

---

# CENG 483

## Introduction to Computer Vision

Fall 2021-2022

### Take Home Exam 2

### Object Recognition

Student Random ID:

---

Please fill in the sections below only with the requested information. If you have additional things to mention, you can use the last section. Please note that all of the results in this report should be given for the **validation set**. Also, when you are expected to comment on the effect of a parameter, please make sure to fix other parameters.

## 1 Local Features (25 pts)

- Explain SIFT and Dense-SIFT in your own words. What is the main difference?
- Put your quantitative results (classification accuracy) regarding 5 values of SIFT and 3 values of Dense-SIFT parameters here. In SIFT change each parameter once while keeping others same and in Dense-SIFT change size of feature extraction region. Discuss the effect of these parameters by using 128 clusters in k-means and 8 nearest neighbors for classification.

Sift is a keypoint detector that is designed to capture features that are independent from the scale of the image. Its scale invariance property stems from difference of gaussian operations with different kernel size resolutions which captures keypoints at different scales. Then we can quantize an orientation for each keypoint and when orientations of all features are described wrt the keypoint orientation, then it'll additionally have rotation invariance alongside scale invariance. These are all desired invariances that are needed to select interesting or descriptive points in images in order to identify, match, classify etc. Dense-Sift is a kind of modification on top of conventional Sift method. In Sift locations of keypoint descriptors are determined by the default algorithm. On the contrary, in Dense-Sift center-positions (x,y) can be fed to Sift algorithm to create descriptors around that point, or equivalently one can partition the image into grid like portions and feed them individually through default Sift which in a way has similar behaviour, resulting in specific sift vectors around that spatial region.

Sift: Parameters vs Accuracy					
Accuracy (1.0)	nfeatures	nOctaveLayers	contrastThreshold	edgeThreshold	sigma
0.17	0	3	0.04	10	1.6
0.10	1	3	0.04	10	1.6
0.18	20	3	0.04	10	1.6
0.16	100	3	0.04	10	1.6
0.19	0	1	0.04	10	1.6
0.16	0	5	0.04	10	1.6
0.17	0	7	0.04	10	1.6
0.0	0	3	1	10	1.6
0.16	0	3	0.1	10	1.6
0.17	0	3	0.01	10	1.6
0.0	0	3	0.04	1	1.6
0.19	0	3	0.04	50	1.6
0.19	0	3	0.04	100	1.6
0.0	0	3	0.04	10	0.1
0.15	0	3	0.04	10	1
0.17	0	3	0.04	10	2
0.17	0	3	0.04	10	3
0.0	0	3	0.04	10	10

Dense-Sift: Parameters vs Accuracy (default sift parameters)	
Accuracy (1.0)	grid.size
-	8
-	16
-	32

In this experiment, k-value in kmeans is set to 128 cluster and k value of k\_nn nearest neighbor is set to 8 as required. By changing the parameters of SIFT, I've obtained significant changes on the accuracy of our classifier.

Firstly, nfeatures determines the number of top features in terms of local contrast scores to keep in the sift descriptor. Default nfeatures value takes all of the features in the descriptor. When nfeatures is set very low like 1, I've observed that accuracy drops drastically which makes sense because we are dropping discriminative sift vectors and this information loss has a cost on the overall accuracy. However, nfeatures=20 or 100, more or less yields the same accuracy which probably means that top 20 sift vectors were discriminative enough to make classification equivalent to the default param that is used in sift.

Secondly, nOctaveLayers are number of difference of gaussian layers which helps capturing features at varying scales. If we have more nOctaveLayers then that would be more features at varying scales but with exponentially (with order of  $k^2$ ) more blurring. This creates new features at very high scale due to higher blurring as a result of increased nOctaveLayers value. As far as I understand, in sift they have probably chosen the elbow value which should be somewhere around 3, because after/smaller than 3 it starts to decline very sharply. In addition, having less nOctaveLayers might be disadvantageous too

depending on the problem and image properties, hence I think 3 would be a better or less risky choice when compared to 1,2 which results in less scale invariant descriptors.

Thirdly, `contrastThreshold` parameter filters weak feature below a threshold value measured in contrast. When it's set higher, we would get less features and when it's set lower, we would get more features. However, more features don't always imply more accuracy, some redundant features can hurt our performance in classification as we'll see. I've observed that increasing `contrastThreshold`, results in less sift descriptors. For instance, `contrastThreshold=1` results no sift vector, because it's too high and `contrastThreshold=0.001` results in most number of sift descriptors. `ContrastThreshold` is a hyperparameter that we should tune for our problem and in our case 0.04 has resulted in the best accuracy, because it has discarded less important features and emphasized on more important features.

Fourthly, `edgeThreshold` is a hyperparameter that basically filters edge like features, its meaning is slightly different than `contrastThreshold` in the sense that, higher `edgeThreshold` value results in more features, because it basically eliminates features above that threshold not below. So, higher threshold results in more features, but once again that doesn't imply more accuracy, we still need to tune our param to give priority to important edges. As it's clear from our experiments, low `edgeThreshold` values yields less features, `edgeThreshold=1` is an extreme case which results in no sift descriptors. However, after some point accuracy doesn't improve as threshold gets larger.

Finally, `sigma` parameter is the sigma of the gaussians applied at each octave layers in sift. Higher sigma implies more blurring, which would encourage sift to capture features at higher scales and lower sigma would encourage sift to capture features at smaller scale. As sigma increases, number of sift descriptors decrease, because blurring kills some details. Also, low sigma values have low accuracies, because we capture features at lower scales more frequently. That's why we need to form a balance for the problem by tuning the hyperparameter.

## 2 Bag of Features (45 pts)

- How did you implement BoF? Briefly explain.
- Give pseudo-code for obtaining the dictionary.
- Give pseudo-code for obtaining BoF representation of an image once the dictionary is formed.
- Put your quantitative results (classification accuracy) regarding 3 different parameter configurations for the BoF pipeline here. Discuss possible reasons for each one's relatively better/worse accuracy. You are suggested to keep  $k \leq 1024$  in k-means to keep experiment durations manageable. You need to use the best feature extractor you obtained in the previous part together with the same classifier.

## 3 Classification (30 pts)

- Put your quantitative results regarding k-Nearest Neighbor Classifier for k values 16, 32 and 64 by using the best k-means representation and feature extractor. Discuss the effect of these briefly.
- What is the accuracy values, and how do you evaluate it? Briefly explain.
- Give confusion matrices for classification results of these combinations.

## 4 Additional Comments and References

(if there any)