



HA NOI UNIVERSITY OF SCIENCE AND TECHNOLOGY
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Computer Vision

Chapter 7 (part 1): Object recognition

Contents

- Overview of 'semantic vision'?
- Image classification/ recognition
- Bag-of-words
 - Recall
 - Vocabulary tree
- Classification
 - K nearest neighbors
 - Naïve Bayes
 - [Support vector machine]



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Is this a street light?
(Recognition / classification)



3

Where are the people?
(Detection)



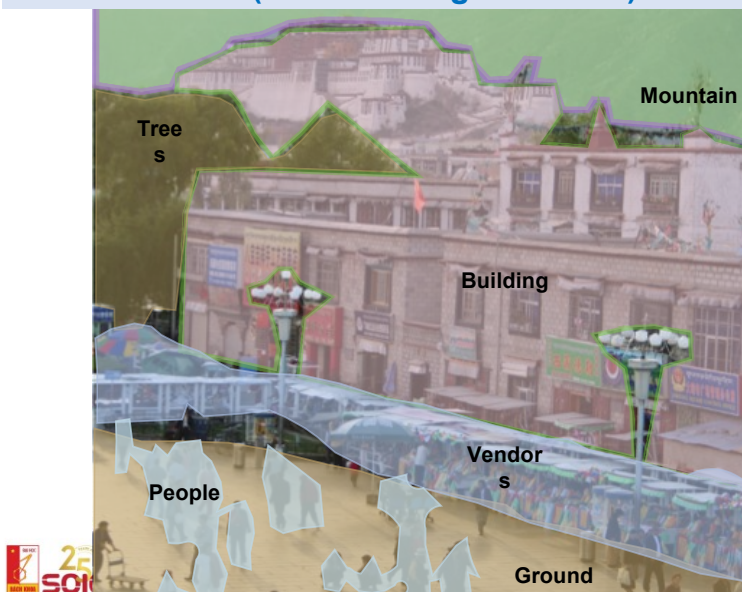
4

Is that Potala palace? (Identification)



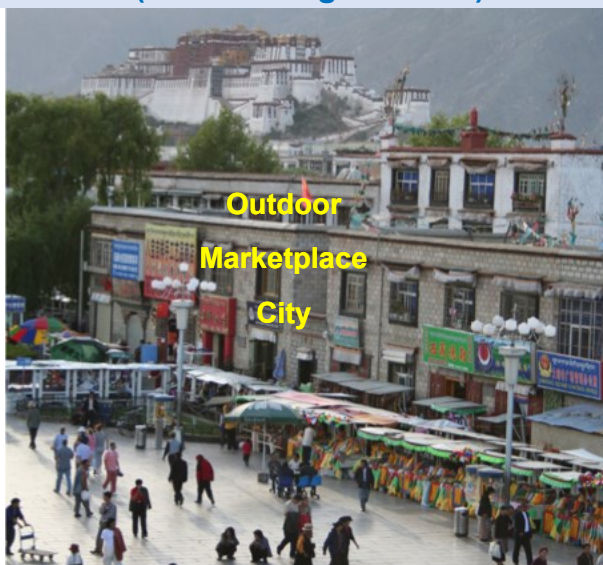
5

What's in the scene? (semantic segmentation)



6

What type of scene is it?
(Scene categorization)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

7

How many object categories are there?



Challenge: variable viewpoint

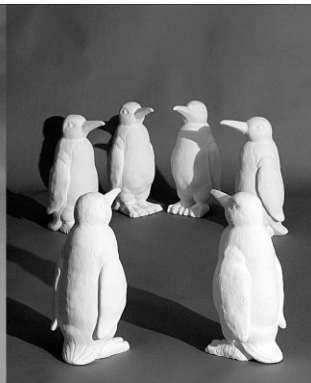
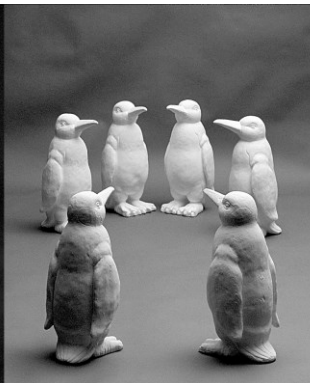


Michelangelo 1475-1564



9

Challenge: variable illumination



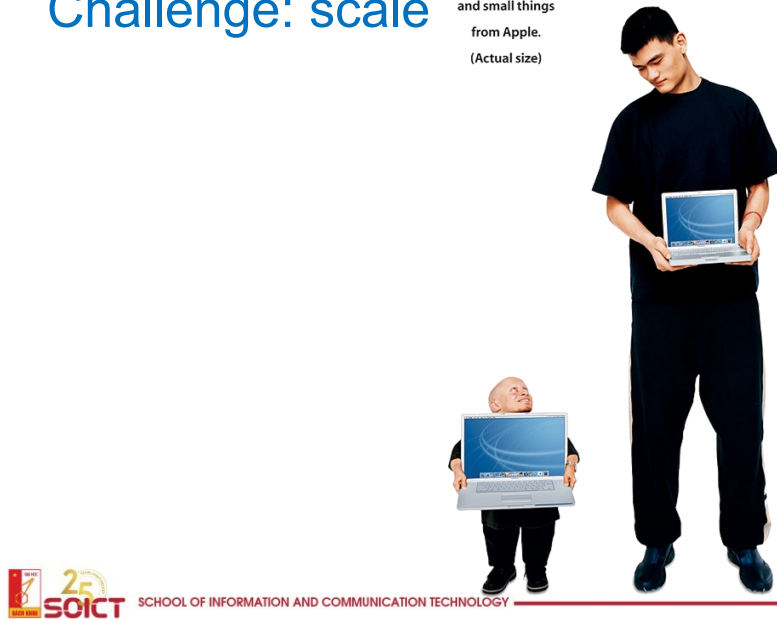
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

image credit: J. Koenderink

10

Challenge: scale

and small things
from Apple.
(Actual size)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

11

Challenge: deformation



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

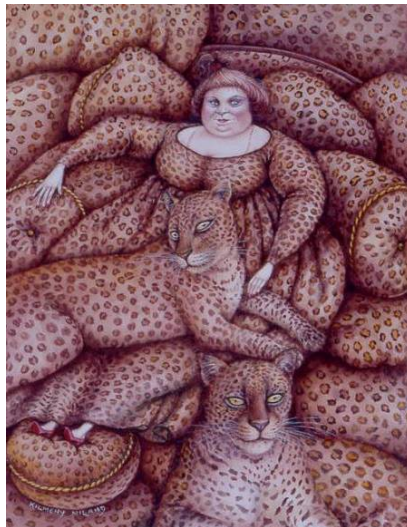
12

Challenge: Occlusion



Magritte, 1957

Challenge: background clutter



Kilmeny Niland. 1995



Challenge: intra-class variations



Image Classification/ Recognition



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

17

Image Classification/ Recognition



(assume given set of discrete labels)
{dog, cat, truck, plane, ...}



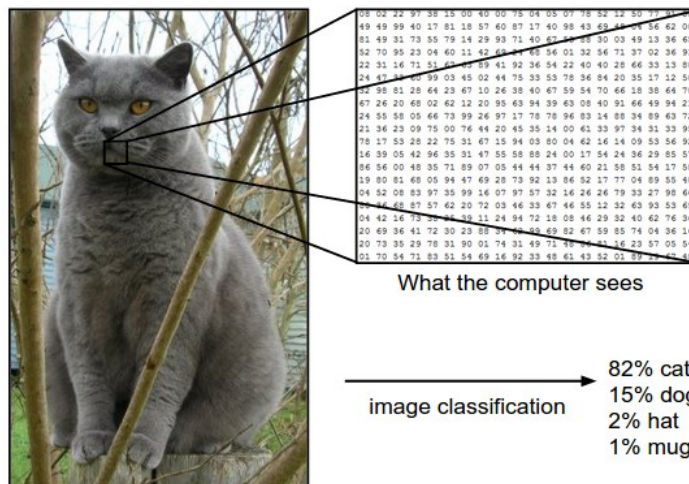
cat



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

18

Image Classification: Problem



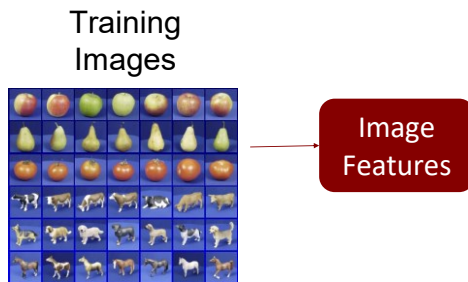
Data-driven approach

- Collect a database of images with labels
- Use ML to train an image classifier
- Evaluate the classifier on test images

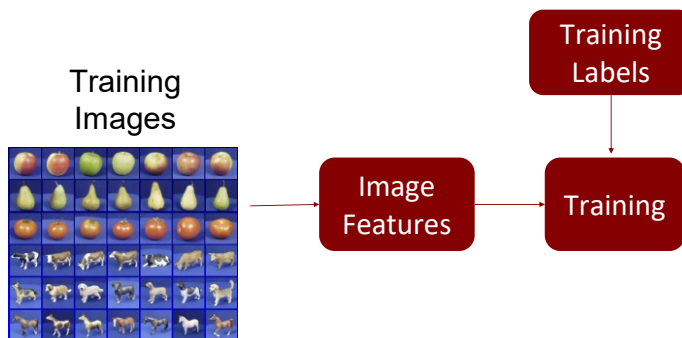
Example training set



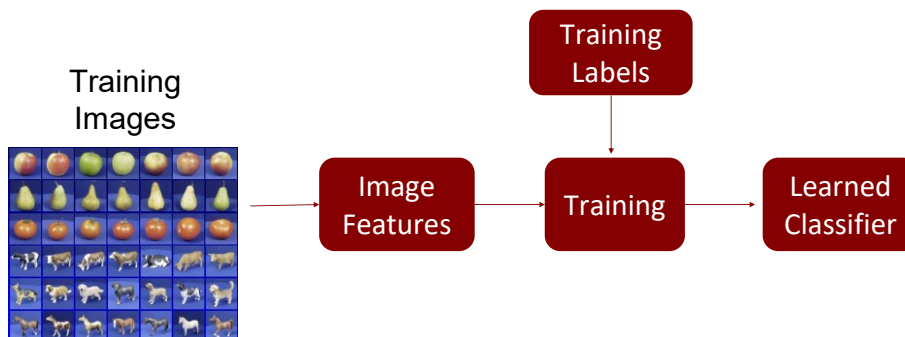
A simple pipeline - Training



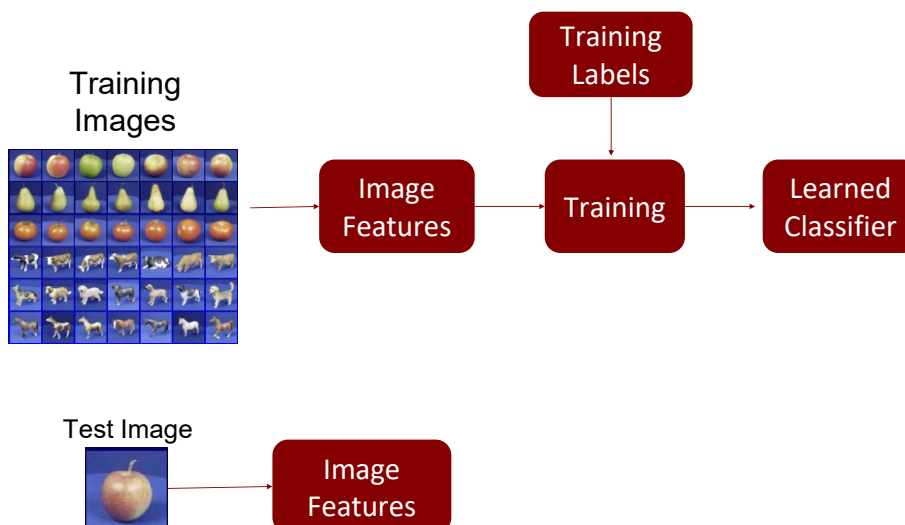
A simple pipeline - Training



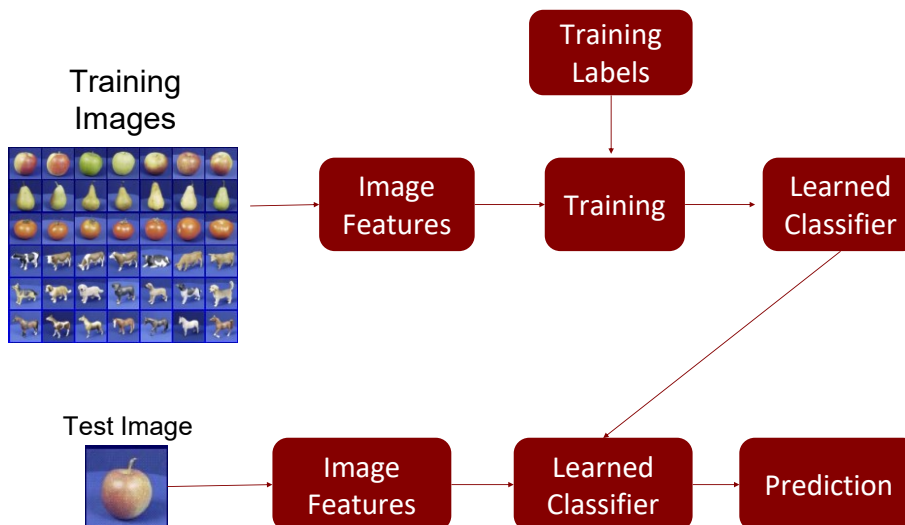
A simple pipeline - Training



A simple pipeline - Training



A simple pipeline - Training



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

25

Bag of words

Basic model

Vocabulary tree

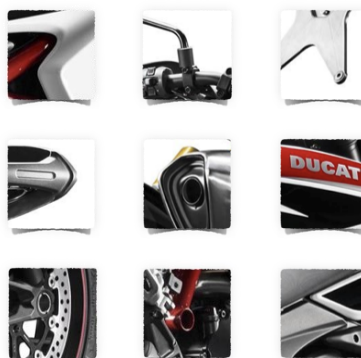


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

26

Some local feature are
very informative

An object as



a collection of local features
(bag-of-features)

- deals well with occlusion
- scale invariant
- rotation invariant

27

Bag-of-words

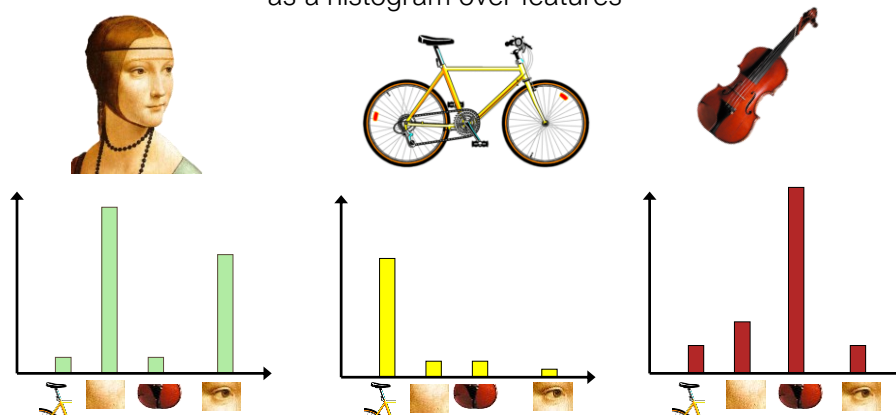
- Local feature ~~ a word
- An image ~~ a document
- Apply a technique for textual document representation:
vector model



28

Bag-of-words

represent a data item (document, texture, image)
as a histogram over features



Standard BOW 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



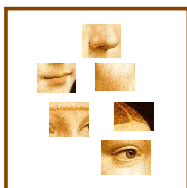
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

31

Dictionary Learning:

Learn Visual Words using clustering

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

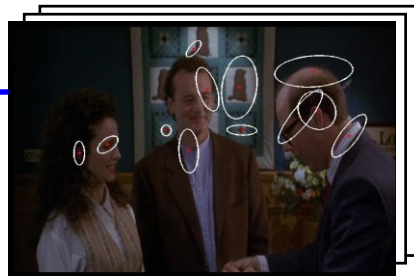
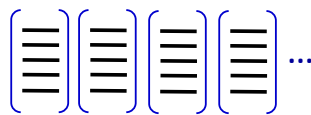
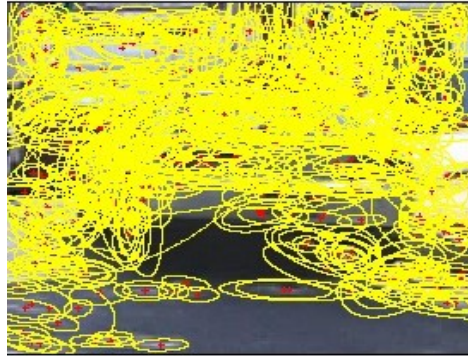


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

32

What kinds of features can 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)

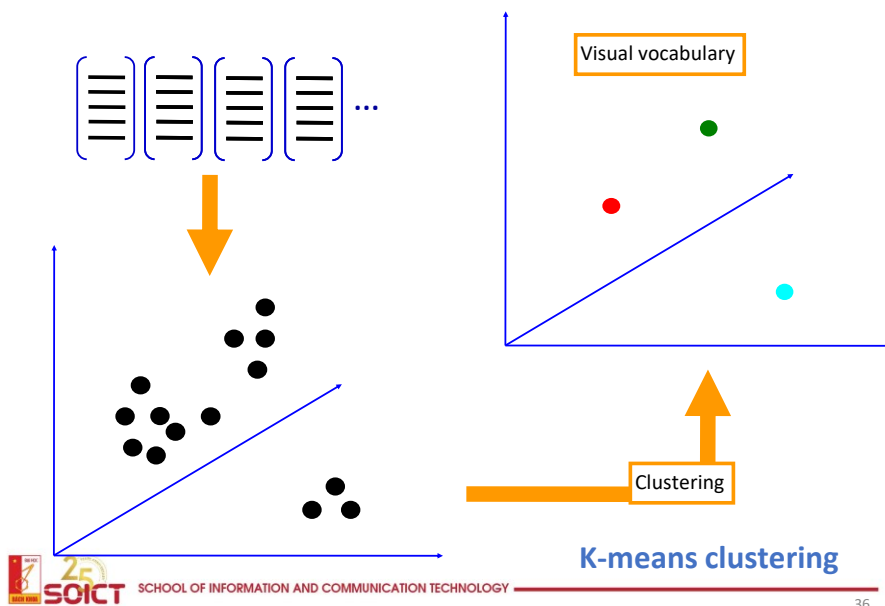


Dictionary Learning: Learn Visual Words using clustering

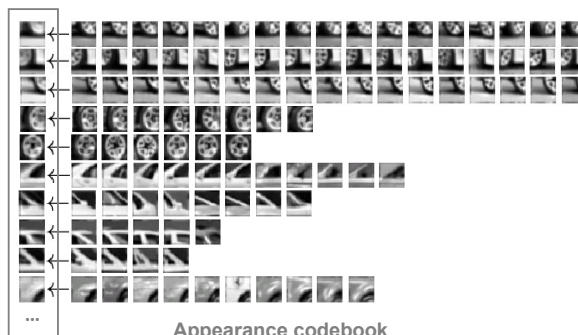
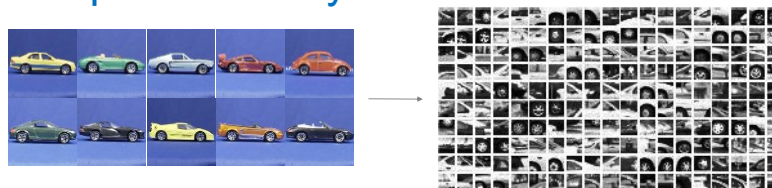
2. Learn visual dictionary (e.g., K-means clustering)



How do we learn the dictionary?



Example dictionary



Appearance codebook



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Source: B. Leibe

37

Dictionary Learning:

Learn Visual Words using clustering

Encode:

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

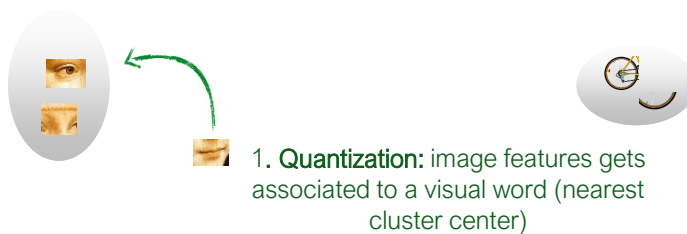
Classify:

Train and test data using BOWs



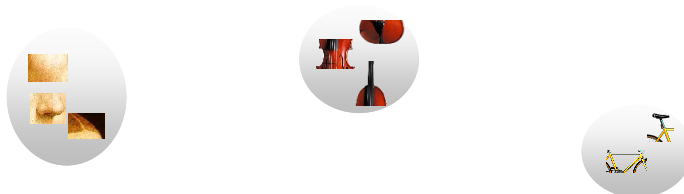
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

38



Encode:

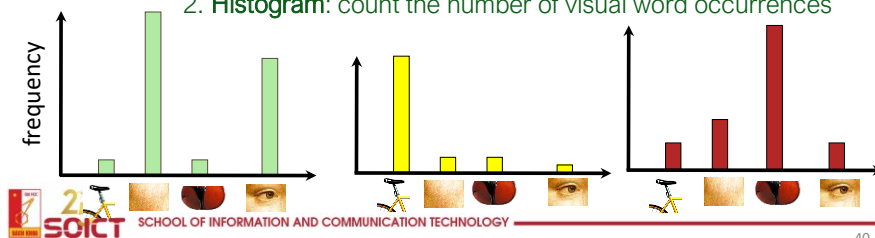
build Bags-of-Words (BOW) vectors for each image

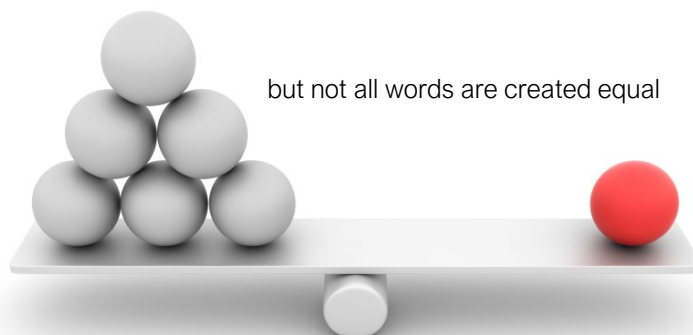


Encode:

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

2. **Histogram:** count the number of visual word occurrences





TF-IDF

Term **F**requency Inverse **D**ocument **F**requency

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

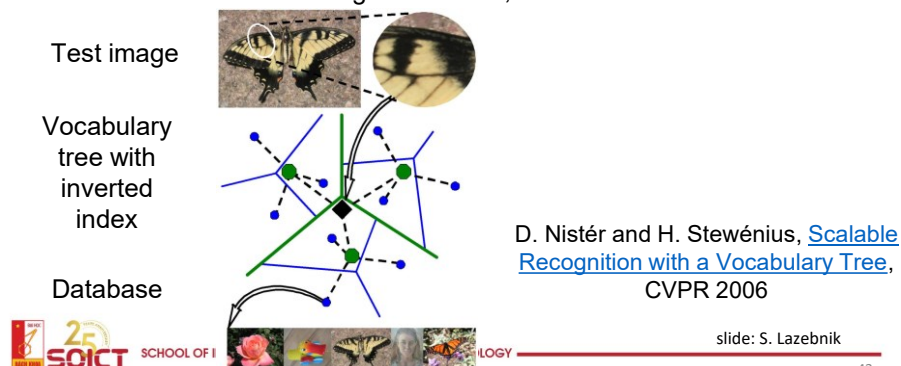
weight each word by a heuristic

$$\mathbf{v}_d = [n(w_{1,d})\alpha_1 \quad n(w_{2,d})\alpha_2 \quad \cdots \quad n(w_{T,d})\alpha_T]$$

$$n(w_{i,d})\alpha_i = n(w_{i,d}) \log \left\{ \frac{\text{inverse document frequency}}{\text{term frequency}} \right\}$$

Scalability: Alignment to large databases

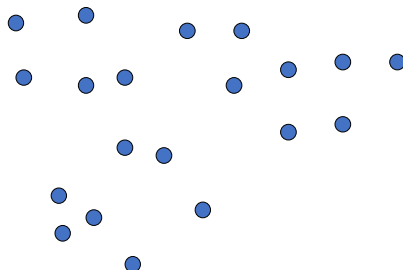
- What if we need to align a test image with thousands or millions of images in a model database?
 - Efficient putative match generation
 - Fast nearest neighbor search, inverted indexes



43

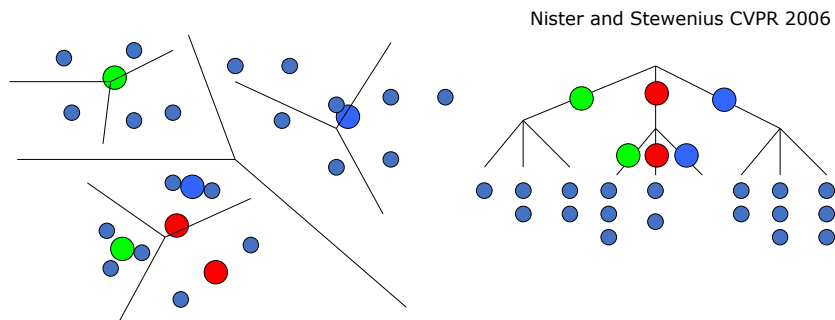
What is a Vocabulary Tree?

Nister and Stewenius CVPR 2006

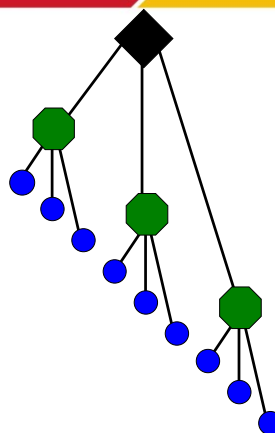


44

What is a Vocabulary Tree?

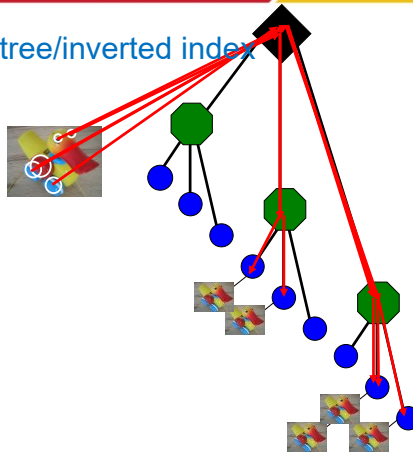


- Multiple rounds of K-Means to compute decision tree (offline)
- Fill and query tree online



Populating the vocabulary tree/inverted index

Model images

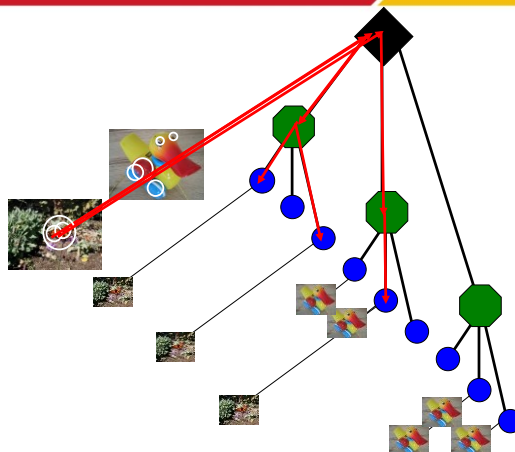


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Slide credit: D. Nister

47

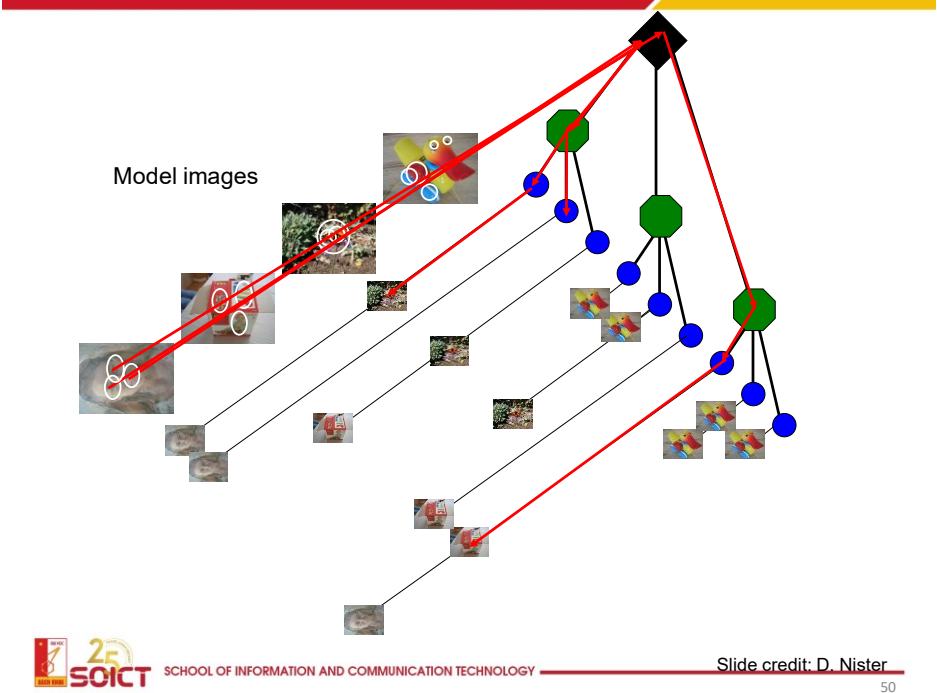
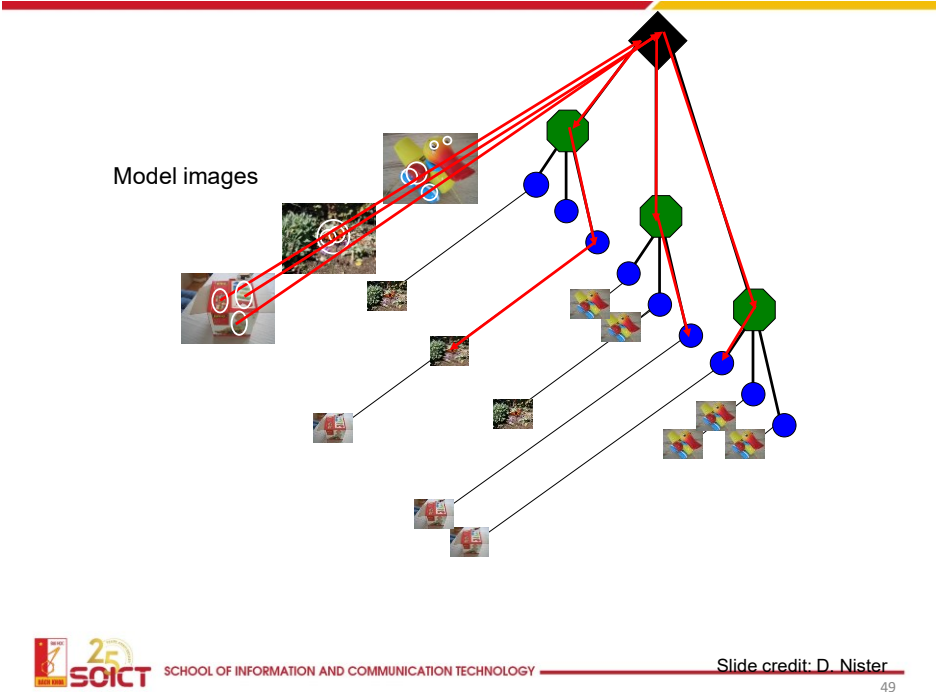
Model images

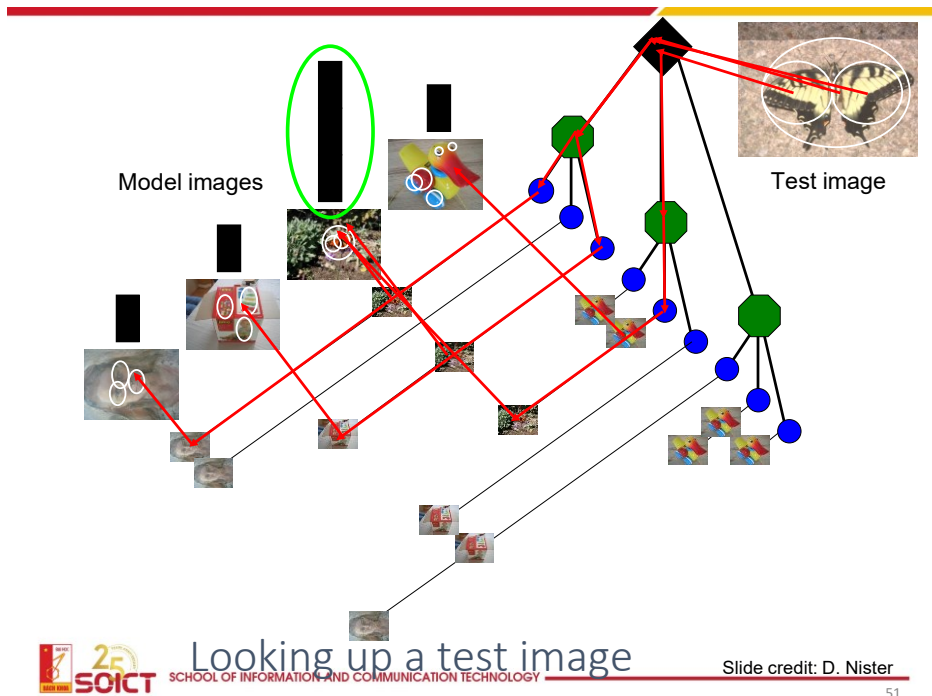


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Slide credit: D. Nister

48

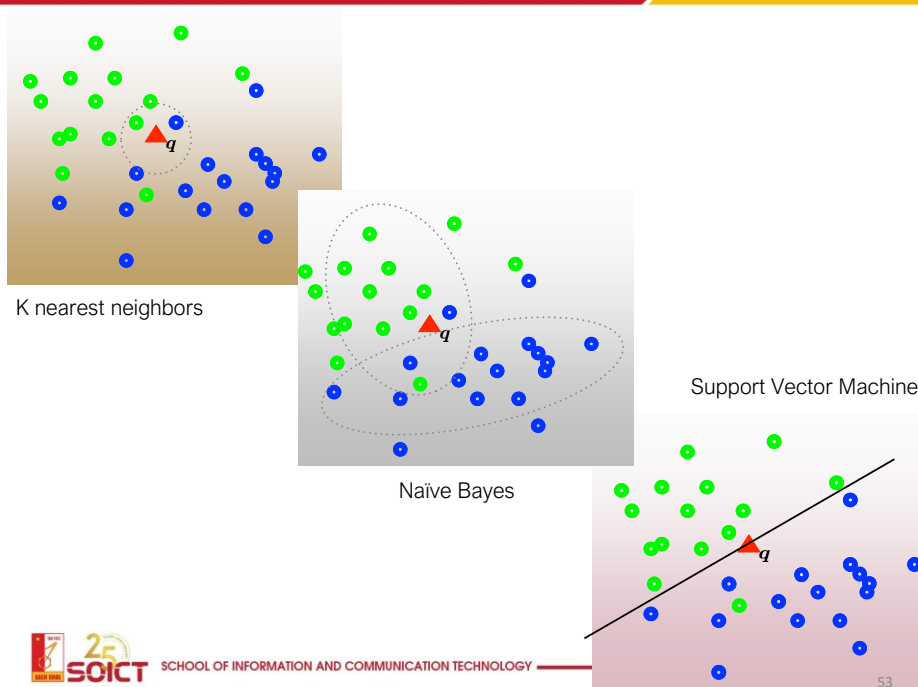




Dictionary Learning:
Learn Visual Words using clustering

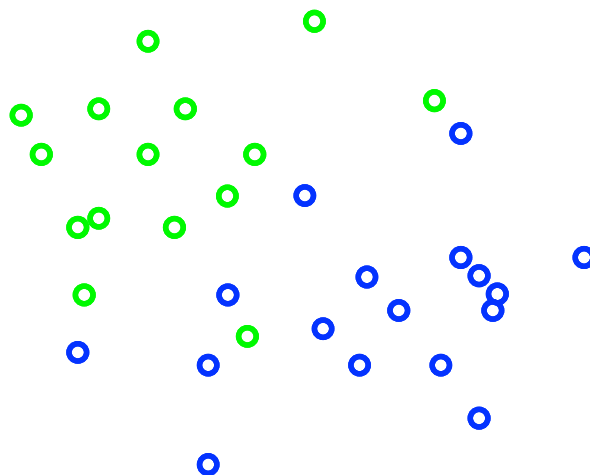
Encode:
build Bags-of-Words (BOW) vectors
for each image

Classify:
Train and test data using BOWs

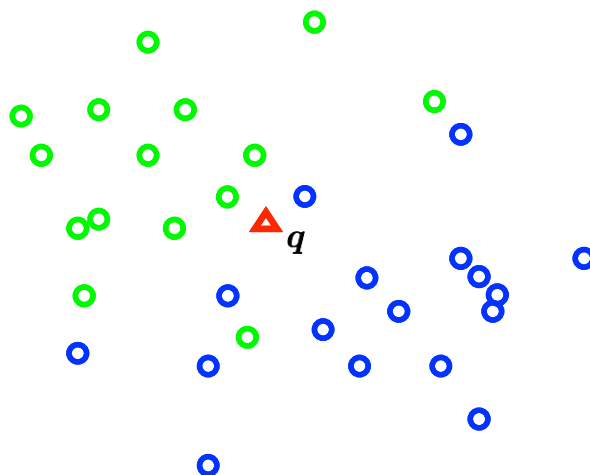


K nearest neighbors

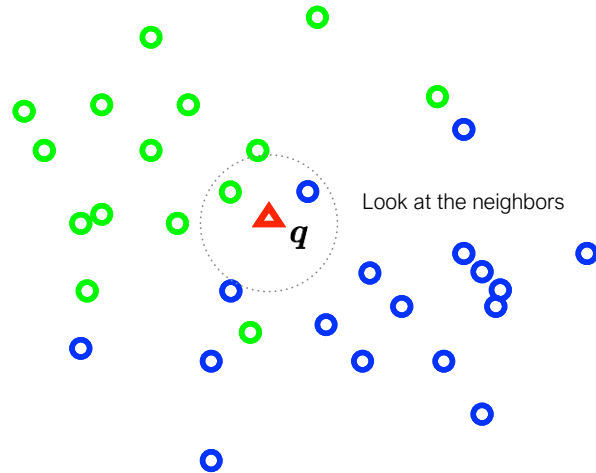
Distribution of data from two classes



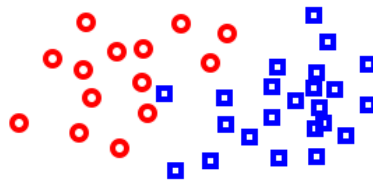
Distribution of data from two classes



Distribution of data from two classes



K-Nearest Neighbor (KNN) Classifier

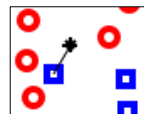


Non-parametric pattern classification approach

Consider a two class problem where each sample consists of two measurements (x,y).

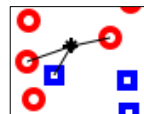
For a given query point q , assign the class of **the nearest neighbor**

$k = 1$



Compute the **k nearest neighbors** and assign the class by majority vote.

$k = 3$



Nearest Neighbor is competitive

		Test Error Rate (%)
Linear classifier (1-layer NN)		12.0
K-nearest-neighbors, Euclidean		5.0
K-nearest-neighbors, Euclidean, deskewed		2.4
MNIST Digit Recognition		
– Handwritten digits	K-NN, Tangent Distance, 16x16	1.1
– 28x28 pixel images: $d = 784$	K-NN, shape context matching	0.67
– 60,000 training samples	1000 RBF + linear classifier	3.6
– 10,000 test samples	SVM deg 4 polynomial	1.1
	2-layer NN, 300 hidden units	4.7
	2-layer NN, 300 HU, [deskewing]	1.6
	LeNet-5, [distortions]	0.8
	Boosted LeNet-4, [distortions]	0.7

Yann LeCunn



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

59

What is the best distance metric between data points?

- Typically Euclidean distance
- Locality sensitive distance metrics
- Important to normalize.
Dimensions have different scales

How many K?

- Typically $k=1$ is good
- Cross-validation (try different k !)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

60

Distance metrics

$$D(\mathbf{x}, \mathbf{y}) = \sqrt{(x_1 - y_1)^2 + \cdots + (x_N - y_N)^2} \quad \text{Euclidean}$$

$$D(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|} = \frac{x_1 y_1 + \cdots + x_N y_N}{\sqrt{\sum_n x_n^2} \sqrt{\sum_n y_n^2}} \quad \text{Cosine}$$

$$D(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \sum_n \frac{(x_n - y_n)^2}{(x_n + y_n)} \quad \text{Chi-squared}$$



Distance metrics

L1 (Manhattan) distance

L2 (Euclidean) distance

$$d_1(I_1, I_2) = \sum_p |I_1^p - I_2^p|$$

$$d_2(I_1, I_2) = \sqrt{\sum_p (I_1^p - I_2^p)^2}$$

- Two most commonly used special cases of p-norm

$$\|\mathbf{x}\|_p = (|x_1|^p + \cdots + |x_n|^p)^{\frac{1}{p}} \quad p \geq 1, \mathbf{x} \in \mathbb{R}^n$$



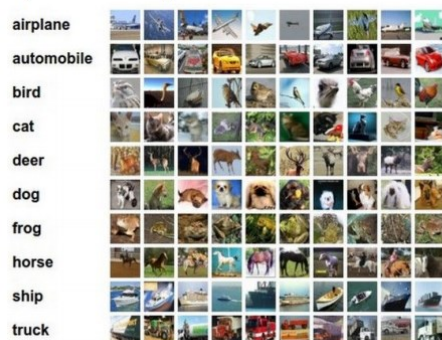
CIFAR-10 and NN results

Example dataset: **CIFAR-10**

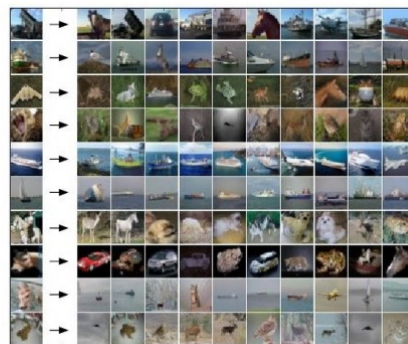
10 labels

50,000 training images

10,000 test images.



For every test image (first column),
examples of nearest neighbors in rows

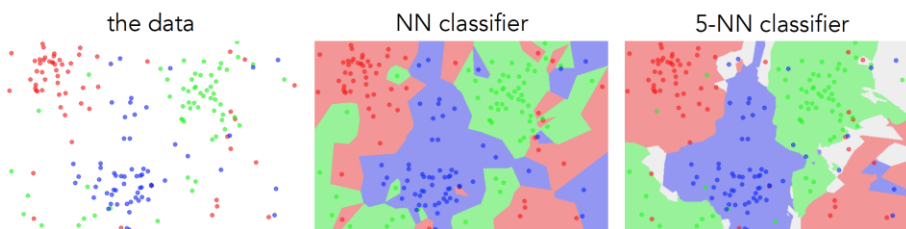


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

64

k-nearest neighbor

- Find the k closest points from training data
- Labels of the **k points “vote”** to classify



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

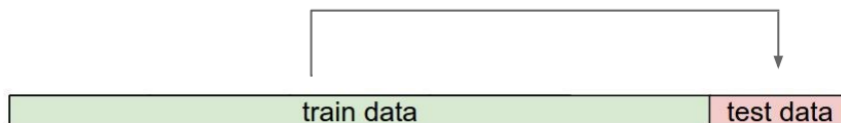
65

Hyperparameters

- What is the best distance to use?
- What is the best value of k to use?
- i.e., how do we set the hyperparameters?
- Very problem-dependent
- Must try them all and see what works best



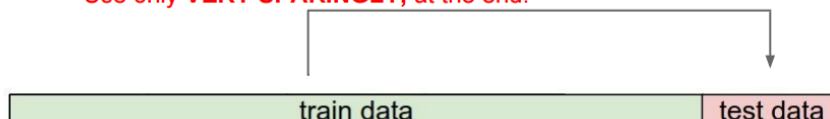
Try out what hyperparameters work best on test set.



Trying out what hyperparameters work best on test set:

Very bad idea. The test set is a proxy for the generalization performance!

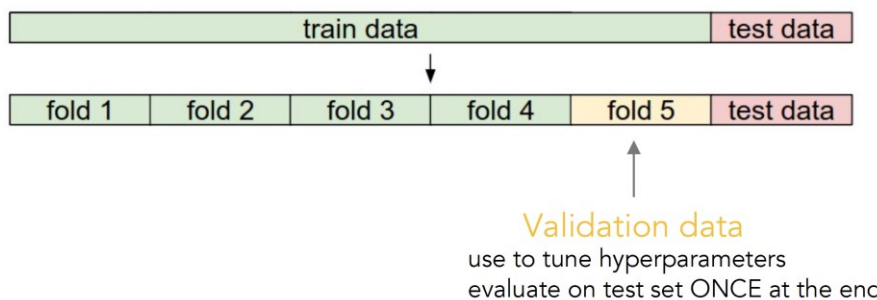
Use only **VERY SPARINGLY**, at the end.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

68

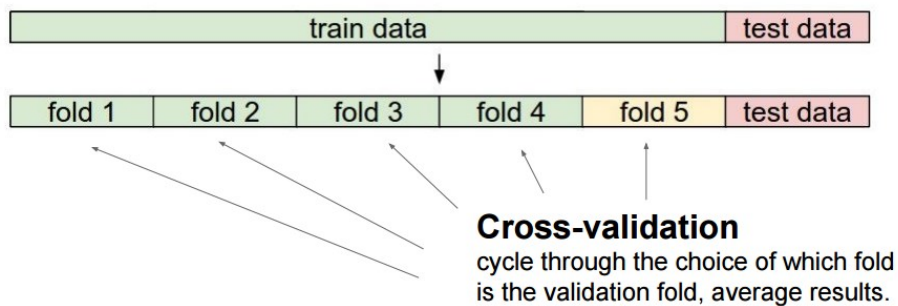
Validation



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

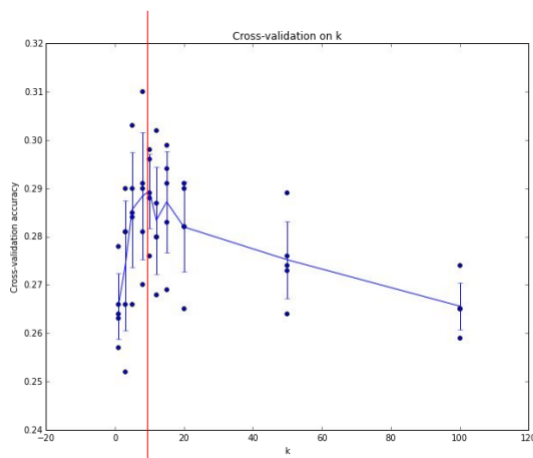
69

Cross-validation



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

70



Example of
5-fold cross-validation
for the value of k .

Each point: single
outcome.

The line goes
through the mean, bars
indicated standard
deviation

(Seems that $k \approx 7$ works best
for this data)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

71

How to pick hyperparameters?

- Methodology
 - Train and test
 - Train, validate, test
- Train for original model
- Validate to find hyperparameters
- Test to understand generalizability



kNN

Pros

- simple yet effective

Cons

- search is expensive (can be speed-up)
- storage requirements
- difficulties with high-dimensional data



kNN -- Complexity and Storage

- N training images, M test images
- Training: $O(1)$
- Testing: $O(MN)$
- Hmm...
 - Normally need the opposite
 - Slow training (ok), fast testing (necessary)



Other classifiers

- Naïve Bayes
- SVM
- Random Forest
- Neural Network
- ...



References

Most of these slides were adapted from:

1. Ioannis Yannis, Gkioulekas (16-385 Computer Vision, Spring 2020, CMU)
2. Kristen Grauman (CS 376: Computer Vision, Spring 2018, The University of Texas at Austin)
3. Noah Snavely (Cornell University)
4. Fei-Fei Li (Stanford University)

