

# Bag of Visual Words Classifier

Computer Vision and Pattern Recognition Exam

University of Trieste (UniTS)

Marco Tallone

January 2025

## Abstract

This report presents the implementation of a Bag of Visual Words (BoW) image classifier. The objective is to build a classifier for scene recognition by building a visual vocabulary from a set of images and performing multi-class classification. The visual vocabulary is built by clustering SIFT descriptors extracted from the images, and classification is performed comparing K-Nearest Neighbors (KNN) and Support Vector Machines (SVM) classifiers with different kernels. Results show that the best performance is achieved by the SVM classifiers, in particular when implemented with a spatial pyramid feature representation to add spatial information to the classic BoW approach.

## 1 Introduction

The Bag of Visual Words (BoW) model is a popular computer vision technique used for image classification or retrieval. Inspired by the analogous model used in natural language processing, this approach is based on the idea of treating images as collections of visual words belonging to a visual vocabulary, which is obtained by clustering local features extracted from a (possibly different) set of representative images.

This report presents an implementation of the BoW image classifier for scene recognition by first building a visual vocabulary from a set of test images and then performing multi-class classification using K-Nearest Neighbors (KNN) and Support Vector Machines (SVM) classifiers.

In particular, the visual vocabulary is built by clustering SIFT descriptors extracted from the test images using the K-Means algorithm. Descriptors are computed both from keypoints detected with the SIFT detector and from ones obtained by dense sampling of the images with a fixed grid.

In the classification phase, the performance of the KNN classifier is compared with that of different SVM classifiers, all using the “*one-vs-all*” strategy with different kernels.

Additionally, different ways to represent images as input feature vectors are tested. These include the classic representation as normalized histograms of visual words, the implementation of the *soft assignment* techniques proposed by *Van Gemert et al.* [7] and the use of the *spatial pyramid feature representation* proposed by *Lazebnik et al.* [3].

The objectives of this study are to compare the performance of the different classifiers and image representations and to reproduce the results obtained by *Van Gemert et al.* [7] and *Lazebnik et al.* [3] on the *15-Scenes* dataset.

The report is organized as follows. Section 2 presents the dataset used and analyzes the images distribution. Section 3 explains how the visual vocabulary is built from the sampled descriptors, while in Section 4 the different image representation techniques are compared. Then, Section 5 presents the classifiers used for the multi-class classification task, while in Section 6 the results obtained in this study are summarized, followed by some final considerations in the last section.

From a technical point of view, all the classifiers have been implemented in Python 3.12 adopting the `scikit-learn` library [4] (version 1.5.2). In particular, the SVM classifier has been implemented using the `SVC` class, which internally relies on `libsvm` [1]. However, with the aim of maintaining a clear and concise report of the work, the following sections will only present concepts from a theoretical point of view and the most noticeable results of the study, while further analysis and implementation details will be available on the author’s GitHub repository [6].

## 2 Dataset Description

The dataset used in this study is the *15-Scenes* dataset from *Lazebnik et al.* [3]. This dataset consists of a total of 4485 grayscale images, each belonging to one of 15 different scene categories (*Office, Kitchen, LivingRoom, Bedroom, Store, Industrial, TallBuilding, InsideCity, Street, Highway, Coast, OpenCountry, Mountain, Forest, Suburb*). The images are divided into a training set of 1500 images, with 100 images per category, and a test set of 2985 images with a non-uniform distribution as shown in Figure 2.1. Moreover, it's worth mentioning that the images in the dataset have different sizes and aspect ratios, with the majority of them being  $256 \times 256$  pixels.

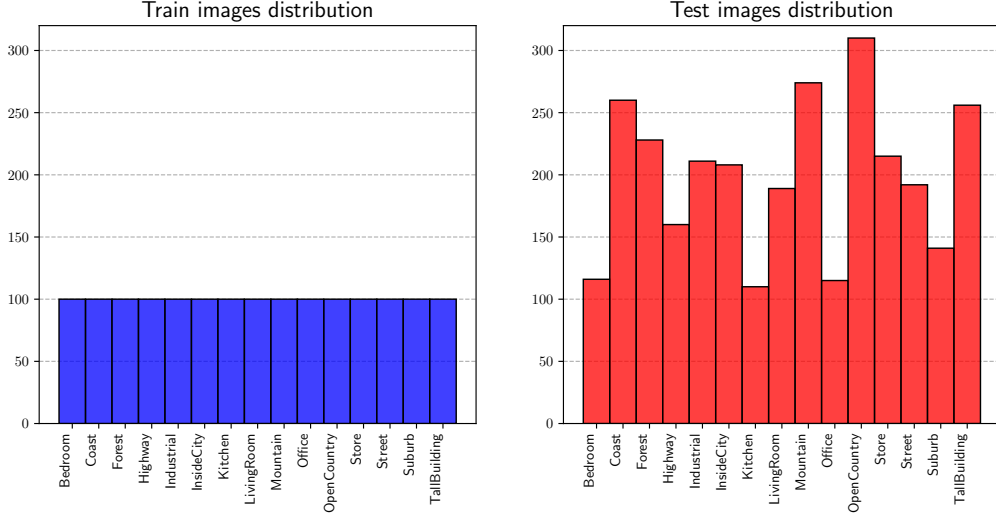


Figure 2.1: Distribution of images in train (*blue*) and test (*red*) sets.

## 3 Visual Vocabulary Construction

The process of building a visual vocabulary can be split in two main steps:

1. **Feature extraction:** sampling of SIFT descriptors from train images.
2. **Clustering:** grouping of the extracted descriptors into  $K$  clusters.

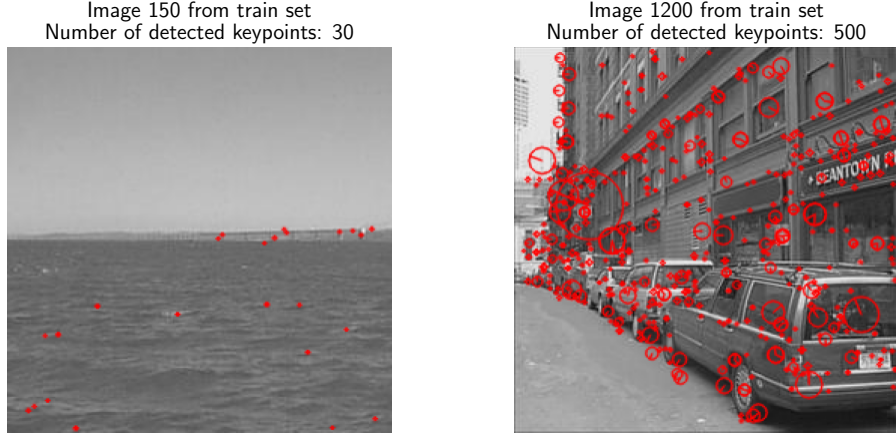
### 3.1 Feature Extraction

The first step involves the extraction of local features from the images in the training set. In this case, SIFT descriptors have been used as features, and they have been extracted following two different approaches:

1. Using the **SIFT detector**: the SIFT detector is used to detect keypoints in the images and the corresponding SIFT descriptors are computed on these keypoints.
2. Sampling on a **dense regular grid**: following the approach adopted by *Lazebnik et al.* [3], SIFT descriptors can also be computed from keypoints sampled on a regular dense grid over each image. Following the approach of the paper, a regular grid with spacing of 8 pixels between each keypoint has been used.

Both sampling methods have been tested and compared, with the results shown in later section 6. For the first approach, the SIFT detector has been initialized to retain only the best 500 features in each image, resulting in a total of 593 006 keypoints, while for the second approach 1 482 434 keypoints have been collected.

Notice that with both approaches the exact number of keypoints extracted in each image can change from one image to another. In the first case, this is due to the content of the images, i.e. the actual number of features that can be detected by the SIFT detector. In the second case instead, this is due to the different sizes and aspect ratios of the images. Figure 3.1 shows an example of this phenomenon, where a different number of keypoints has been detected by the SIFT detector in two different images, while Appendix A1 shows the distribution descriptors sampled from the train set with the two methods.



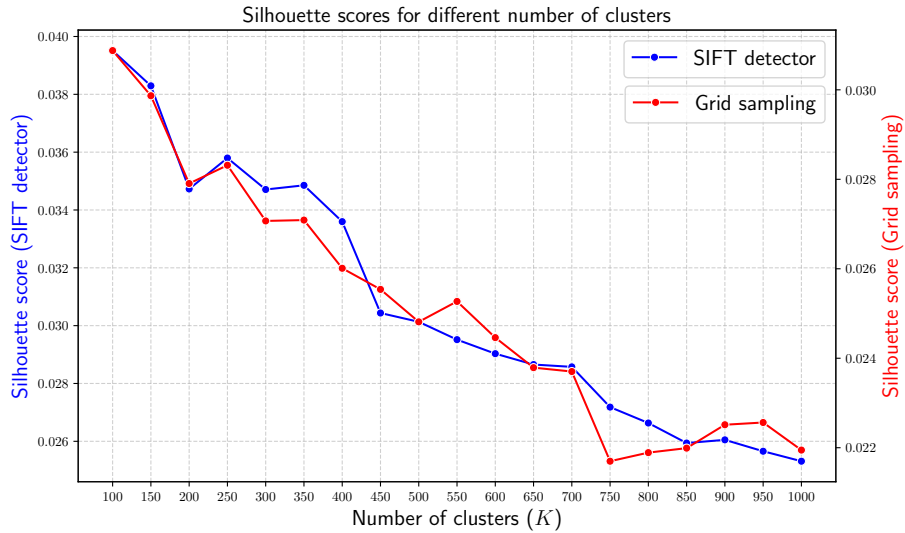
**Figure 3.1:** Example of keypoints detection with the SIFT detector. In the left side image only 30 keypoints were detected, while in the right side one all the desired 500 keypoints were found.

### 3.2 Clustering

Following feature extraction, clustering has been performed on a random subset of extracted SIFT descriptors to build the visual vocabulary. Hence, the resulting centroids, representing the visual words of the vocabulary, are 128-dimensional vectors.

The size of the random subset has been chosen to guarantee the same ratio between total number of descriptors and number of descriptors used for clustering in the two feature extraction methods. Hence, this has been set to 10 000 for descriptors sampled using the SIFT detector and to 25 000 for descriptors sampled on a regular grid.

Clustering has then been performed using the *k-means* algorithm to group descriptors into  $K$  clusters, where  $K$  becomes the size of the visual vocabulary. The number of clusters is an hyperparameter of the clustering process which has to be chosen. An analysis of the average silhouette score [5] has been performed for values of  $K$  ranging from 100 to 1000 with the results shown in figure 3.2. Ultimately, the value of  $K = 400$  has been chosen as a good trade-off between the quality of the clustering and the overall computational cost of the clustering process. Moreover, the choice is also motivated by the final results obtained in Section 6 and by the fact that, given the objectives of this project, this value allows for a direct comparison with the results reported by *Lazebnik et al.* [3] and *Van Gemert et al.* [7], where the same vocabulary size was used.



**Figure 3.2:** Silhouette scores resulting from clustering SIFT descriptors sampled using the SIFT detector (*blue*) and on a regular grid (*red*). As expected, the value decreases with the increase of  $K$ . Noticeable decreases occur after  $K \approx 400$  and  $K \approx 700$ .

## 4 Image Representation

A important step in the Bag of Visual Words (BoW) pipeline is the representation of images as fixed-length feature vectors. In this study, 4 different techniques have been implemented to perform this task.

The classic approach is to represent images as *normalized histograms of visual words (HIST)*. For each image, the previously extracted descriptors are assigned to the closest visual word in the vocabulary, and a histogram is built counting the occurrences of each visual word in the image. Formally, the  $i$ -th bin of each histogram is computed as

$$h_i = \frac{n_{id}}{n_d}, \quad \forall i \in \{0, \dots, K-1\}, \forall d \in \{0, \dots, N-1\} \quad (4.1)$$

where  $n_{id}$  is the number of occurrences of the  $i$ -th visual word in image  $d$ ,  $n_d$  is the total number of visual words in the image, and  $N$  is the total number of images in the dataset. After normalization, the result is a fixed-length representation in which images are seen as normalized histograms having  $K$  bins, corresponding to visual words frequencies.

In order to account for the relevance of the visual words, another idea is to use the *term frequency-inverse document frequency (TF-IDF)* weighting scheme. In this case, given the number of images  $N_i$  containing the  $i$ -th visual word, the histogram bins are

$$h_i = \frac{n_{id}}{n_d} \cdot \log \left( \frac{N}{N_i} \right), \quad \forall i \in \{0, \dots, K-1\}, \forall d \in \{0, \dots, N-1\} \quad (4.2)$$

A third approach is to use the *soft assignment* techniques proposed by Van Gemert *et al.* [7], which argue that the hard assignment of visual words to descriptors gives rise to two types of ambiguity: *codeword uncertainty* when a descriptor is similarly close to two or more visual words, and *codeword plausibility* when a descriptor is relatively far from all visual words. To address this issue, the authors propose the usage of a Gaussian kernel  $K_\sigma(x)$  density estimator<sup>1</sup> based on the Euclidean distance  $D(w_i, x_j)$  between each  $i$ -th visual word  $w_i$  and  $j$ -th descriptor  $x_j$  to approximate the probability density function of visual words in the images. Hence, by replacing the histogram estimator with the kernel density estimator the result is the *kernel codebook (KCB)* representation, in which each descriptor can contribute to multiple bins according to the equation

$$h_i = \frac{1}{n_{ri}} \sum_{j=0}^{n_{ri}} K_\sigma(D(w_i, x_j)), \quad \forall i \in \{0, \dots, K-1\}, \forall d \in \{0, \dots, N-1\} \quad (4.3)$$

A visual comparison between the representations obtain with hard assignment and soft assignment methods for a sample image is shown in Figure 4.1. According to the authors, this soft assignment strategy models the two types of ambiguity at the same time, but it's also possible to consider them separately by adopting either the *codeword uncertainty (UNC)* representation

$$h_i = \frac{1}{n_{ri}} \sum_{j=0}^{n_{ri}} \frac{K_\sigma(D(w_i, x_j))}{\sum_{l=0}^{K-1} K_\sigma(D(w_l, x_j))}, \quad (4.4)$$

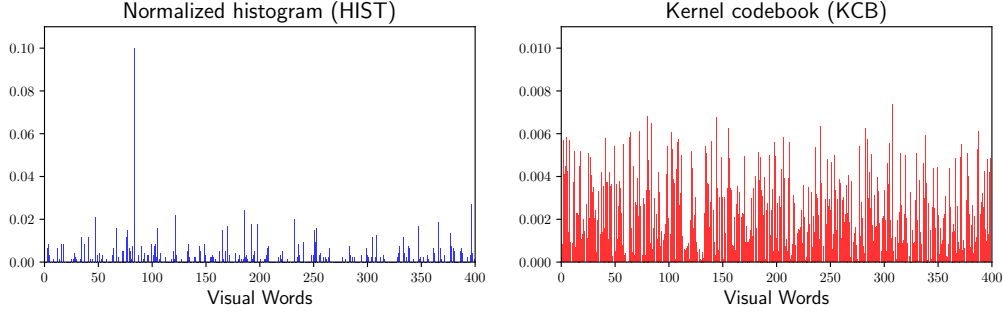
or the *codeword plausibility (PLA)* representation

$$h_i = \frac{1}{n_{ri}} \sum_{j=0}^{n_{ri}} \begin{cases} K_\sigma(D(w_i, x_j)), & \text{if } w_i = \underset{l \in \{0, \dots, K-1\}}{\operatorname{argmin}} D(w_l, x_j) \\ 0, & \text{otherwise} \end{cases} \quad (4.5)$$

Notice that in equations 4.3, 4.4, and 4.5 the sums are not performed for all the  $n_d$  descriptors in the image, but only over the  $n_{ri}$  descriptors in the region around the  $i$ -th visual word. In the original paper, in fact, the authors employed *radius based clustering* to construct the visual vocabulary. Since in this study *k-means* has been used, an approximation of each cluster radius has been obtained as the distance<sup>2</sup> between the cluster centroid and the furthest descriptor assigned to it. The region around each visual word is hence defined as the set of descriptors whose distance from the centroid is less than the corresponding radius. Such definition possibly implies overlapping regions, meaning that a descriptor can belong to more than one region and hence contribute to more than one histogram bin. Empirically, this led to better results than including all the descriptors of the image in the computation.

<sup>1</sup>See Appendix A2.

<sup>2</sup>In Euclidean norm



**Figure 4.1:** Comparison between normalized histogram (*blue*) and kernel codebook (*red*) representations for image 1200 in train set.

The last approach used to fulfill the image representation task is the use of *spatial pyramid matching* (SPM) method proposed by *Lazebnik et al.* [3]. The idea is to repeatedly subdivide each image into subsequently finer grids at different resolution levels  $\ell \in \{0, \dots, L-1\}$ , and to compute the histograms  $H_\ell(g)$  of visual words for each  $g$ -th grid at level  $\ell$ . Precisely, at level  $\ell$  each image is divided into  $2^\ell$  cells along each dimension, and the count of visual words is computed for each cell. All the histograms are then weighted according to the weighting scheme in equation 4.6, that favors features counts computed at higher levels of the pyramid.

$$\omega_0 = \frac{1}{2^L}, \quad \omega_\ell = \frac{1}{2^{L-\ell+1}}, \quad \forall \ell \in \{0, \dots, L-1\} \quad (4.6)$$

The resulting weighted histograms are then stacked together to form, for each image, a single *extended multi-level descriptor* that encodes both visual words and positional information at different resolutions and can be used as input to SVM classifiers. The idea of using SVM classifiers specifically for this representation comes from the possibility of adopting the *spatial pyramid kernel* in order to compute the similarity between these extended descriptors by correctly measuring the information at different levels of the pyramid. Formally, given two images  $X$  and  $Y$  represented by their extended descriptors of elements  $H_{\ell,i}^X(g)$  and  $H_{\ell,i}^Y(g)$ , the *pyramid match kernel* between the two is

$$K^L(X, Y) = \sum_{\ell=0}^{L-1} \sum_{i=0}^{K-1} \sum_{g=0}^{2^{2\ell}-1} \min(H_{\ell,i}^X(g), H_{\ell,i}^Y(g)) \quad (4.7)$$

which can be efficiently computed as a simple histogram intersection kernel.

Finally, it's important to notice that for this last representation the features are extracted from images only by using the grid approach. This has been done both to follow the same approach of the original paper and also to avoid the possibility of having empty cells which might happen for those containing uniform regions of the image<sup>3</sup>.

## 5 Classification

In order to perform the scene recognition task, two main classes of classifiers have been implemented: K-Nearest Neighbors (KNN) and Support Vector Machines (SVM).

### 5.1 K-Nearest Neighbors (KNN)

A K-Nearest Neighbors classifier taking as input features the normalized histograms of visual words has been implemented. The classifier assigns to each image in the test set the label of the majority class among its  $k$  nearest neighbors<sup>4</sup> in the training set.

In the simple case of  $k = 1$ , the label of the closest histogram in the training set is assigned to the test image. A slightly better result can be achieved by performing a linear search for the hyperparameter  $k$  over the range  $[1, 50]$  using the average accuracy as assessment metric. For each value of  $k$  in the range, the accuracy of the corresponding KNN classifier is computed and stored. At the end, the value of  $k$  that maximizes the accuracy is selected and the performance of the resulting classifier is evaluated.

<sup>3</sup>E.g. cells positioned in the top region of the left side image or in the bottom left corner of the right side image in Figure 3.1.

<sup>4</sup>Notice  $k$  is the number of KNN neighbors, not to be confused with the number of clusters  $K$ .

## 5.2 Support Vector Machines (SVM)

A series of multi-class Support Vector Machine (SVM) classifiers have been implemented following the “one-vs-all” strategy. For each possible class, a single binary classifier is trained and the final prediction is obtained by selecting the class with the highest confidence score. Each binary classifier is trained taking as input features one of the presented image representations and modified ground truth labels where the class of interest is labeled as +1 and all other classes are labeled as -1.

Different kernels for these SVM classifiers have been tested and compared. Initially, the default radial basis function (RBF) kernel has been adopted. Its performance has then been compared with the generalized Gaussian kernel (with  $\gamma = 1/2$ ) based on the  $\chi^2$  distance and the histogram intersection kernel in equation 5.1, both widely adopted in histogram comparison tasks.

$$k_{\chi^2}(\mathbf{x}, \mathbf{x}') = \exp \left( -\gamma \sum_i \frac{(x_i - x'_i)^2}{x_i + x'_i} \right) \quad k_{\cap}(\mathbf{x}, \mathbf{x}') = \sum_i \min(x_i, x'_i) \quad (5.1)$$

For the *SPM* feature representation, the SVM classifiers only used the histogram intersection kernel as anticipated in the previous section and pyramid levels up to  $L = 2$ .

## 6 Results

The results obtained with all the classifiers, the two feature extraction methods and the different image representations are summarized in Table 6.1 while Figure 6.1 compares the results for the grid sampling method.

The average accuracy over all the classes has been used as the main assessment metric for all the models. All the experimental results have been obtained by running the models on a machine equipped with an Intel® Core™ i7-8565U CPU @ 4.60 GHz and 8 GB of RAM.

The dummy classifier, that always predicts the most frequent class (*OpenCountry*) on the test set, has been used as baseline for comparison for the other models and the average accuracy it achieved is 10.39 %.

First of all, it’s possible to notice that the results obtained using grid sampling consistently surpass the ones resulting from the SIFT detector. This confirms the idea presented in *Fei-Fei and Perona* [2] and *Lazebnik et al.* [3] for which grid sampling works better for scene classification tasks, as it allows to capture uniform regions such as sky, water, or road surfaces that might be crucial in discriminating between classes. Looking at the performance of the different classifiers instead, the KNN ones surpassed the baseline, recording average accuracies between 40 % and 50 % with a small gain in performance given by the multiple neighbors approach over the single neighbor classifier. Undoubtably, much better results have been achieved by the SVM classifiers with accuracies ranging from 50 % to 70 %. In particular, the adoption of specialized kernels significantly boosted the performance of the SVM classifiers, at least in the *HIST* and *TF-IDF* representations. Indeed, the 71.42 % accuracy obtained with the  $\chi^2$  kernel and the 70.62 % with the intersection one are in agreement with the findings of *Van Gemert et al.* [7] for the hard assignment case.

However, from the comparison of the different image representation techniques it emerges that, while little to no differences are observed between the *HIST* and *TF-IDF* representations, the implementation of the soft assignment techniques failed in improving the classification performances as reported in *Van Gemert et al.* [7], except for some isolated cases. In contrast with the results of the paper, the accuracies measured for these representations are mostly comparable, if not worse with respect to the ones of the hard assignment cases.

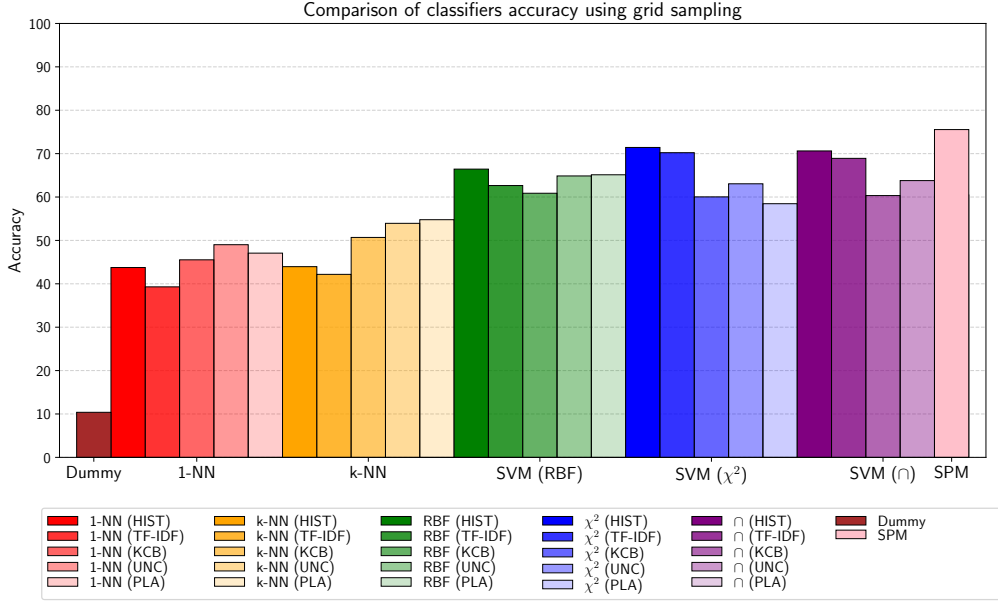
A further level of assessment is also given by the confusion matrices computed for all the classifiers<sup>5</sup> and reported in appendix A3. These, not only confirm the superiority of the SVM classifiers with respect to the KNN ones, but also highlight the most difficult scenes to recognize. In detail, in most of the classifiers a noticeable (but understandable) confusion arises between the *Open Country* and *Coast* classes as well as among the *LivingRoom*, *Bedroom*, *Kitchen* and *Office* classes. Consistently, the hardest scene to classify has been the *Industrial* one, which matches with the results from literature. By far, the best accuracy has been obtained by the SVM classifiers using the spatial

<sup>5</sup>Only for *HIST* representation and grid sampling method as explained in appendix A3.

pyramid matching approach. In this case, the combination of sampling features from a regular grid and adopting the *SPM* representation for the images allowed to achieve an accuracy of 75.54 %. This result is significantly better than the other classifiers and is a clear indication that the spatial information is crucial for the scene recognition task. The superiority of this approach is also confirmed by the almost completely diagonal confusion matrix (Figure A.7), in which only few of the previous misclassifications remain, mainly between the *Living Room* and *Bedroom* classes. Although higher than the other classifiers, the lowest accuracy is still measured for the *Industrial* class.

Classifier	Image Representation	Feature Extraction	
		<i>SIFT detector</i>	<i>Grid sampling</i>
Dummy Classifier	<i>HIST</i>	10.39 %	10.39 %
	<i>TF-IDF</i>	10.39 %	10.39 %
	<i>KCB</i>	10.39 %	10.39 %
	<i>UNC</i>	10.39 %	10.39 %
	<i>PLA</i>	10.39 %	10.39 %
1-NN	<i>HIST</i>	31.49 %	43.75 %
	<i>TF-IDF</i>	31.83 %	39.30 %
	<i>KCB</i>	33.40 %	45.53 %
	<i>UNC</i>	35.38 %	49.01 %
	<i>PLA</i>	36.72 %	47.07 %
k-NN	<i>HIST</i>	37.19 %	43.95 %
	<i>TF-IDF</i>	36.42 %	42.18 %
	<i>KCB</i>	38.93 %	50.69 %
	<i>UNC</i>	42.31 %	53.94 %
	<i>PLA</i>	43.05 %	54.77 %
SVM (RBF)	<i>HIST</i>	50.65 %	66.43 %
	<i>TF-IDF</i>	50.89 %	62.65 %
	<i>KCB</i>	54.77 %	60.87 %
	<i>UNC</i>	56.08 %	64.86 %
	<i>PLA</i>	51.39 %	65.13 %
SVM ( $\chi^2$ )	<i>HIST</i>	52.66 %	71.42 %
	<i>TF-IDF</i>	52.03 %	70.22 %
	<i>KCB</i>	53.94 %	60.03 %
	<i>UNC</i>	54.37 %	63.05 %
	<i>PLA</i>	45.06 %	58.46 %
SVM ( $\cap$ )	<i>HIST</i>	50.95 %	70.62 %
	<i>TF-IDF</i>	50.79 %	68.91 %
	<i>KCB</i>	52.23 %	60.34 %
	<i>UNC</i>	53.37 %	63.79 %
	<i>PLA</i>	49.15 %	60.47 %
SVM ( <i>Spatial Pyramid</i> )	<i>SPM</i>	—	75.54 %

**Table 6.1:** Average accuracies for all the implemented classifiers comparing feature extraction methods and image representations.



**Figure 6.1:** Comparison of the classifiers accuracies using grid sampling.

## 7 Conclusions

In conclusion, this study has shown the implementation of a Bag of Visual Words (BoW) classifier for scene recognition from the construction of a visual vocabulary by performing clustering on SIFT descriptors to the classification of images using machine learning techniques such as K-Nearest Neighbors (KNN) and Support Vector Machines (SVM). From the results, it emerges that the SVMs have clearly outperformed the KNN classifiers in terms of accuracy, but the latter at least surpassed the dummy classifier performance used as baseline. In particular, the usage of specialized kernels for the SVM classifiers has also shown to be beneficial with respect to the default RBF kernel.

Concerning the worse results recorded for the soft assignment techniques, these might be explained the different approach used in this study with respect to the original paper. Specifically, since *k-means* has been used instead of the *radius based* clustering implemented by the original authors, it's possible that the approximation used to define the region around each centroid might have led to worse results. Moreover, *Van Gemert et al.* [7] conducted *k-fold* cross-validation on the training set to find the optimal value of  $\sigma$  for the Gaussian kernel, while in this study the value has been set to  $\sigma = 200$  based on the results of the paper due to computational and time constraints. Additionally, as done in this study, the authors also performed clustering on a random subset of the extracted descriptors, but it's possible that the differences in the subset size and in the representativeness of the sampled descriptors might have led to different results.

An improvement and further extension of this project could in fact be to repeat the random sampling and clustering multiple times with different random seeds in order to collect significant statistics about the performance of all the techniques and hence perform a more robust comparison.

Ultimately, the best accuracy has been achieved by the SVM classifiers using the spatial pyramid approach proposed by *Lazebnik et al.* [3]. Although the final accuracy of 75.54 % doesn't quite match the one obtained by the authors of the original paper, in this case the results are satisfactory since they are, by a good margin, better than the ones of the other classifiers. This confirms the importance of adding spatial information to the classic BoW approach. Also in this case, the slight differences in the final results might be due to the different parameters used in the clustering process. Nevertheless, the results obtained are still satisfactory and confirm the validity of the BoW approach for the scene recognition task.

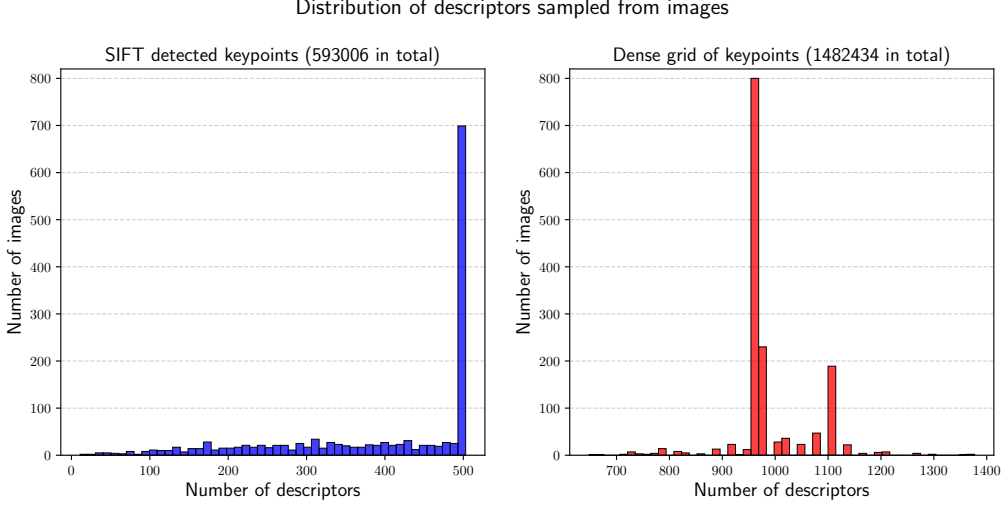


# References

- [1] Chih-Chung Chang and Chih-Jen Lin. Libsvm: A library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011.
- [2] L. Fei-Fei and P. Perona. A bayesian hierarchical model for learning natural scene categories. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 2, pages 524–531 vol. 2, 2005. doi:10.1109/CVPR.2005.16.
- [3] S. Lazebnik, C. Schmid, and J. Ponce. Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, volume 2, pages 2169–2178, 2006. doi:10.1109/CVPR.2006.68.
- [4] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [5] Peter J. Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20:53–65, 1987. URL: <https://www.sciencedirect.com/science/article/pii/0377042787901257>, doi:10.1016/0377-0427(87)90125-7.
- [6] Marco Tallone. Bow classifier. URL: <https://github.com/marcotallone/bow-classifier>.
- [7] Jan C. van Gemert, Jan-Mark Geusebroek, Cor J. Veenman, and Arnold W. M. Smeulders. Kernel codebooks for scene categorization. In David Forsyth, Philip Torr, and Andrew Zisserman, editors, *Computer Vision – ECCV 2008*, pages 696–709, Berlin, Heidelberg, 2008. Springer Berlin Heidelberg.

# A Appendix

## A1 Descriptors Distribution



**Figure A.1:** Distribution of keypoints sampled from the train set using the SIFT detector (*blue*) and the dense grid approach (*red*). For the first method, the desired number of keypoints (500) has been detected in the majority of the images, while for the second method the number of keypoints per image is concentrated around 1000 with small differences due to the different sizes and aspect ratios of the images.

## A2 Kernel Density Estimation

The Gaussian density estimator  $K_\sigma(x)$  is defined as

$$K_\sigma(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{x^2}{2\sigma^2}\right) \quad (\text{A.1})$$

Whereas, given the fact that both visual words and SIFT descriptors are 128-dimensional vectors, the Euclidean distance  $D(w_i, x_j)$  between the  $i$ -th visual word  $w_i$  and the  $j$ -th descriptor  $x_j$  can be introduced as

$$D(w_i, x_j) = \|w_i - x_j\| = \sqrt{\sum_{k=1}^{128} (w_{ik} - x_{jk})^2} \quad (\text{A.2})$$

Hence the soft assignment approach proposed by *Van Gemert et al.* [7] relies on the density estimator

$$K_\sigma(D(w_i, x_j)) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{D(w_i, x_j)^2}{2\sigma^2}\right) \quad (\text{A.3})$$

### A3 Confusion Matrices

Here the confusion matrices for all the classifiers are reported. These matrices refer to the case in which features have been extracted by a SIFT descriptor from keypoints sampled on a dense grid on the images and the input representation is the normalized histograms of visual words (*HIST*). This has been done to maintain the results consistent with the ones obtained using the spatial pyramid matching approach and to be able to compare the performances of the different classifiers. However, similar matrices can be computed in the other cases as well.

		Confusion Matrix of 1-NN Classifier															
True Labels		Predicted Labels															
		Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding	
Bedroom		22	0	0	0	2	1	23	12	0	23	0	1	1	3	0	
Coast		8	31	1	34	3	2	3	1	1	16	10	3	1	26	0	
Forest		2	0	79	0	0	1	1	0	2	0	0	11	0	9	0	
Highway		3	2	0	68	4	0	2	1	1	4	1	1	0	13	1	
Industrial		9	0	1	4	18	1	21	11	1	25	2	7	7	21	4	
InsideCity		3	0	0	1	10	31	29	4	0	19	0	7	6	13	3	
Kitchen		5	0	0	0	0	1	47	12	0	21	0	0	0	1	0	
LivingRoom		17	0	0	0	3	1	61	29	0	27	0	3	1	0	0	
Mountain		7	0	7	6	9	2	5	3	45	3	6	7	6	21	2	
Office		4	0	0	0	0	0	19	10	0	57	0	0	0	1	0	
OpenCountry		3	10	10	12	9	1	3	1	7	2	32	13	2	41	2	
Store		3	1	4	1	8	4	22	7	1	10	1	41	2	11	3	
Street		4	0	0	2	10	3	4	6	1	6	1	6	54	8	1	
Suburb		1	0	0	0	2	0	7	3	0	8	0	1	2	77	0	
TallBuilding		5	1	1	2	9	4	28	11	0	16	0	10	4	1	45	

**Figure A.2:** Confusion matrix for the 1-NN classifier.

True Labels	Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding
Bedroom	24	0	0	0	0	0	27	14	0	23	0	0	0	1	0
Coast	20	35	0	33	5	2	3	1	1	12	7	0	1	23	0
Forest	1	0	87	0	0	0	0	0	1	0	0	7	1	4	0
Highway	6	4	0	67	1	1	3	1	0	3	0	1	0	14	0
Industrial	16	0	0	8	17	0	25	13	1	35	0	6	2	16	2
InsideCity	6	0	0	2	15	33	29	2	0	23	0	5	4	9	2
Kitchen	6	0	0	0	0	0	52	9	0	24	0	0	0	0	0
LivingRoom	17	0	0	0	1	1	62	30	0	30	0	2	0	0	0
Mountain	11	1	6	9	8	2	2	6	45	4	2	7	4	26	0
Office	3	0	0	0	0	0	21	8	0	63	0	0	0	0	0
OpenCountry	7	9	15	22	8	0	0	1	6	3	25	9	1	51	0
Store	9	0	5	0	7	1	21	10	1	15	0	46	1	9	0
Street	5	0	0	1	10	2	5	8	0	10	0	4	51	11	1
Suburb	6	0	0	1	1	0	7	3	0	7	0	0	0	76	0
TallBuilding	16	0	1	2	8	2	42	12	1	14	0	9	4	2	36
Predicted Labels	Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding

**Figure A.3:** Confusion matrix for the k-NN classifier with  $k = 5$ .

True Labels	Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding
Bedroom	43	0	0	0	2	1	11	7	1	10	0	1	2	1	4
Coast	1	80	0	13	0	0	0	0	2	0	7	0	0	1	0
Forest	0	2	92	0	0	0	0	0	3	0	0	0	0	1	1
Highway	0	7	0	79	0	0	0	0	1	2	1	1	1	2	0
Industrial	7	3	1	6	27	6	4	2	4	8	4	8	7	11	10
InsideCity	0	3	0	1	1	58	6	2	0	6	2	4	9	9	5
Kitchen	8	0	0	0	1	2	54	4	0	14	0	1	2	0	2
LivingRoom	30	0	0	1	1	2	20	32	0	22	0	6	2	5	5
Mountain	0	5	5	7	0	0	0	0	81	1	3	0	0	3	1
Office	5	0	0	0	1	0	5	2	0	76	0	0	0	4	1
OpenCountry	0	20	11	8	2	0	0	0	9	0	57	1	1	6	0
Store	3	2	6	1	1	7	5	2	4	3	2	56	4	5	5
Street	0	0	1	4	1	0	1	1	4	0	0	3	77	5	2
Suburb	1	1	0	1	0	1	1	1	0	2	0	0	2	88	0
TallBuilding	0	3	2	2	1	2	1	1	2	0	0	2	2	0	85
Predicted Labels	Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding

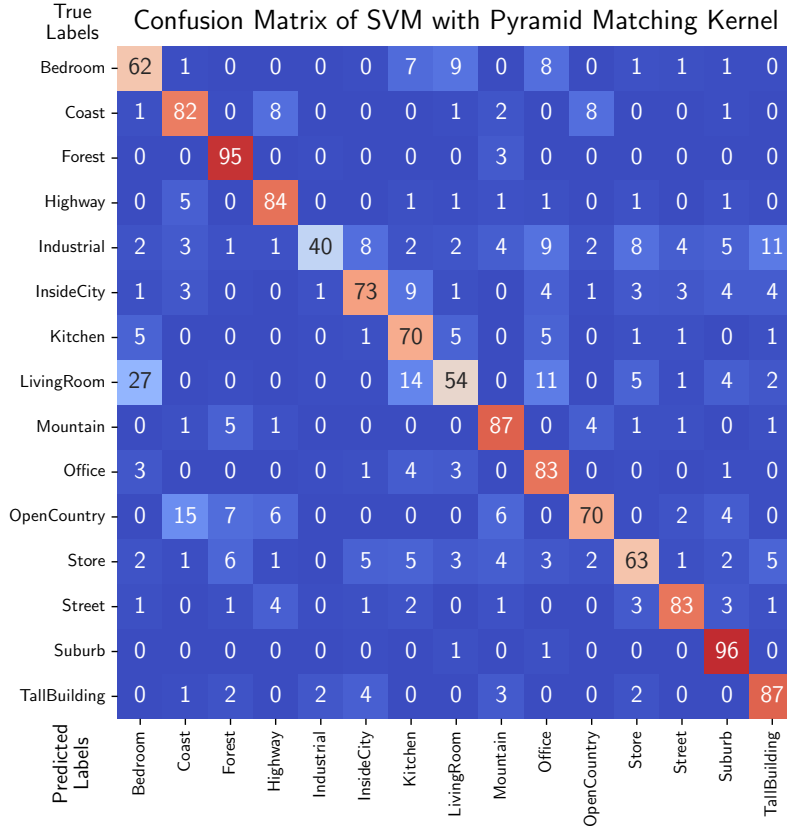
**Figure A.4:** Confusion matrix for the SVM classifier with RBF kernel.

True Labels	Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding
Bedroom	50	0	0	0	2	1	10	9	0	10	0	1	0	1	2
Coast	1	83	1	8	0	0	0	0	1	1	6	0	0	3	0
Forest	0	0	95	0	0	0	0	0	3	0	0	0	0	0	0
Highway	3	5	0	80	0	0	0	0	2	2	1	0	1	1	0
Industrial	4	3	1	6	34	5	3	4	3	12	3	7	5	6	10
InsideCity	0	1	0	1	0	72	5	2	0	4	1	4	5	8	4
Kitchen	9	1	0	1	1	2	52	5	0	13	0	1	2	0	1
LivingRoom	27	0	0	0	2	1	15	38	0	20	0	6	2	4	5
Mountain	0	2	5	3	0	0	0	0	86	0	4	0	0	1	0
Office	4	0	0	0	0	0	3	2	0	83	0	0	0	3	0
OpenCountry	0	16	9	5	1	0	0	0	8	0	65	1	1	6	0
Store	3	1	6	1	0	9	4	3	3	2	1	59	3	4	5
Street	0	0	1	3	0	0	2	1	2	0	1	3	79	9	1
Suburb	1	0	0	0	0	0	0	1	0	3	0	0	1	93	0
TallBuilding	0	0	2	0	2	2	0	1	3	0	0	2	1	0	88
Predicted Labels	Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding

**Figure A.5:** Confusion matrix for the SVM classifier with  $\chi^2$  kernel.



True Labels	Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding
Bedroom	45	0	0	0	3	1	9	8	1	12	0	1	1	2	3
Coast	1	85	1	8	0	0	0	0	1	1	5	0	0	2	0
Forest	0	0	95	0	0	0	0	0	3	0	0	0	0	1	0
Highway	2	5	0	80	0	0	0	0	2	2	1	1	1	1	0
Industrial	4	2	1	5	34	7	1	4	3	14	3	7	6	9	9
InsideCity	0	2	0	1	0	75	7	2	0	3	1	2	4	5	4
Kitchen	10	1	0	1	2	1	53	4	0	11	0	0	2	0	2
LivingRoom	29	0	0	1	3	1	15	33	0	22	0	6	2	6	5
Mountain	0	2	6	3	0	0	0	0	85	0	4	0	1	1	0
Office	3	0	0	0	0	1	5	2	0	81	0	0	0	2	0
OpenCountry	0	17	13	7	2	0	0	0	7	0	59	1	2	6	0
Store	4	1	6	1	0	9	5	2	4	2	1	58	3	4	5
Street	0	0	1	3	1	1	1	1	3	0	1	4	78	6	1
Suburb	0	0	0	0	0	0	0	1	0	2	0	0	0	96	0
TallBuilding	2	0	2	0	2	3	0	1	3	0	0	2	1	0	87
Predicted Labels	Bedroom	Coast	Forest	Highway	Industrial	InsideCity	Kitchen	LivingRoom	Mountain	Office	OpenCountry	Store	Street	Suburb	TallBuilding


**Figure A.6:** Confusion matrix for the SVM classifier with histogram intersection kernel.



**Figure A.7:** Confusion matrix for the SVM classifier with spatial pyramid matching approach.

#### A4 Notice on Generative Tools Usage

Generative AI tools have been used as a support for the development of this project. In particular, the Copilot  generative tool based on OpenAI GPT 4o  model has been used as assistance medium in performing the following tasks:

- writing documentation and comments in implemented functions by adhering to the NumPy Style Guide 
- improving variable naming and code readability
- minor bug fixing in implemented functions
- grammar and spelling check both in the repository README [6] and in this report
- tweaking aesthetic improvements in the plots of this report
- formatting table of results in this report