increase in the performances from target $\alpha$ to $\gamma$ is noted.

Analysis with input B.

In tests on target $\alpha$ (Figure 5a), Error Histogram has most of instances closed to zero, mostly in the case with low hidden neuron (error of 0.0023) and an irregular path. In Confusion Matrixes, the number of false values and null true ones is quite high. ROC Curves largely intersect diagonal, and class 1 performs worst.

In tests on target $\beta$ (Figure 5b), the Error Histogram shows an alignment of more instances to -0.013. Confusion Matrixes report overall performances of 100% in all cases. ROC Curves perfectly fit, having an

AUC value of 1 for each one; only in validation phase class 2 performs worst.

In tests on target $\gamma$ (Figure 5c), the Error Histogram has most of instances with a small error. The Confusion Matrixes report an overall performance greater than 97% with 100% in specific phases; usually, lower performances are reported for class 2 and 6. ROC Curves have averagely good paths. Although the ANN shows lower performances with more hidden neurons, these are good.

Summarizing, results regarding the model with 5 typological features (input B) show better performances compared to the model with input A. This happens probably because the randomness of targets decreased, being uniformly assigned to buildings belonging to the same census section. Also, by increasing number of features, better ANN performances occur.

Analysis with input C.

In tests with target $\alpha$ (Figure 6a), the Error Histograms (both with 8 and 16 hidden neurons) have again a very fluctuating path but with errors mostly lower than -0.05. The Confusion Matrixes show worsening results from the training to the test sets; they also report few good predictions in several classes. ROC Curves widely cross the diagonal with worst performance for class 2 and, unusually, for class 1.

In tests with target β (Figure 6b), in the Error Histograms most of instances report errors of more than 4 (8 hidden neurons) and 7 (16 hidden neurons). However, Confusion Matrixes have 100% of good overall predictions in all sets; with less hidden neurons performances among classes drop from training to test phases. ROC Curves perfectly fit (AUC equals 1), except for classes 2, in both cases with 8 and 16 hidden neurons, and 1, in testing with 8 hidden neurons.

In tests with target γ (Figure 6c), in Error Histogram with 8 hidden neurons most of instances have an error of less than -1, while close to 7 with 16 hidden neurons. Confusion Matrixes fit with overall values of 100% positive in all phases. ROC Curves perfectly fit (AUC equals 1), except for class 3 in validation, class 6 in testing (8 hidden neurons), and class 2 in testing (16 hidden neurons).

Summarizing, results of using 8 combined features (input C) show that the maximum number of input data associated to the "analytical crossed assignment" labelling approach, which correlates age and size of buildings, allows to better fit the model finding a pattern in random data.

Application on buildings classification at urban scale.

Considering the previous tests, the ANN with input C and target γ, due to best performances, was retrained until reaching the least number of hidden neurons (i.e., 3), although maintaining optimal performances and avoiding overfitting. In Figure 7, the Error Histogram shows most of instances with a small error; Confusion Matrixes have overall predictions equal to 100%; ROC Curves perfectly fit and an AUC of 1 is reported in overall phase for all classes.

Then, the targets were determined by retraining the ANN with new input data (52423 samples' rows * 8 features' columns) to estimate the buildings classification at urban level i.e., for the whole Milan city. As shown in Figure 8, the model tried to predict the typology of buildings putting some emphasis on the proximity among buildings (probably because the radius was used as feature and, more, because similar features occur within each census section) and this is reasonable. A decrease in buildings period of construction from city centre toward the edges is reported, although slightly evident. However, an unusual prevalence of recent buildings is returned, therefore for applications on a larger scale the model should be improved, maybe using more samples.
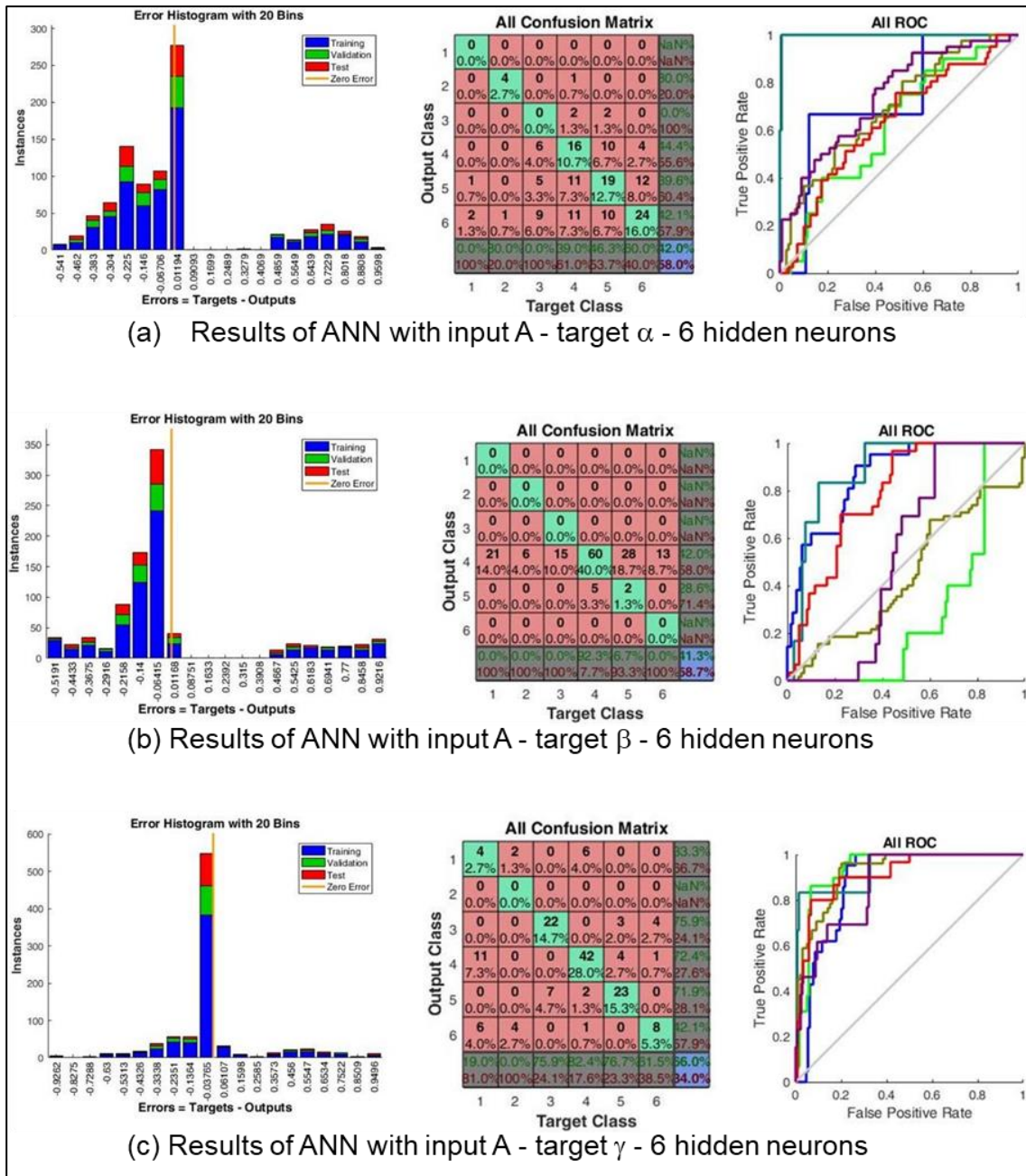
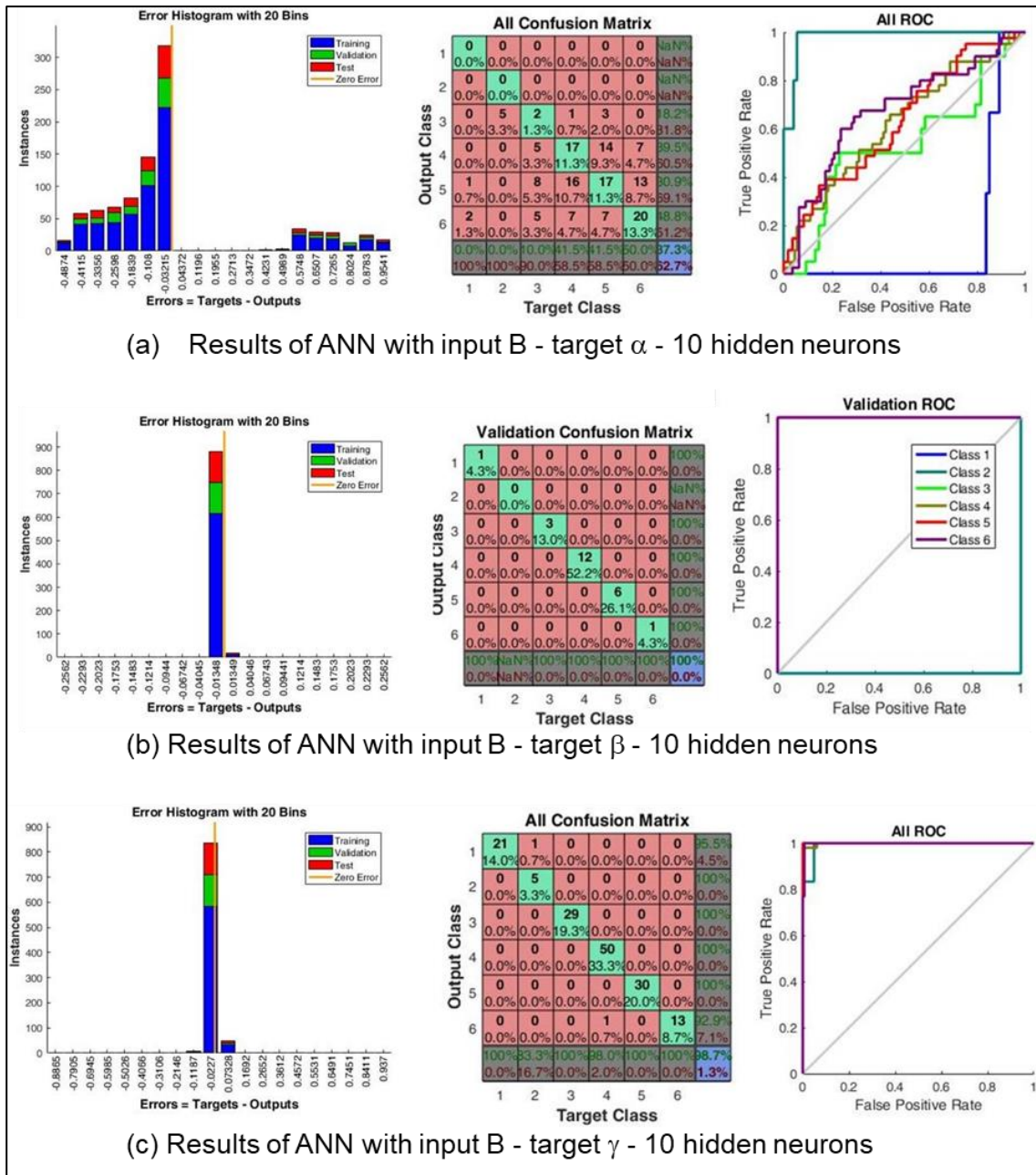Figure 4. Results of analysis with input A.

(a)    Results of ANN with input B - target $\alpha$ - 10 hidden neurons
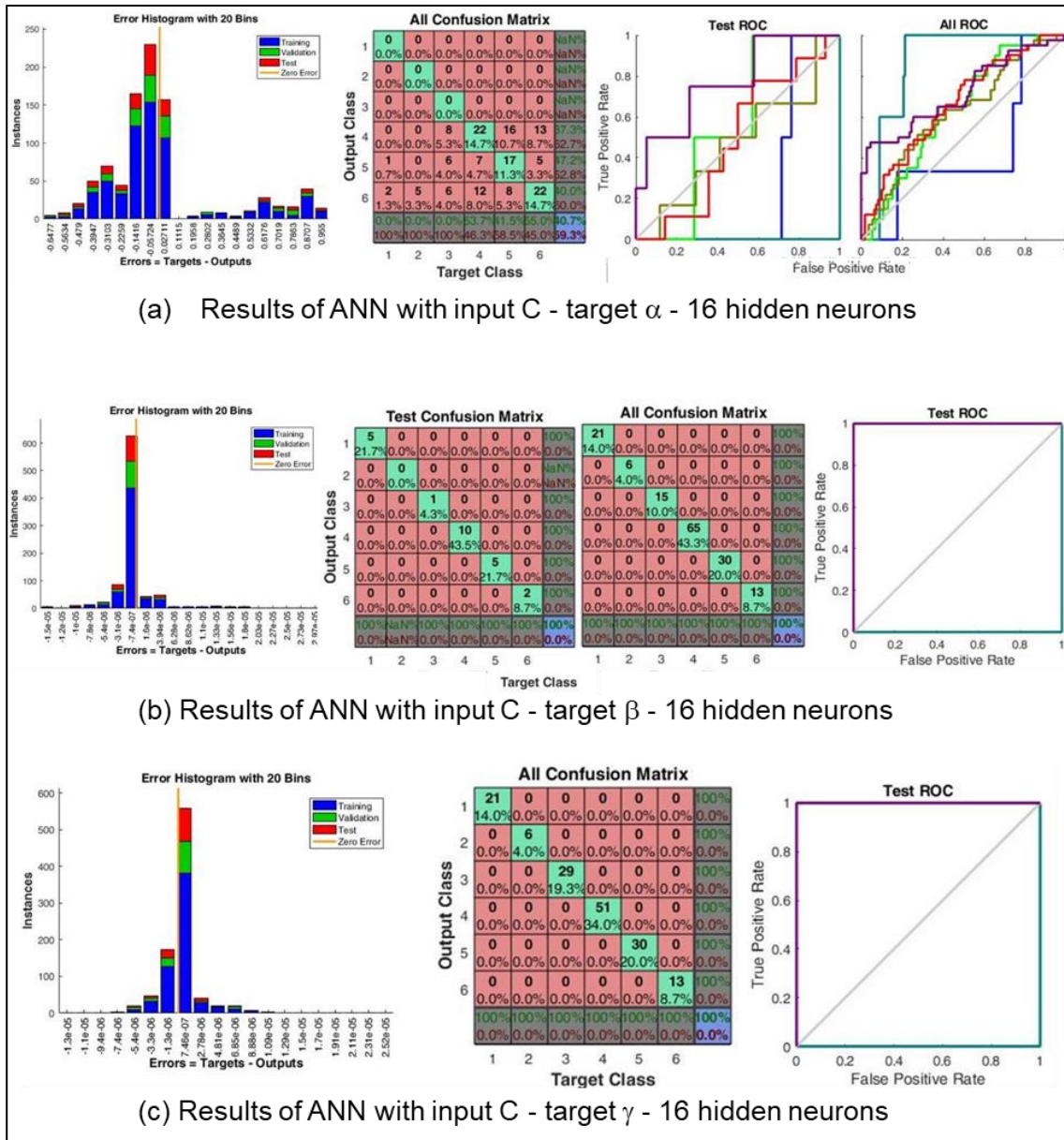
(b) Results of ANN with input B - target $\beta$ - 10 hidden neurons

(c) Results of ANN with input B - target $\gamma$ - 10 hidden neurons

*Figure 5. Results of analysis with input B.*

(a) Results of ANN with input C - target $\alpha$ - 16 hidden neurons

(b) Results of ANN with input C - target $\beta$ - 16 hidden neurons

(c) Results of ANN with input C - target $\gamma$ - 16 hidden neurons

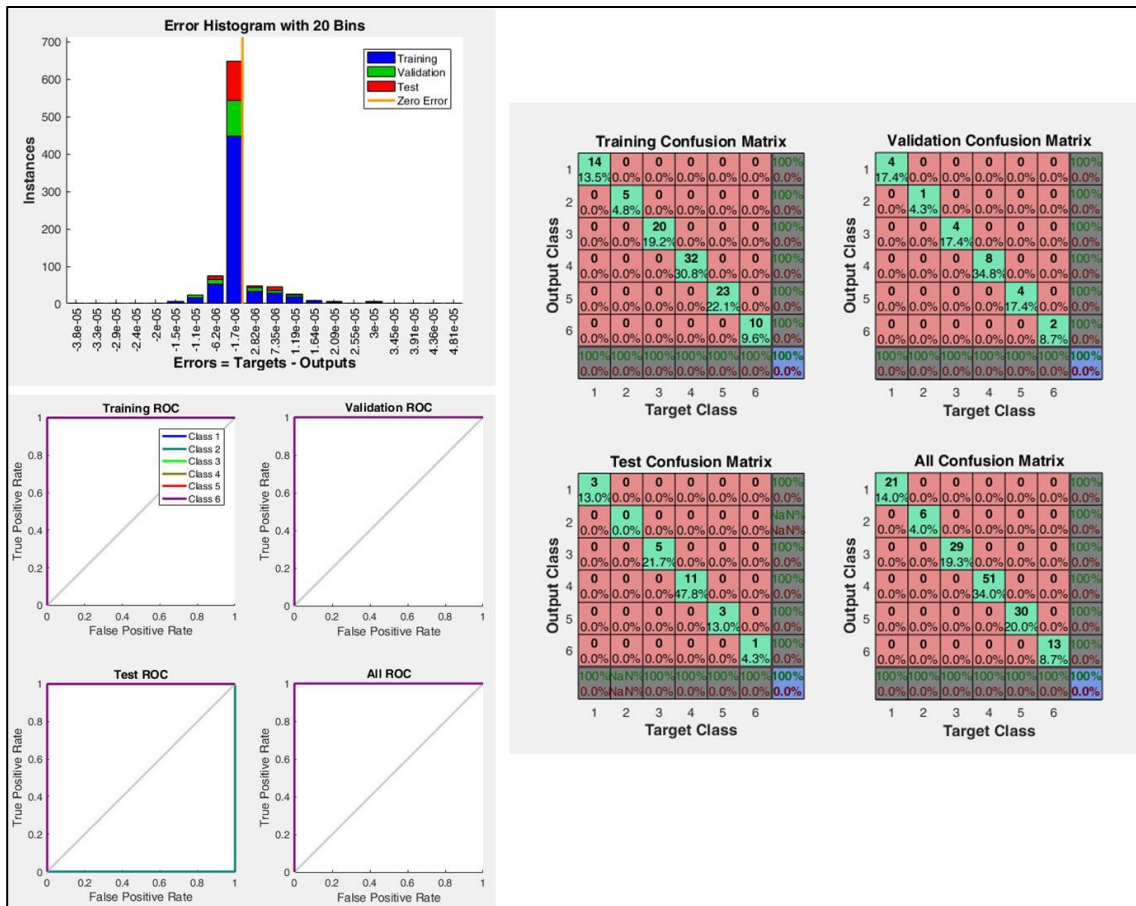*Figure 6. Results of analysis with input C.*

*Figure 7. Results of ANN with input C - target $\gamma$ - 3 hidden neurons.*



*Figure 8. Buildings classification application at urban level for Milan city, as result of best ANN retrain.*

## Conclusions

Starting from data on building geometry and typological characteristics, a sample of buildings was classified and their correlation with defined features was tested with an ANN. Several combinations, in terms of input, targets, number of hidden neurons, and labelling approaches, were tried. Results show that, by increasing the number of features, the model performs better. Also, using processed targets, based on statistical data, in place of manual labels, makes the model more stable. Therefore, last performed combination, with more features and targets assigned based on statistical crossed data, revealed optimal performances so was chosen for an urban level application. A balance is among the three test sets (training, validation, and test) even if training performance are usually slightly higher. Six classes were defined, based on three periods of construction and two uses. Class 2, corresponding to the post-war residential buildings, performs worse than the other, probably because it is the most consistent and widespread.

Generally, this study shows that it is possible reaching a good prediction of buildings spatial classification based on statistical available data, rather than individual manual labelling (a possible reason could be that the latter method is subjected to highest variability). Also, this approach is suitable because it allows to better clustering neighborhoods with an acceptable approximation, regardless single variations.

**This study relevance lays in the fact that provides an interesting tool for urban planning purpose, as it allows to estimate buildings characteristics based on a small sample and to visualize the distribution of buildings typologies at urban level. Hence, the joint action of ANN and GIS provides a powerful tool, which is at the same time accurate and intuitive. The model is also versatile as it can be integrated with other data, based on local needs and availability.** Several applications can derive from this study. For instance, it could be adopted to map buildings features in unknown areas. Based on the urban energy balance top-down data and bottom-up building features, the neighborhoods energy consumption could be estimated. Finally, as a potential further development, this methodology could represent a tool to identify the optimal districts for implementing renovation or energy strategies, aiding the public administrations in better targeting financial supporting schemes for envelope retrofit on effectively energivorous buildings, towards a more efficient and impacting decarbonization of the urban building stock.

# References

[1] European Union (EU) Commission and Parliament, Directive 2010/31/EU of the European Parliament and of the Council of 19 May 2010 on the energy performance of buildings (EPBD Recast).

[2] Keirstead J., Jennings M., Sivakumar A. A review of urban energy system models: Approaches, challenges and opportunities. Renewable and Sustainable Energy Reviews 16 (2012) 3847-3866.

[3] Caputo P., Costa G., Ferrari S. A supporting method for defining energy strategies in the building sector at urban scale. Energy Policy 55 (2013) 261–270.

[4] Corgnati S. P., Fabrizio E., Filippi M., Monetti V. Reference buildings for cost optimal analysis: Method of definition and application. Applied Energy 102 (2013) 983–993.

[5] Fonseca J. A., Schlueter A. Integrated model for characterization of spatiotemporal building energy consumption patterns in neighborhoods and city districts. Applied Energy 142 (2015) 247–265.

[6] Rafe Biswas M.A., Melvin D. Robinson, Fumo Ne. Prediction of residential building energy consumption: A neural network approach. Energy 117 (2016) 84-92.

[7] Ekici B. B., Teoman Aksoy U. Prediction of building energy consumption by using artificial neural networks. Advances in Engineering Software 40 (2009) 356–362

[8] Khayatian F., Sarto L., Dall'Ò G. Application of neural networks for evaluating energy performance certificates of residential buildings. Energy and Buildings 125 (2016) 45–54.

[9] Ascione F., Bianc N., De Stasio C., Mauro G. M., Vanoli G. P. Artificial neural networks to predict energy performance and retrofit scenarios for any member of a building category: A novel approach. Energy 118 (2017) 999-1017.

[10] Melchiorre C., Matteucci M., Azzoni A., Zanchi A. Artificial neural networks and cluster analysis in landslide susceptibility zonation. Geomorphology 94 (2008) 379–400.

[11] Kumar R., Aggarwal R.K., Sharma J.D. Energy analysis of a building using artificial neural network: A review. Energy and Buildings 65 (2013) 352–358.

[12] Bishop C. M. Neural Networks for Pattern Recognition. Clarendon Press, Oxford, 1995.

[13] Swets J. A. Measuring the Accuracy of Diagnostic Systems. Science, 240 (1988) 1285-1293.