

THE UNIVERSITY OF BRITISH COLUMBIA

CPSC 425: Computer Vision



Lecture 29: Image Classification

Menu for Today (November 16, 2018)

Topics:

- Scene Classification
- Bag of Words Representation

Redings:

- Today's Lecture: Forsyth & Ponce (2nd ed.) 16.1.3, 16.1.4, 16.1.9
- Next Lecture:

Reminders:



Decision Tree Boosting

Forsyth & Ponce (2nd ed.) 17.1–17.2

Assignment 5: Scene Recognition with Bag of Words due last day of classes





Today's "fun" Example: CVPR Deadline is Today

Lecture 28: Re-cap

Classifiers take as input a set of features and output (predict) a class label

Classifiers need to take into account "**loss**" associated with each kind of classification error

A **receiver operating characteristic** (ROC) curve plots the trade-off between false negatives and false positives

Non-parametric classifiers, like k-nearest neighbour, are data driven. New data points are classified by comparing to the training examples directly.

Parametric classifiers, like support vector machines, are model driven. New data points are classified by evaluating a model learned from the training examples.

Lecture 28: Re-cap

majority vote.

various dimensions

k-Nearest Neighbor (kNN) Classifier



Given a new data point, find the k nearest training examples. Assign the label by

Simple method that works well if the distance measure correctly weights the

Figure credit: Hastie, Tibshirani & Friedman (2nd ed.)



Lecture 28: Re-cap

The decision boundary is parameterized as a **separating hyperplane** in feature space. — e.g. a separating line in 2D

We choose the hyperplane that is as far as possible from each class - that maximizes the distance to the closest point from either class.

Support Vector Machines (SVM)





.....

 $\Big)$

Image Classification

We next discuss **image classification**, where we pass a whole image into a classifier and obtain a class label as output.

What Makes Image Classification Hard?





Intra-class variation, viewpoint, illumination, clutter, and occlusion (among others!)



Figure source: Jianxiong Xiao. Original credit: ?

Image Classification

applied to classify natural scenes (e.g. beach, forest, harbour, library).

Why might classifying scenes be useful?

In addition to images containing single objects, the same techniques can be

Image Classification

applied to classify natural scenes (e.g. beach, forest, harbour, library).

Why might classifying scenes be useful?

Visual perception is influenced by expectation. Our expectations are often conditioned on the **context**.

- In addition to images containing single objects, the same techniques can be













Visual Words

Many algorithms for image classification accumulate evidence on the basis of **visual words**.

To classify a text document (e.g. as an article on sports, entertainment, business, politics) we might find patterns in the occurrences of certain words.

Vector Space Model

G. Salton. 'Mathematics and Information Retrieval' Journal of Documentation, 1979





working with collaborators for the robotic device to its at Harvard University, the achieve natural motions in beh University of Southern the ankle.

0

Tartan

http://www.fodey.com/generators/newspaper/snippet.asp

California, MIT and



1	6	2	1	0	0	0	1
Tartan	robot	CHIMP	CMU	bio	soft	ankle	sensor



Vector Space Model

A document (datapoint) is a vector of counts over each word (feature)

 $n(\cdot)$ counts the number of occurrences

What is the similarity between two documents?

 $\boldsymbol{v}_d = [n(w_{1,d}) \ n(w_{2,d}) \ \cdots \ n(w_{T,d})]$

just a histogram over words





Vector Space Model

A document (datapoint) is a vector of counts over each word (feature)

 $n(\cdot)$ counts the number of occurrences

What is the similarity between two documents?

Use any distance you want but the cosine distance is fast and well designed for high-dimensional vector spaces:

$$egin{aligned} d(oldsymbol{v}_i,oldsymbol{v}_j) &= \cos heta \ &= rac{oldsymbol{v}_i \cdot oldsymbol{v}_i}{\|oldsymbol{v}_i\|} \end{aligned}$$

 $oldsymbol{v}_d = [n(w_{1,d}) \quad n(w_{2,d}) \quad \cdots \quad n(w_{T,d})]$

just a histogram over words





 \boldsymbol{v}_{i} $oldsymbol{v}_i \| \| oldsymbol{v}_j \|$



Visual Words

In images, the equivalent of a word is a local image patch. The local image patch is described using a descriptor such as SIFT.

We construct a vocabulary or codebook of local descriptors, containing representative local descriptors.

What **Objects** do These Parts Belong To?































Some local feature are very informative

An object as





- deals well with occlusion
- scale invariant
- rotation invariant

(not so) Crazy Assumption



spatial information of local features can be ignored for object recognition (i.e., verification)

Recall: Texture Representation













Visual Words

patch is described using a descriptor such as SIFT.

We construct a **vocabulary** or **codebook** of local descriptors, containing representative local descriptors.

SIFT descriptors, say 1 million, how can we choose a small number of 'representative' SIFT codewords, say 1000?

- In images, the equivalent of a word is a local image patch. The local image

Question: How might we construct such a codebook? Given a large sample of

Standard **Bag-of-Words** Pipeline (for image classification)

Dictionary Learning: Learn Visual Words using clustering

Encode: build Bags-of-Words (BOW) vectors for each image

Classify: Train and test data using BOWs

1. Dictionary Learning: Learn Visual Words using Clustering

1. extract features (e.g., SIFT) from images











1. Dictionary Learning: Learn Visual Words using Clustering

2. Learn visual dictionary (e.g., K-means clustering)







What **Features** Should We Extract?

- Regular grid Vogel & Schiele, 2003 Fei-Fei & Perona, 2005
- Interest point detector Csurka et al. 2004 Fei-Fei & Perona, 2005 Sivic et al. 2005
- Other methods Random sampling (Vidal-Naquet & Ullman, 2002) Segmentation-based patches (Barnard et al. 2003)



Extracting SIFT Patches



Compute SIFT descriptor

Normalize patch

[Lowe'99]



Detect patches

[Mikojaczyk and Schmid '02] [Mata, Chum, Urban & Pajdla, '02] [Sivic & Zisserman, '03]

Extracting SIFT Patches







Creating **Dictionary**



Creating **Dictionary**





Creating **Dictionary**





K-means clustering

Lecture 26: Re-cap

K-means is a clustering technique that iterates between

- **1**. Assume the cluster centers are known. Assign each point to the closest cluster center.
- **2.** Assume the assignment of points to clusters is known. Compute the best cluster center for each cluster (as the mean).
- K-means clustering is initialization dependent and converges to a local minimum



Example Visual Dictionary



Example Visual Dictionary







Source: B. Leibe

Example Visual Dictionary





Source: B. Leibe

Standard **Bag-of-Words** Pipeline (for image classification)

Classify: Train and test data using BOWs

Dictionary Learning: Learn Visual Words using clustering

Encode: build Bags-of-Words (BOW) vectors for each image

2. Encode: build Bag-of-Words (BOW) vectors for each image



1. Quantization: image features gets associated to a visual word (nearest cluster center)













2. Encode: build Bag-of-Words (BOW) vectors for each image

2. Histogram: count the number of visual word occurrences







2. Encode: build Bag-of-Words (BOW) vectors for each image







frequency

codewords





Standard **Bag-of-Words** Pipeline (for image classification)

Classify: Train and test data using BOWs

Dictionary Learning: Learn Visual Words using clustering

Encode: build Bags-of-Words (BOW) vectors for each image

3. Classify: Train and text classifier using BOWs



K nearest neighbors



Bag-of-Words Representation

Algorithm:

Initialize an empty K -bin histogram, where K is the number of codewords Extract local descriptors (e.g. SIFT) from the image For each local descriptor **x**

Map (Quantize) **x** to its closest codeword \rightarrow **c**(**x**) Increment the histogram bin for c(x)Return histogram

We can then classify the histogram using a trained classifier, e.g. a support vector machine or k-Nearest Neighbor classifier

Please get your iClickers — Quiz

Spatial Pyramid

The bag of words representation does not preserve any spatial information The **spatial pyramid** is one way to incorporate spatial information into the image descriptor.

A spatial pyramid partitions the image and counts codewords within each grid box; this is performed at multiple levels

Spatial Pyramid



Fig. 16.8 in Forsyth & Ponce (2nd ed.). Original credit: Lazebnik et al., 2006

50

Please get your iClickers — Quiz

Summary

Factors that make image classification hard — intra-class variation, viewpoint, illumination, clutter, occlusion...

A codebook of **visual words** contains representative local patch descriptors — can be constructed by clustering local descriptors (e.g. SIFT) in training images

The **bag of words** model accumulates a histogram of occurrences of each visual word

The **spatial pyramid** partitions the image and counts visual words within each grid box; this is repeated at multiple levels