# THE SIGNIFICANCE AND APPLICATION OF COMPACTNESS MEASURES IN MACHINE LEARNING

*Davrbek Xudayorovich Tursunmurotov*
*Teacher of the Department of Fundamentals of Informatics,*
*Tashkent University of Information Technologies*
*named after Muhammad al-Khwarizmi*

**Abstract:** This article explores the concept of compactness as one of the factors that contribute to the generalization ability of models in machine learning. Compactness is a characteristic that reflects the closeness, density, and organization of data, significantly impacting a model's performance on test data. The article provides both an intuitive and formal description of compactness measures, explains how they can be evaluated, and in which scenarios their use is justified. It also discusses the relationship between generalization ability and compactness using practical examples. The results of the study open new possibilities for improving model quality through compactness assessment.

**Keywords:** compactness, generalization ability, overfitting, underfitting.

The success of machine learning systems directly depends on their ability to generalize. This ability ensures high model accuracy not only on training data but also on previously unseen test examples. The level of generalization is one of the key criteria determining the practical utility of a model[1,2].

Among the factors influencing model quality are the quality of training data, model complexity, overfitting and underfitting, as well as the choice of hyperparameters. However, in addition to these aspects, there is another important but often overlooked concept—compactness.

Compactness reflects how data is structured within the model and the nature of distances between individual samples. When compactness is high, the model typically shows better generalization ability, as such a structure more accurately reflects the natural grouping or similarity of data.

This article aims to highlight the importance of compactness measures in machine learning, demonstrate methods for their evaluation, and analyze the conditions under which using this metric is most appropriate. The article examines the relationship between compactness and generalization ability, as well as practical applications of this concept in machine learning algorithms.

## The Concept of Compactness

In machine learning, compactness is an important characteristic that shows how well-organized examples are in a dataset, how closely they are grouped within a class, and the structure of their mutual positioning. Intuitively, if examples of the same class are located close to each other and far from other classes, such a structure is considered "compact." With high compactness, classification or clustering results are usually more accurate and reliable.

## General Definition of Compactness

From a mathematical point of view, compactness is often expressed as the ratio between the average intra-class distance (intra-class distance) of examples in one class and their distance from other classes. Based on these two criteria, various metrics have been developed to assess cluster quality and the overall degree of compactness[5].

## Geometric and Statistical Interpretation of Compactness

**Geometric approach**: This approach takes into account the shape of clusters, their size, and position in space. It relies on distance-based metrics (e.g., Euclidean distance, Mahalanobis distance).

**Statistical approach**: This method assesses variance, covariance, and the probability distribution of each cluster. In statistical models (e.g., Gaussian Mixture Models), compactness is analyzed using probability density functions[3,7].

**Compactness and Generalization Ability**

In machine learning, generalization refers to a model's ability to perform well on new, unseen data. In other words, the model must produce accurate results not only on training data but also on test data, with minimal error. This characteristic determines the model's ability to solve real-world problems.

**Relation to Overfitting and Underfitting**

In the case of **overfitting**, where a model memorizes the training data too precisely (including noise), its generalization ability decreases.

In the case of **underfitting**, the model fails to capture even the basic patterns in the data. In both cases, performance on test data suffers.

Compactness helps to find a balance between these two extremes. It preserves the structure and closeness within a class without excessively complicating the model.

**Improving Generalization Through Increased Compactness**

Compact structures are more resilient to variations in test data. This is because they rely not on noisy examples but on typical representatives of a class. If samples within a class are tightly clustered, the model is more likely to classify similar future examples correctly. Therefore, increasing compactness is directly related to improving generalization[8,9].

**Application of Compactness Measures in Machine Learning**

**Supervised Learning**

In supervised learning, compactness helps identify the internal structure of classes on which the model is trained. When class members are close to each other and form dense groups, model accuracy improves. This is evident in the following popular models:

*k*-**Nearest Neighbors (KNN)**: KNN heavily relies on intra-class compactness. If objects of a class are tightly grouped, their nearest neighbors are likely to belong to the same class, reducing classification errors. For optimal compactness, intra-class distances should be minimized, and inter-class distances maximized[10,11].

**Support Vector Machines (SVM)**: SVM seeks to place the decision boundary as far as possible from each class. If the classes are compact internally, the margin (width of the separating band) increases, enhancing generalization. Compact classes make SVM more robust and accurate[14,15].

**Naive Bayes Classifier**: This method models each class as a separate probability distribution. If the class distribution is compact (i.e., low variance), the overlap between class distributions decreases, simplifying class separation.

**Unsupervised Learning**

In unsupervised learning, where classes are not predefined, compactness measures are vital for data structure analysis and proper clustering.

**K-Means**: Aims to minimize the distance between points and the cluster center. The algorithm effectively maximizes compactness. Metrics like Within-Cluster Sum of Squares (WCSS) or Silhouette Score are often used for evaluation.

**DBSCAN**: Compactness is defined through density. Points at roughly equal distances from each other are grouped into dense clusters. Compact clusters have a dense core, while low-density points are considered noise.

**Hierarchical Clustering**: At each step, the closest groups are merged—essentially using compactness as the criterion. The choice of distance metric (e.g., single-linkage, complete-linkage) affects the internal compactness of clusters.

As noted earlier, generalization ability is one of the most crucial factors defining model performance on test data. The way a model learns the structure and density of samples within a class—i.e., the degree of compactness—is directly related to its generalization capability[9].

**Key Metrics for Measuring Compactness**

**Intra-class distance**: This metric measures the average distance between all points within the same class (or cluster). A small value indicates that the samples are located close to each other, suggesting a high degree of compactness.

$$Intra-class\,Dis = \frac{1}{\Sigma_k \binom{nk}{2}} \sum_{k=1}^{C} \sum_{i<j} \left\| x_i - x_j \right\| \qquad (1)$$

Where:
$C$ — number of classes
$n_k$ — number of elements in the $k$-th class
$d(x_i, x_j)$ — Euclidean distance between points $i$ and $j$

**Inter-class distance**: This metric measures the distance between the centers (or representatives) of different classes. The greater this distance, the more clearly distinguishable the classes are, indicating good separation.

The formula is as follows:

$$Inter-class\,Dis = \frac{2}{C(C-1)} \sum_{i<j} \left\| \mu_i - \mu_j \right\| \qquad (2)$$

Where:
$\mu_i$ – centroid of the $i$ - th class
$d$ – Euclidean distance between the centers of two classes

**Silhouette Score**: For each point, the average intra-cluster distance and the average distance to the nearest other cluster are computed. The score ranges from –1 to 1. A value close to 1 indicates high compactness and good cluster separation.

The formula is:

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))} \qquad (3)$$

Where:
$a(i)$ — average distance between point $i$ and all other points in the same cluster
$b(i)$ — average distance from point $i$ to the nearest neighboring cluster

The average of $s(i)$ over all points is used as the overall Silhouette Score for the dataset.

There are several widely used metrics for evaluating compactness, such as Intra-class Distance, Inter-class Distance, and Silhouette Score. These are primarily based on statistical and geometric analysis. While these metrics are effective in expressing intra-cluster/class closeness, inter-class separation, and shape consistency, they do not always correlate directly with a model's generalization ability. This is particularly true when they fail to account for data structure complexity or the specifics of the learning algorithm.

Moreover, the very concept of compactness lacks a strictly fixed definition — it can be interpreted differently depending on the nature of the data, the model's working principle, or the analysis goal. This leads us to the following important conclusion: there is no universal approach to measuring compactness. Instead, in each specific context, adaptive and goal-oriented metrics must be developed to match the task at hand.

Based on this, in this article, I propose a new approach for measuring compactness. This method retains the strengths of existing techniques while aiming to overcome their limitations. In particular, it enables a more precise assessment of a model's generalization capability.

**COMPUTATIONAL EXPERIMENT**

The table below presents the results of experiments conducted on three popular datasets — Iris, Wine, and Breast Cancer. The values in the datasets were normalized to the range [0;1]. Each dataset was analyzed from the perspective of clustering, and its compactness was evaluated using three classical metrics (Intra-class Distance, Inter-class Distance, Silhouette Score), as well as a new proposed compactness indicator based on connectivity ratios.

In addition, classification was performed on each dataset using the KNN, SVM, and Naive Bayes algorithms, and the models' generalization abilities were evaluated using accuracy[4].

The table provides both overall compactness values for each metric and individual values per class, allowing for a deeper analysis of class differences and structure. The values in parentheses reflect intra-class compactness. From the results obtained, a correlation between compactness values and model accuracy can be observed.

Table 1. Table: Comparison of Compactness Measures and Classification Accuracy Across Datasets

| Dataset name | Intra-class distance | Inter-class distance | Silhouette Score | Accuracy) | | |
|---|---|---|---|---|---|---|
| | | | | KNN | SVM | Naive bayes |
| Iris | 0.8479 (0.69, 0.84, 1.0) | 0.7637 (0.9376, 0.5846, 0.7690) | 0.4570 (0.70,0.38, 0.27) | 1.0 | 1.0 | 1.0 |
| Wine | 0.8636 (0.72,1.0,0.82) | 0.8929 (0.9192,0.7913,0.9681) | 0.2923 (0.40, 0.13,0.38) | 0.94 | 1.0 | 1.0 |
| Breast Cancer | 0.8230 (1.0, 0.71) | 0.9740 (0.97, 0.97) | 0.3389 (0.18,0.43) | 0.96 | 0.98 | 0.96 |

As shown in the table above, each dataset is compared in terms of the accuracy of three popular classification models (KNN, SVM, Naive Bayes) and various compactness metrics (e.g., value (1), value (1), value (2), value (3)).

The analysis shows that for datasets with high compactness metrics, the classification model accuracy is generally also higher.

For example, in the Iris dataset, the accuracy of all models is 1.0. In the Breast Cancer dataset, compactness is also high (0.84), and model accuracy ranges from 0.96 to 0.98. In the Wine dataset, the compactness indicators are slightly lower, which led to a decrease in KNN model accuracy to 0.94.
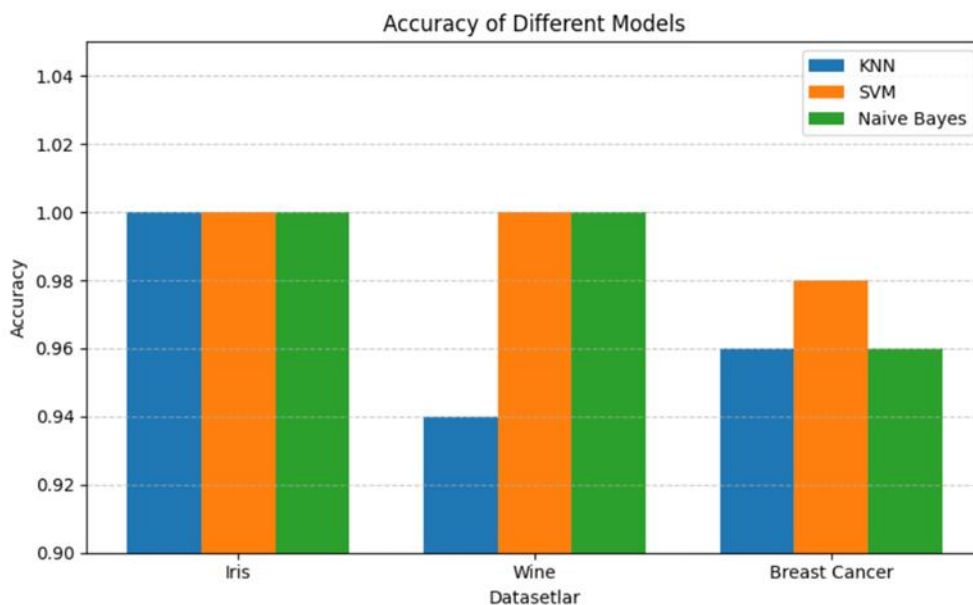
Figure 1 shows the accuracy values for different datasets.

These cases demonstrate a positive correlation between compactness and model accuracy. That is, when the data within classes are denser (i.e., more compact), classifiers are better able to distinguish those classes. This is especially noticeable for models based on geometric or statistical proximity, such as KNN and Naive Bayes.

**CONCLUSION**

This study provides a comprehensive analysis of the concept of *compactness* as one of the key factors influencing a model's generalization ability in machine learning. Compactness refers to the degree to which data within a class are densely, orderly, and tightly grouped. It plays a crucial role in determining how well a model will perform on new, previously unseen data.

Furthermore, the results of this study show that the use of new compactness evaluation methods enables the early assessment and selection of appropriate models. The outcomes of the experiments confirm that evaluating compactness can be an effective tool for model selection and can serve as a means of predicting a model's generalization performance in machine learning tasks.

**REFERENCES**

1. Zagoruiko N.G. Hypotheses of compactness and λ-compactness in data analysis methods // Sib. Zh. Industr. Mathematics, Vol. 1. – No. 1. – 1998. – P. 1 14–126.

2. Zhuravlev Yu.I. On algebraic methods in recognition and classification problems // Recognition, classification, forecasting. Mathematical methods and their application. – 1989. − P. 9-16.

3. Jambeu M. Hierarchical cluster - analysis and correspondence // Translated from French. - M.: Finance and Statistics, - 1988. - 342 p.

4. https://www.kaggle.com/datasets

5. Duda R. O., Hart P. E., Stork D. G. Pattern Classification. 2nd ed. - Wiley-Interscience, 2001. - 654 p

6. Bishop, C. M. (2006). Pattern Recognition and Machine Learning. Springer.

7. Hastie, T., Tibshirani, R., & Friedman, J. (2009). The Elements of Statistical Learning: Data

Mining, Inference, and Prediction. Springer.

8. Tan, P.-N., Steinbach, M., & Kumar, V. (2019). Introduction to Data Mining. Pearson.

9. Xu, R., & Wunsch, D. (2005). Survey of clustering algorithms. IEEE Transactions on Neural Networks, 16(3), 645–678.

10. Rousseeuw, P. J. (1987). Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. Journal of Computational and Applied Mathematics, 20, 53–65.

11. Biehl, M., Hammer, B., & Villmann, T. (2016). Prototype-based models in machine learning. Wiley Interdisciplinary Reviews: Cognitive Science, 7(2), 92–111.

12. Goodfellow, I., Bengio, Y., & Courville, A. (2016). Deep Learning. MIT Press.

13. Schölkopf, B., & Smola, A. J. (2002). Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond. MIT Press.

14. Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. Journal of Machine Learning Research, 12, 2825–2830.