



# Master in Computer Vision *Barcelona*

**Module 3:** Machine learning for computer vision

**Project:** Bag of Visual Words Image Classification

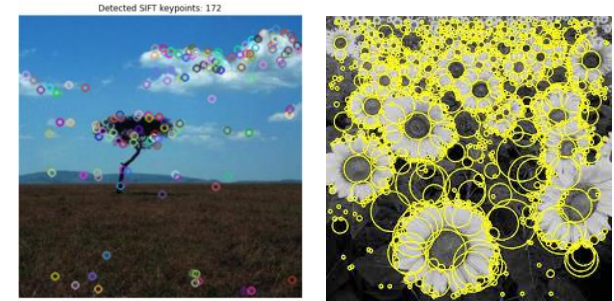
**Lecturer:** Ramon Baldrich, [ramon.baldrich@uab.cat](mailto:ramon.baldrich@uab.cat)

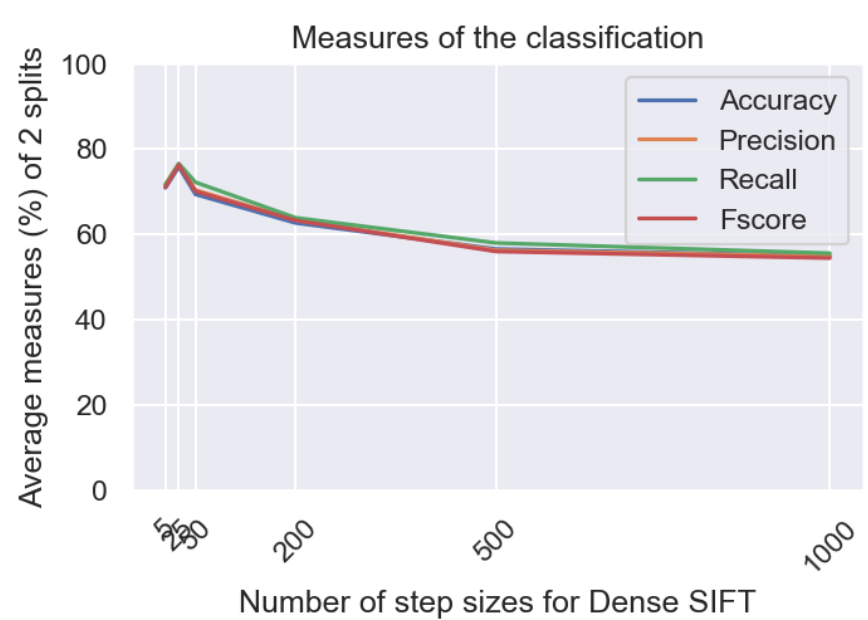
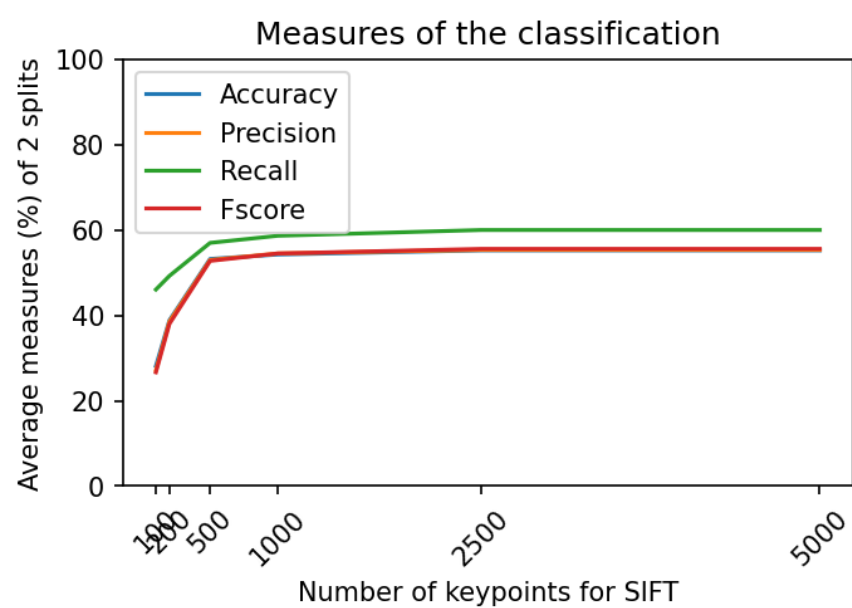
Credit to Marçal Rossinyol



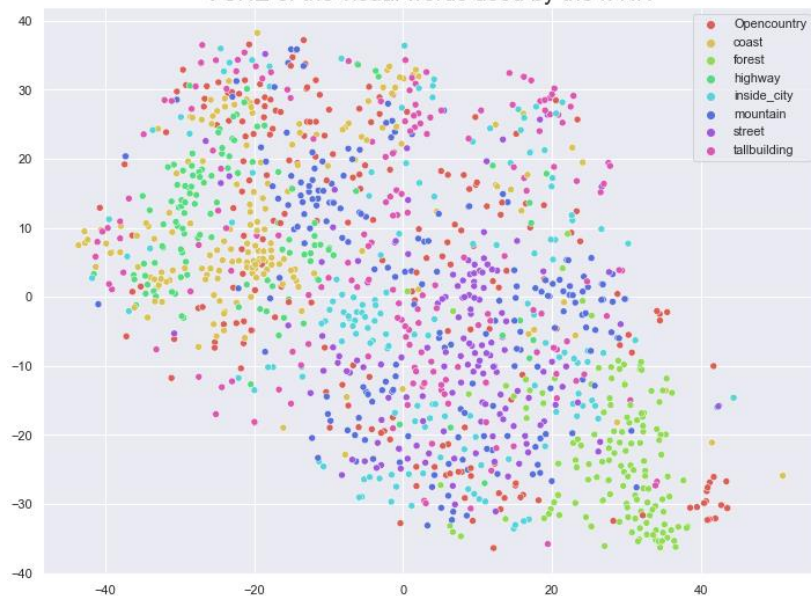
# S01 discussion

- Number of keypoints
  - The more the better
- Dense SIFT
  - Nearly nobody tried the role of scales!
- Codebook sizes / k-nn value
- k-nn and distances
  - Just slight differences found between point-wise distances
  - Which distance would work better for HISTOGRAMS?
- Dimensionality reduction not much effect
- Precompute stuff, store to disk!

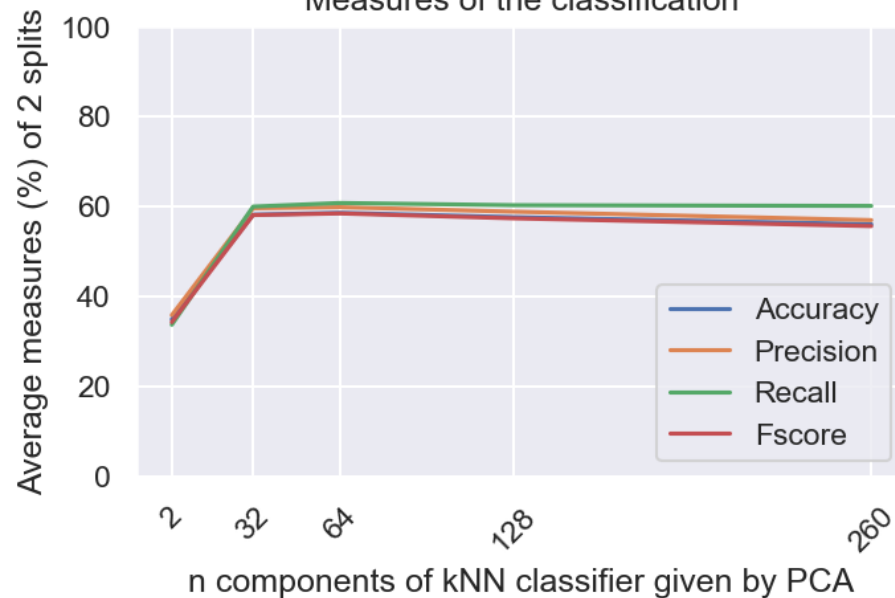




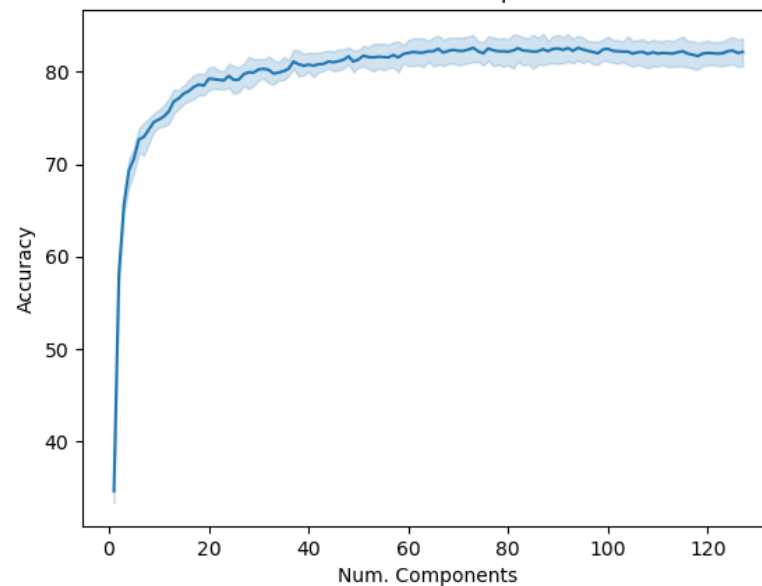
t-SNE of the visual words used by the k-NN



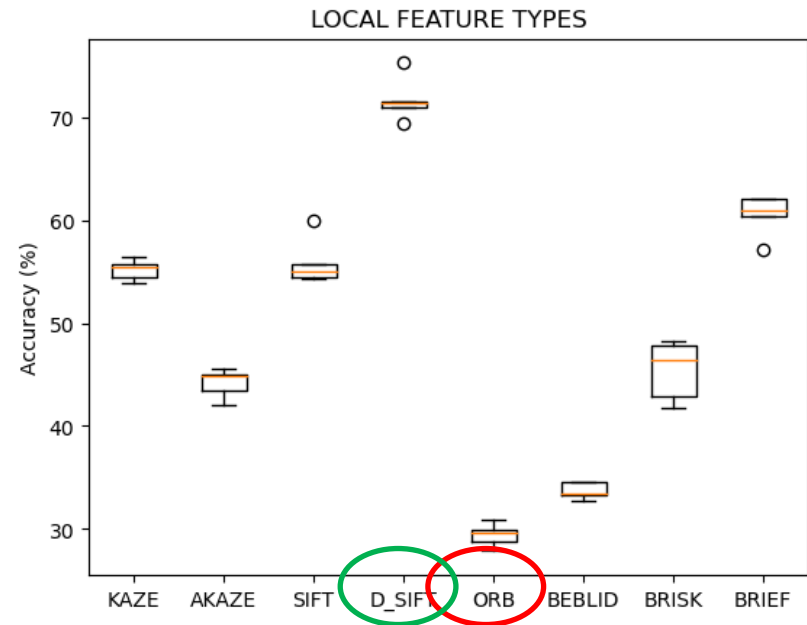
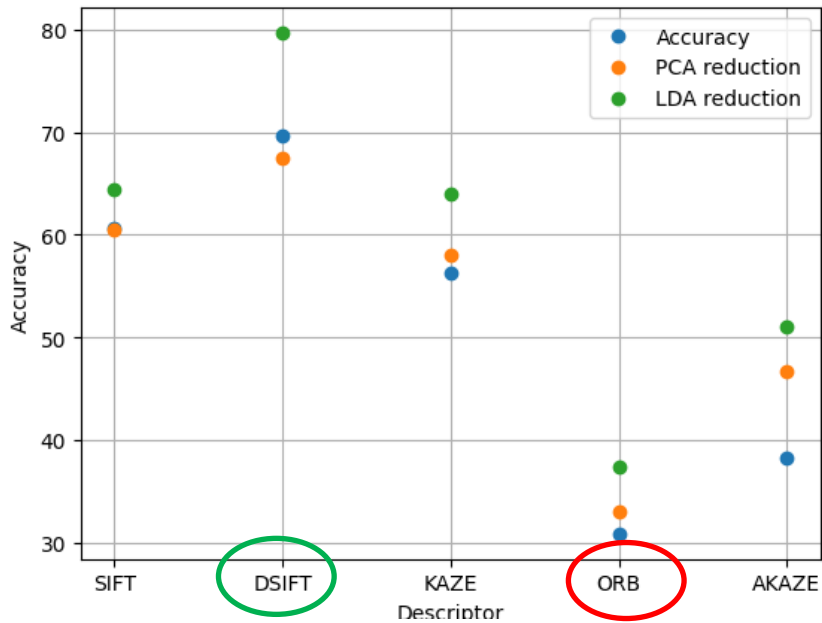
Measures of the classification



PCA: number of components

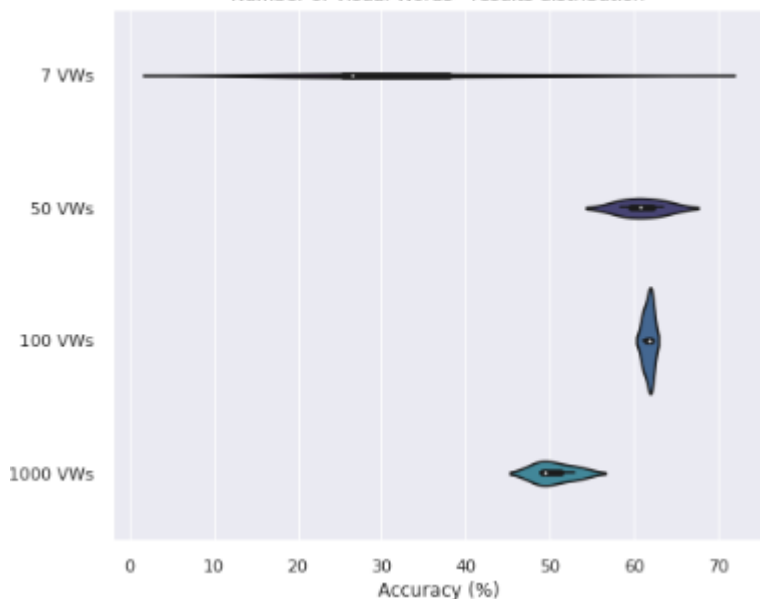


# The role of the descriptor

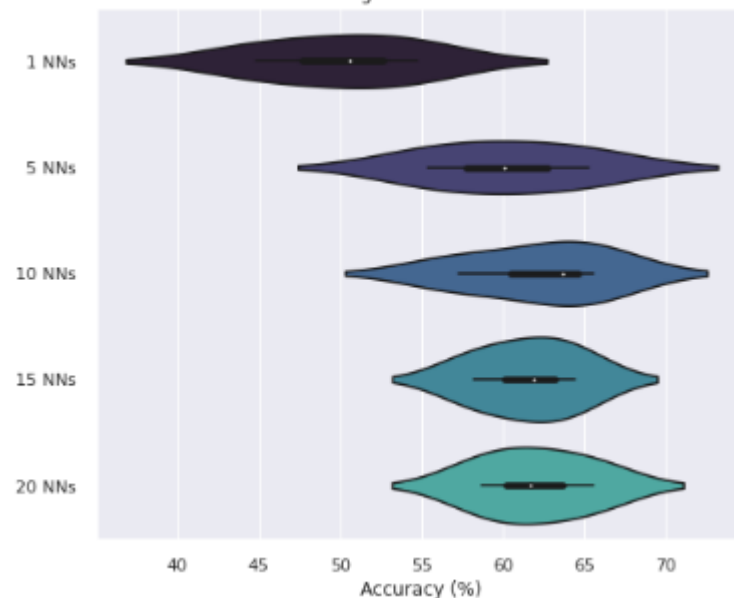


# Understanding parameter effects on performance

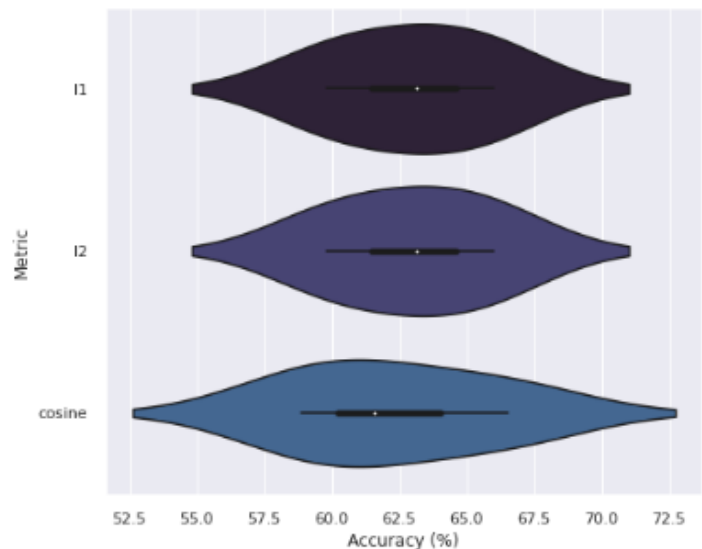
Number of Visual Words - results distribution



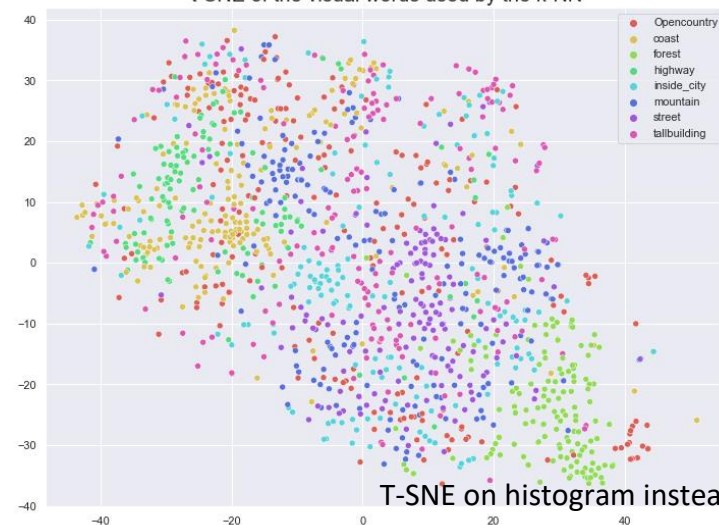
Number of Neighbors - results distribution



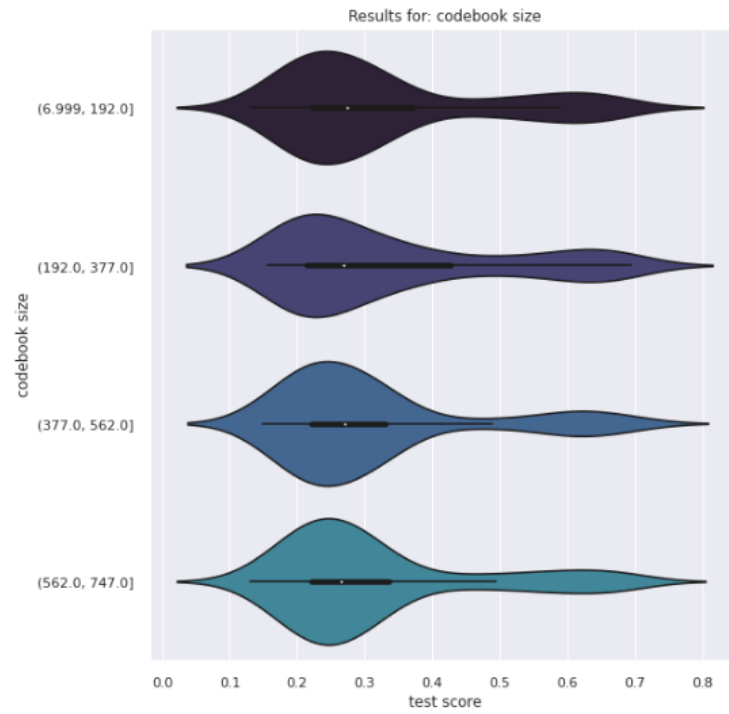
Distance metric - results distribution



t-SNE of the visual words used by the k-NN



T-SNE on histogram instead?



Don't compare parameters where other parameters change



# Visualize results (intermidate and final)

Original\_image

SIFT-Keypoints:151

BRISK-Keypoints:478

AKAZE-Keypoints:115

SURF-Keypoints:252



Opencountry

coast

forest

highway

inside\_city

mountain

street

tallbuilding

ClassAccuracy: 55.93  
GT: Opencountry,  
Pred: coast

ClassAccuracy: 57.76  
GT: coast,  
Pred: highway

ClassAccuracy: 87.13  
GT: forest,  
Pred: tallbuilding

ClassAccuracy: 67.11  
GT: highway,  
Pred: coast

ClassAccuracy: 64.89  
GT: inside\_city,  
Pred: tallbuilding

ClassAccuracy: 62.28  
GT: mountain,  
Pred: Opencountry

ClassAccuracy: 75.00  
GT: street,  
Pred: highway

ClassAccuracy: 62.04  
GT: tallbuilding,  
Pred: Opencountry



GT: Opencountry,  
Pred: coast

GT: coast,  
Pred: inside\_city

GT: forest,  
Pred: coast

GT: highway,  
Pred: coast

GT: inside\_city,  
Pred: tallbuilding

GT: mountain,  
Pred: coast

GT: street,  
Pred: highway

GT: tallbuilding,  
Pred: mountain



GT: Opencountry,  
Pred: coast

GT: coast,  
Pred: forest

GT: forest,  
Pred: tallbuilding

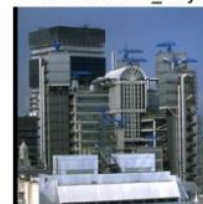
GT: highway,  
Pred: Opencountry

GT: inside\_city,  
Pred: tallbuilding

GT: mountain,  
Pred: Opencountry

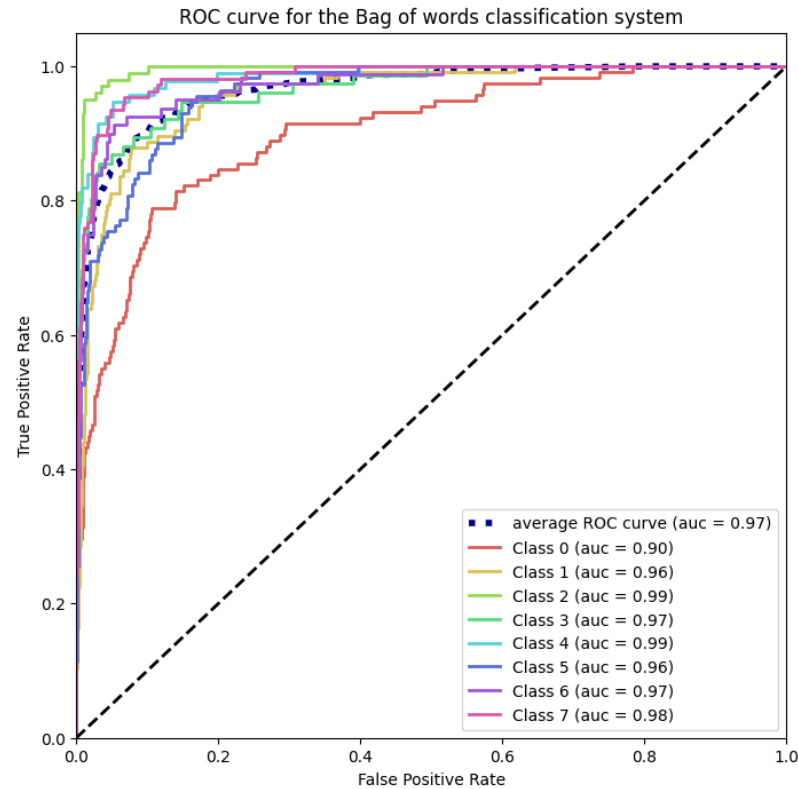
GT: street,  
Pred: mountain

GT: tallbuilding,  
Pred: inside\_city



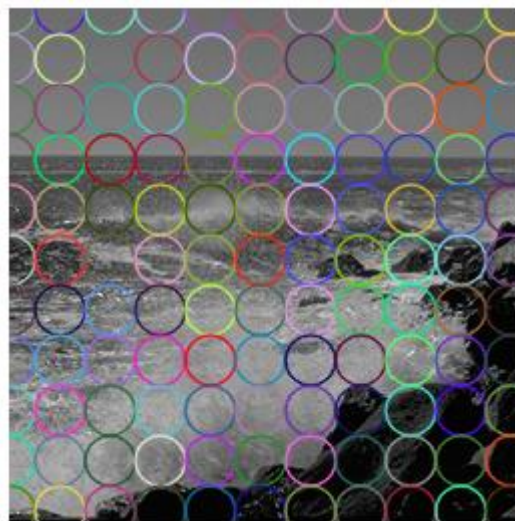
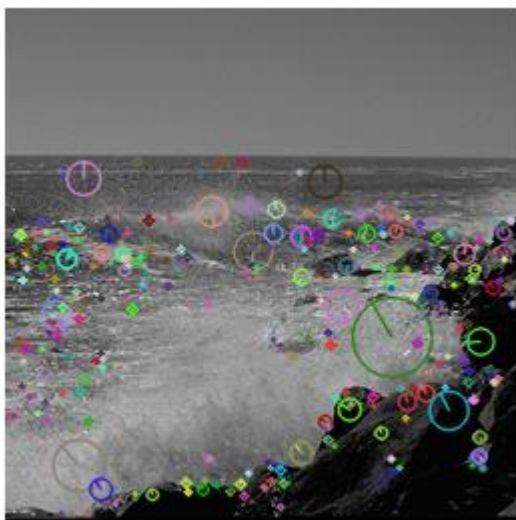


# Final performance



or precision/recall curve depending on the task

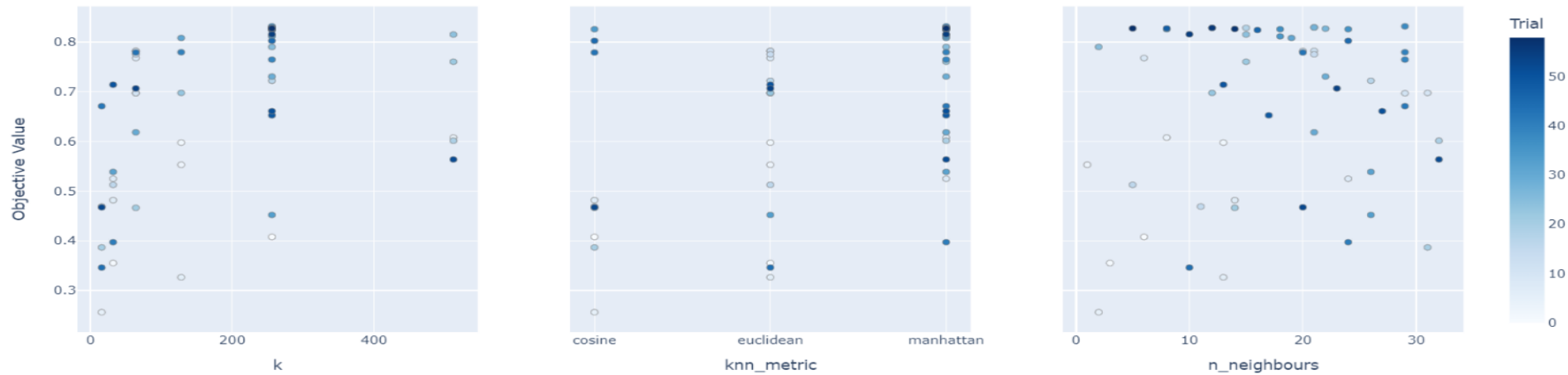
# Does scale play any role?



# Optuna

```
1 optuna.visualization.plot_slice(study, params=["k", "n_neighbours", "knn_metric"])
```

Slice Plot



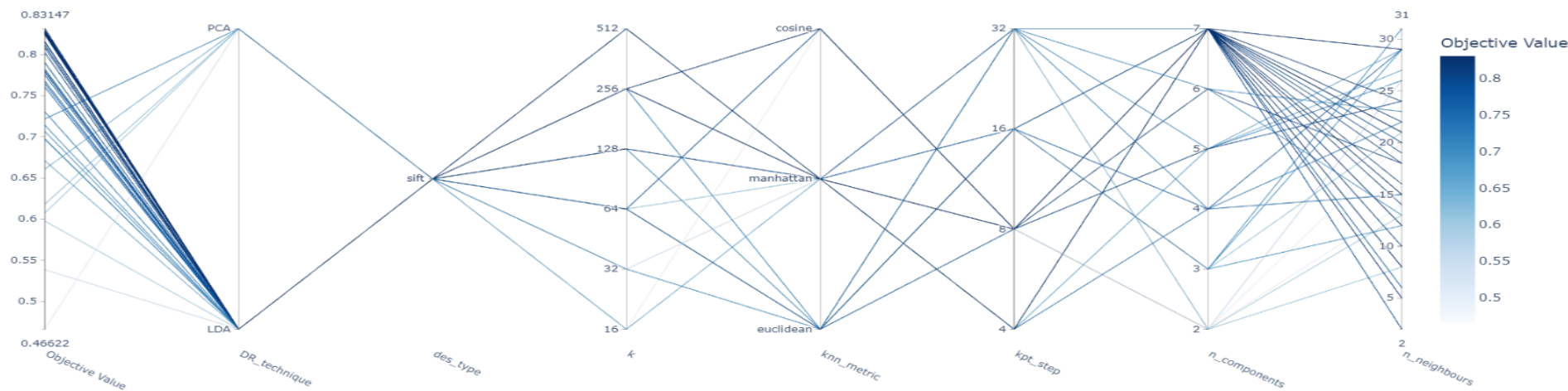
```

fig = optuna.visualization.plot_parallel_coordinate(sift_study, target=lambda t: t.values[0])
fig.data[0].line.reversescale = not fig.data[0].line.reversescale
fig.show()

```

2023-01-08 22:16:13,323] A new study created in memory with name: no-name-d00c3c6e-852e-41d8-90b6-a510fc2f6b9a

Parallel Coordinate Plot



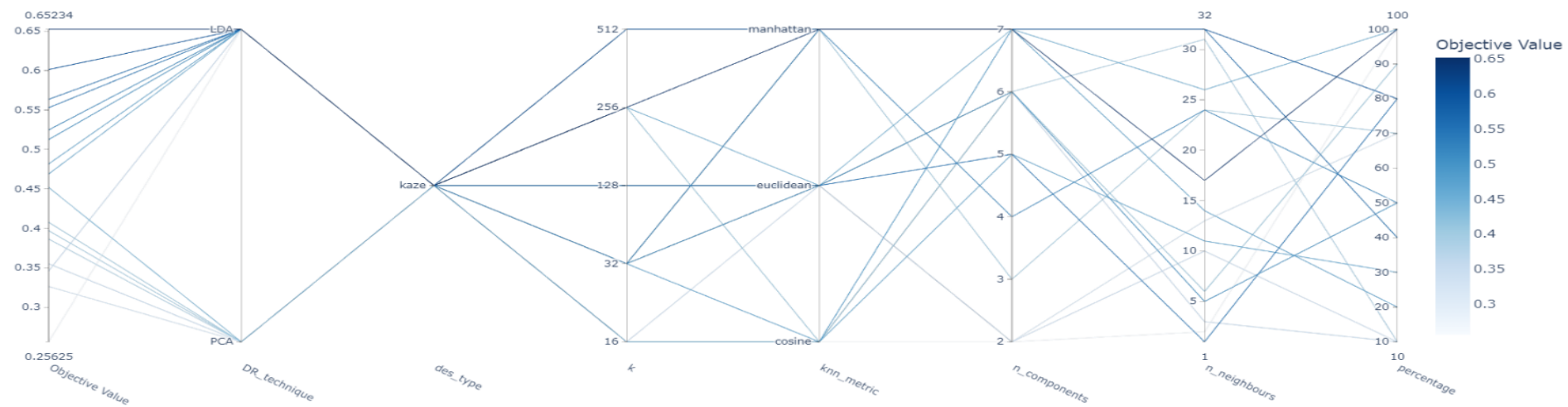
```

7 fig = optuna.visualization.plot_parallel_coordinate(kaze_study, target=lambda t: t.values[0])
8 fig.data[0].line.reversescale = not fig.data[0].line.reversescale
9 fig.show()

```

2023-01-08 22:16:09,923] A new study created in memory with name: no-name-271a6e74-960d-4e71-a91a-0a701868f52a

Parallel Coordinate Plot



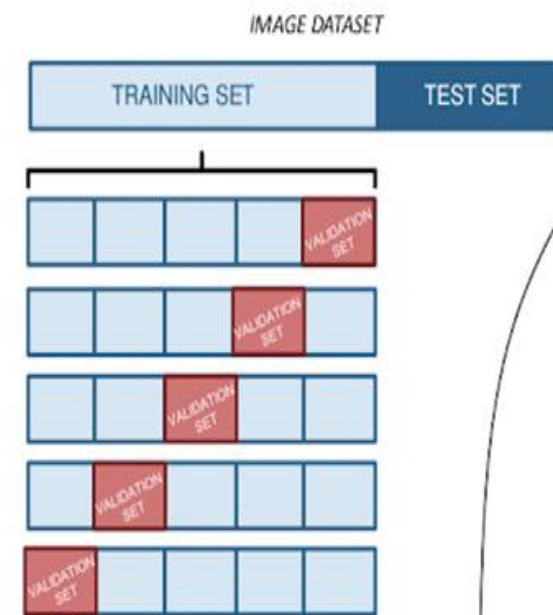
Group	Temptative grade
1	8
2	7
3	8
4	10
5	10
6	7
7	9
8	8
9	9
10	10



# S02

- We'll start with BoVW computed with Dense SIFT with a large enough codebook size
- We'll normalize descriptors
  - L2-norm, Power-norm, etc..
- Cross-validation
  - Sklearn functions: StratifiedkFold, GridsearchCV, Optuna
- Spatial Pyramids
- SVM and kernels
  - Use sklearn standardScaler to project every dimension to [0, 1]!
  - linear kernel
  - RBF kernel
  - our own histogram intersection kernel
- OPTIONAL: Fisher Vectors ([http://yael.gforge.inria.fr/tutorial/tuto\\_imgindexing.html](http://yael.gforge.inria.fr/tutorial/tuto_imgindexing.html))

# Cross Validation

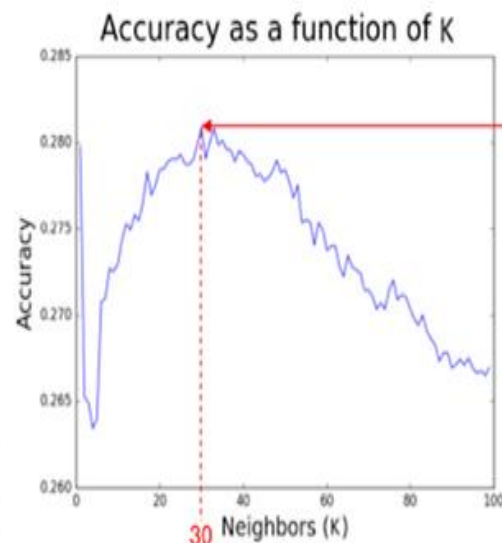


We use K-fold cross validation to pick suitable parameters for the classifier while **avoiding overfitting**.



For instance, in the KNN classifier we want to find the most suitable  $K$  (number of neighbours used to classify each descriptor)

We plotted the average of the accuracies (Y axis) obtained with 5 folds with respect to varying  $K$  from 1 to 100 (X axis).



Best average accuracy with  $K = 30$  !

So we obtain that individual SIFT descriptors are best classified using  $K = 30$  (for another kind of descriptor the 'best'  $K$  will vary)

# Hyperparameter optimization

Journal of Machine Learning Research 13 (2012) 281-305

Submitted 3/11; Revised 9/11; Published 2/12

## Random Search for Hyper-Parameter Optimization

James Bergstra  
Yoshua Bengio

Département d'Informatique et de recherche opérationnelle  
Université de Montréal  
Montréal, QC, H3C 3J7, Canada

JAMES.BERGSTRA@UMONTREAL.CA  
YOSHUA.BENGIO@UMONTREAL.CA

Editor: Leon Bottou

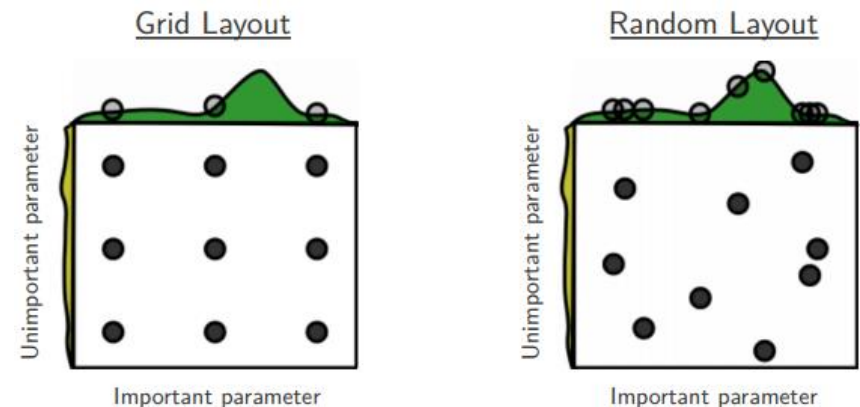
### Abstract

Grid search and manual search are the most widely used strategies for hyper-parameter optimization. This paper shows empirically and theoretically that randomly chosen trials are more efficient for hyper-parameter optimization than trials on a grid. Empirical evidence comes from a comparison with a large previous study that used grid search and manual search to configure neural networks and deep belief networks. Compared with neural networks configured by a pure grid search, we find that random search over the same domain is able to find models that are as good or better within a small fraction of the computation time. Granting random search the same computational budget, random search finds better models by effectively searching a larger, less promising configuration space. Compared with deep belief networks configured by a thoughtful combination of manual search and grid search, purely random search over the same 32-dimensional configuration space found statistically equal performance on four of seven data sets, and superior performance on one of seven. A Gaussian process analysis of the function from hyper-parameters to validation set performance reveals that for most data sets only a few of the hyper-parameters really matter, but that different hyper-parameters are important on different data sets. This phenomenon makes grid search a poor choice for configuring algorithms for new data sets. Our analysis casts some light on why recent "High Throughput" methods achieve surprising success—they appear to search through a large number of hyper-parameters because most hyper-parameters do not matter much. We anticipate that growing interest in large hierarchical models will place an increasing burden on techniques for hyper-parameter optimization; this work shows that random search is a natural baseline against which to judge progress in the development of adaptive (sequential) hyper-parameter optimization algorithms.

**Keywords:** global optimization, model selection, neural networks, deep learning, response surface modeling

### 1. Introduction

The ultimate objective of a typical learning algorithm  $\mathcal{A}$  is to find a function  $f$  that minimizes some expected loss  $\mathcal{L}(x; f)$  over i.i.d. samples  $x$  from a natural (grand truth) distribution  $\mathcal{G}_x$ . A learning algorithm  $\mathcal{A}$  is a functional that maps a data set  $\mathcal{X}^{(\text{train})}$  (a finite set of samples from  $\mathcal{G}_x$ ) to a function



Continuous hyperparameter: distribution over possible values

generate random variable

Discrete hyperparameter: list of discrete choices  
random selection  
(without replacement if all discrete)

Set the number of trials

Bergstra, James, and Yoshua Bengio. "Random search for hyper-parameter optimization." *Journal of Machine Learning Research* 13.Feb (2012): 281-305.

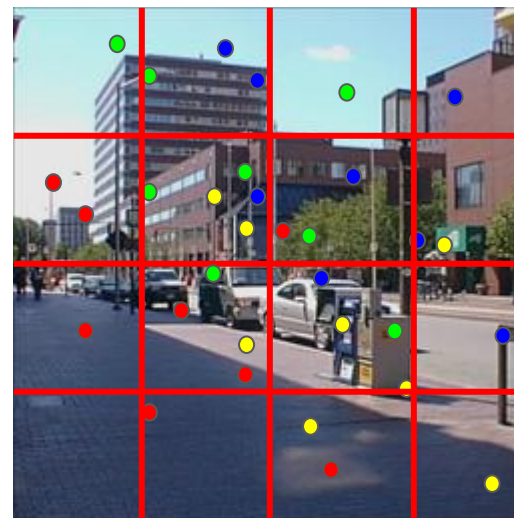
```
10 params = {
11     "C": np.arange(0.001, 1, 0.01),
12     "max_iter": np.arange(50, 550, 50)
13 }
14
15 lg = LogisticRegression(solver='liblinear')
16
17 lg_grid = GridSearchCV(lg, params, cv=8,
18     scoring=["accuracy", "f1_macro"], refit="accuracy", return_train_score=True)
19 lg_grid.fit(visual_words, train_labels)
```

# Spatial Pyramids





# Spatial Pyramids



# Histogram Intersection kernel

```
def histogramIntersection(M, N):
```

$$K_{int}(A, B) = \sum_{i=1}^m \min\{a_i, b_i\}.$$

```
    return K_int
```

# Fisher Vectors

- BoVW is only counting the number of local descriptors assigned to each Voronoi cell
- Why not including higher order statistics?
  - Mean of local descriptors
  - Co-variance of local descriptors
- FV is typically  $2 \times D \times k$  dimensional

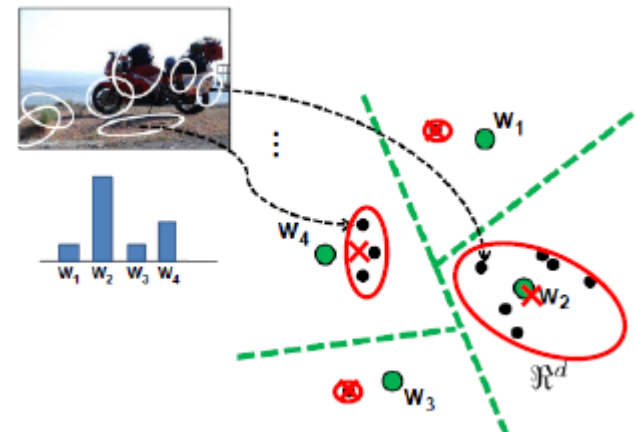
## Fisher Kernels on Visual Vocabularies for Image Categorization

Florent Perronnin and Christopher Dance  
Xerox Research Centre Europe  
6 chemin de Maupertuis, 38240 Meylan, France  
Firstname.Lastname@xrce.xerox.com  
<http://www.xrce.xerox.com>

### Abstract

Within the field of pattern classification, the Fisher kernel is a powerful framework which combines the strengths of generative and discriminative approaches. The idea is to characterize a signal with a gradient vector derived from a generative probability model and to subsequently feed this representation to a discriminative classifier. We propose to apply this framework to image categorization where the in-

the bag-of-keypatches [2] or bag-of-visual-words (BOV) [11]. In the following, we use the latter denomination which is more general (the term keypatches assumes the use of an interest point detector for the extraction of low-level feature vectors). Given a visual vocabulary, the idea is to characterize an image with the number of occurrences of each visual word. Any classifier can then be used for the categorization of this histogram representation. Most of the work on bags-of-visual-words has focused on the estimation of the visual vocabulary. This is done through the extraction of low-level



Slide credit F. Perronnin. Features for Large-Scale Visual Recognition

# Tasks to do

Improve the BoVW code with:

- Dense SIFT (with tiny steps and different scales!)
- L2-norm - power norm
- SVM classifier
- StandardScaler
- Cross-validation
- Linear, RBF and histogram intersection kernels
- Spatial Pyramids
- Fisher Vectors (OPTIONAL)

# Deliverable

- A **single Python notebook file per group** reporting all the work done,
  - with the different experiments,
  - code,
  - plots,
  - explanations, etc.
  - **EVERYTHING EXECUTED!**
- To deliver by Monday 17th @ 10 A.M.
  - Please, state clearly your group.

**Warning:** provided code might not work out of the box depending on the used versions (OpenCV, numpy, sklearn...) do not panic, and read the documentation