



# Computer Vision

## Chapter 7: Object recognition (Part 1)

1

## Chapter 7: Object recognition (Part 1) Content

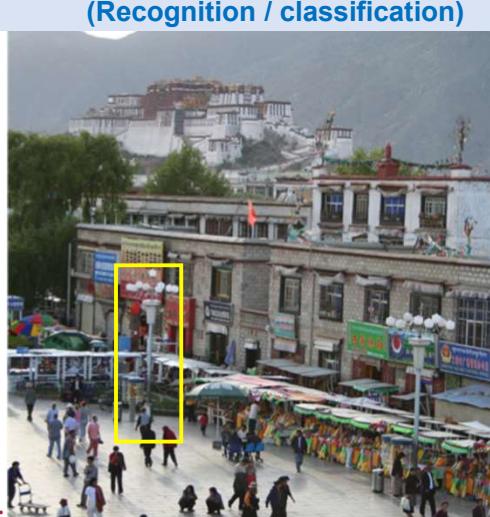
- Overview of ‘semantic vision’?
- Image classification/ recognition
- Bag-of-words
  - Recall
  - Vocabulary tree
- Classification
  - K nearest neighbours
  - Naïve Bayes
  - Support vector machine (SVM)

2

## Where are the people? (Detection)



3



4

**Is that Potala palace?  
(Identification)**

A photograph showing a wide-angle view of a lively outdoor market. In the background, a massive, white-walled building complex is built into a hillside, with numerous smaller structures and flags flying from its roofs. The foreground is filled with people walking through the market stalls, which have colorful awnings.

25 SOICT SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

5

**What's in the scene?  
(semantic segmentation)**

The same photograph as slide 5, but with semantic segmentation applied. The image is divided into several regions, each labeled: 'Mountain' (the hillside), 'Trees' (the green trees on the left), 'Building' (the large white building), 'Vendor(s)' (the market stalls), 'People' (the people walking), and 'Ground' (the paved area). The labels are placed near their respective segments.

25 SOICT SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

6

**What type of scene is it?  
(Scene categorization)**

The same photograph as slides 5 and 6, but with scene categorization labels overlaid. The labels 'Outdoor', 'Marketplace', and 'City' are written in yellow text and positioned over the respective elements in the image: the open air, the market stalls, and the urban buildings.

25 SOICT SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

7

**What are these people doing?  
(Activity / Event Recognition)**

The same photograph as slides 5, 6, and 7, showing a general view of the market square with people walking around and interacting.

25 SOICT SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

8

## Object recognition Is it really so hard?

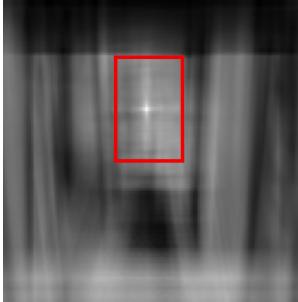
This is a chair



Find the chair in this image



Output of normalized correlation

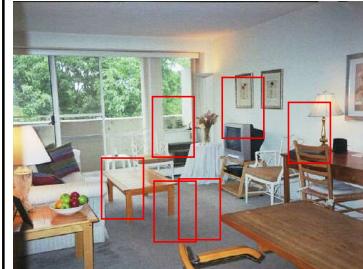


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

9

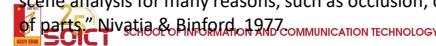
## Object recognition Is it really so hard?

Find the chair in this image



Pretty much garbage

A “popular method is that of template matching, by point to point correlation of a model pattern with the image pattern. These techniques are inadequate for three-dimensional scene analysis for many reasons, such as occlusion, changes in viewing angle, and articulation of parts.” Nivat & Binford, 1977



10

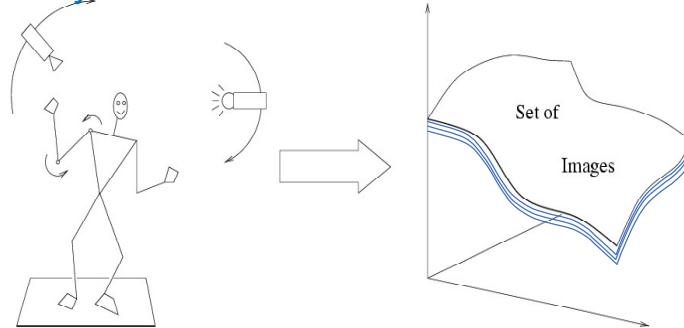
## And it can get a lot harder



Brady, M. J., & Kersten, D. (2003). Bootstrapped learning of novel objects. *J Vis*, 3(6), 413-422

11

## Why is this hard?



Variability:  
Camera position  
Illumination  
Shape parameters



12



13

### Challenge: variable viewpoint



14

### Challenge: variable illumination

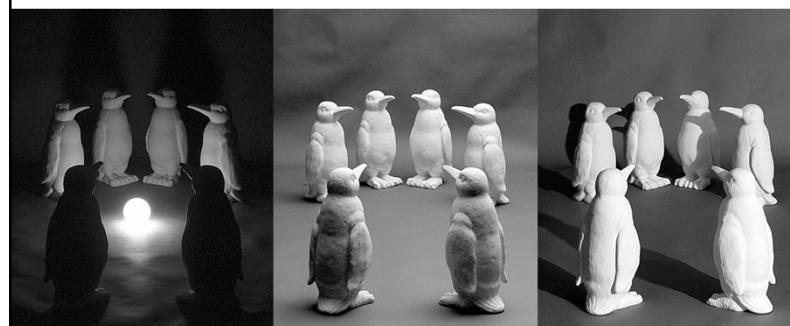


image credit: J. Koenderink

15

### Challenge: scale



16

## Challenge: deformation



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

17

## Challenge: Occlusion



Magritte, 1957

18

## Challenge: background clutter

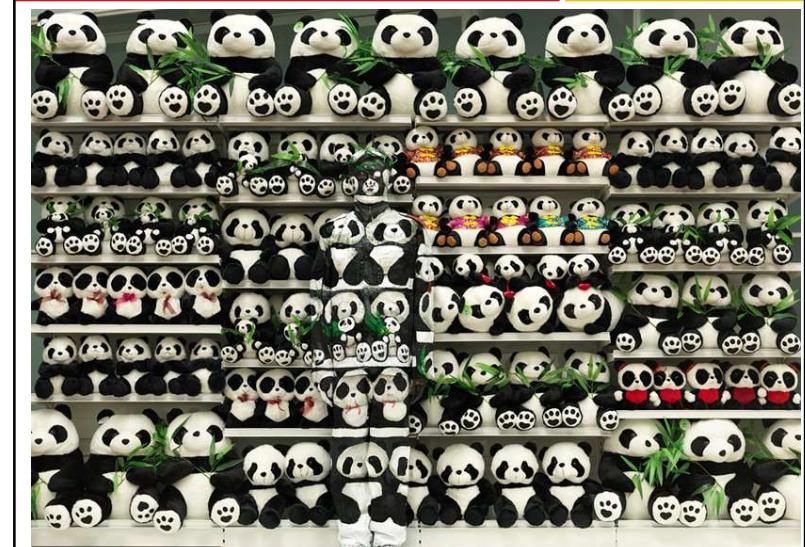


Kilmeny Niland. 1995



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

19



Challenge: Background clutter

20

Challenge: intra-class variations



**SOICT** SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Svetlana Lazebnik

21

## Image Classification/ Recognition

**SOICT** SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

22

## Image Classification/ Recognition



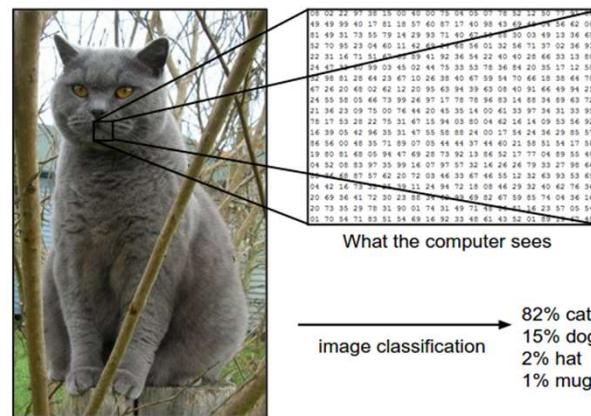
(assume given set of discrete labels)  
{dog, cat, truck, plane, ...}

→ cat

**SOICT** SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

23

## Image Classification: Problem



What the computer sees

0.02	22	91	38	15	0.45	0.05	0.04	0.07	0.18	0.12	0.39	0.15	0.16	0.42	0.01				
33	41	31	73	19	14	20	81	11	0.14	0.04	0.03	0.03	0.13	0.36	0.01				
33	33	33	33	33	33	33	33	33	0.04	0.04	0.04	0.04	0.04	0.04	0.04				
22	33	14	73	53	14	14	42	36	54	22	40	40	28	66	33	13	80		
24	41	17	39	03	45	02	44	73	33	53	78	36	64	20	35	17	12	50	
24	34	91	23	23	23	10	26	38	40	47	59	54	70	46	18	38	44	70	
24	34	34	48	04	12	12	12	12	12	12	12	12	12	12	12	12	12	21	
24	55	58	05	66	73	99	24	97	17	78	78	96	63	14	88	34	69	43	72
23	34	23	09	73	00	76	41	20	45	35	14	00	41	33	97	34	31	33	95
24	57	53	28	22	75	31	61	15	14	00	04	42	14	14	09	53	56	92	
14	14	14	14	14	14	14	14	14	14	14	14	14	14	14	14	14	14	14	95
14	54	00	49	33	71	89	07	05	44	44	37	44	40	21	56	51	54	27	58
19	80	81	49	05	94	47	69	28	73	92	13	32	17	77	04	59	55	40	
24	52	08	83	97	35	99	14	07	97	37	44	36	26	79	33	27	98	46	
24	42	14	11	11	11	11	11	11	11	11	11	11	11	11	11	11	11	11	59
20	69	34	41	72	32	23	85	32	32	49	92	47	59	85	74	04	36	16	
20	73	35	29	78	31	90	01	74	31	49	73	73	73	73	73	73	73	73	
23	70	34	71	03	51	54	49	14	92	33	48	61	43	52	01	31	31	31	

image classification →

82% cat  
15% dog  
2% hat  
1% mug

**SOICT** SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

24

## 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



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

25

## A simple pipeline - Training

Training  
Images



Image  
Features



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

26

## A simple pipeline - Training

Training  
Images

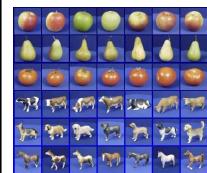


Image  
Features

Training  
Labels

Training



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

27

## A simple pipeline - Training

Training  
Images



Image  
Features

Training  
Labels

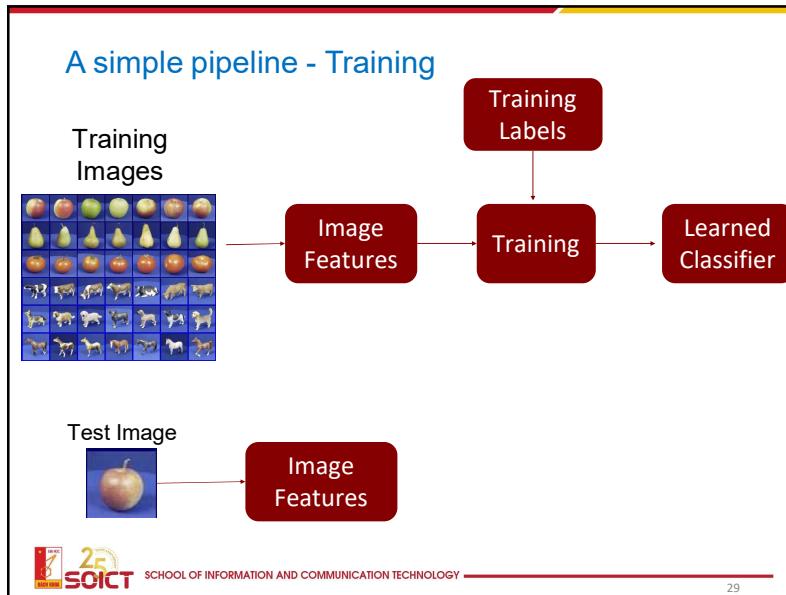
Training

Learned  
Classifier

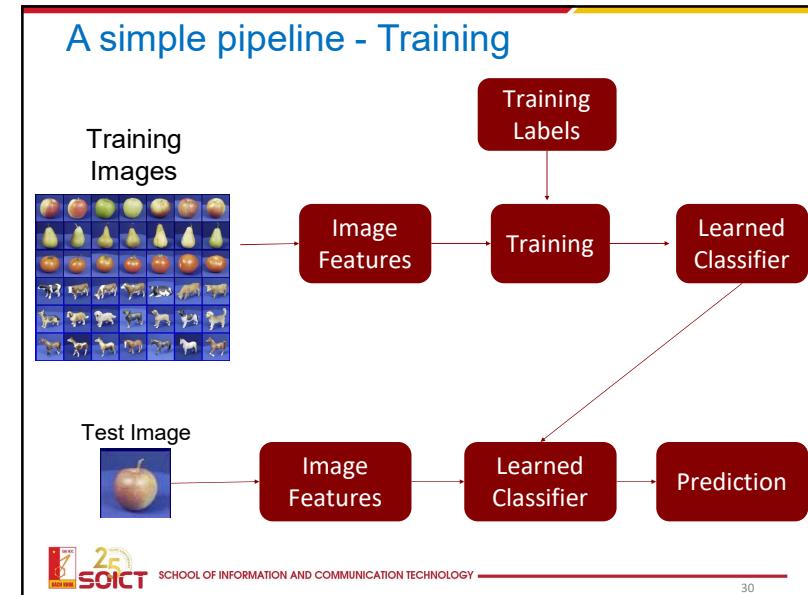


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

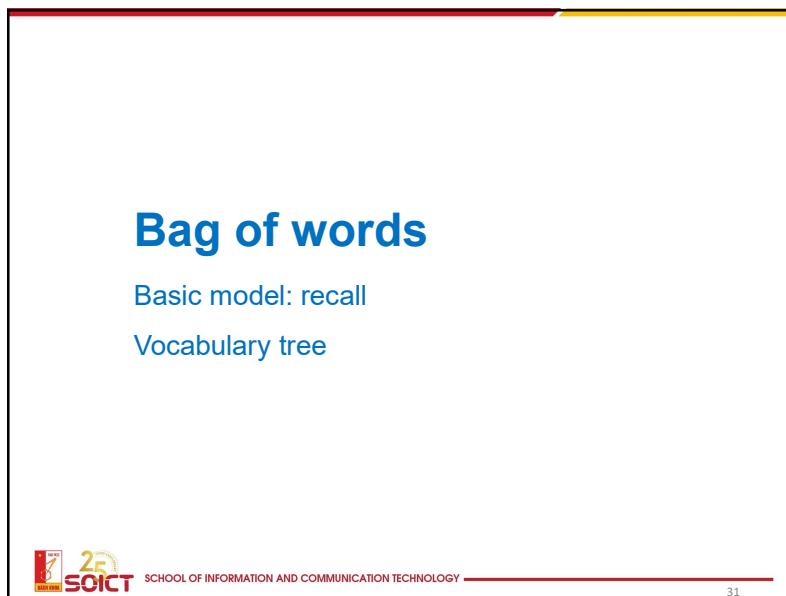
28



29



30



31



32



33



34

CalTech6 dataset						
class	bag of features		bag of features		Parts-and-shape model	
	Zhang et al. (2005)	Willamowski et al. (2004)			Fergus et al. (2003)	
airplanes	<b>98.8</b>		97.1		90.2	
cars (rear)	98.3		<b>98.6</b>		90.3	
cars (side)	<b>95.0</b>		87.3		88.5	
faces	<b>100</b>		99.3		96.4	
motorbikes	<b>98.5</b>		98.0		92.5	
spotted cats	<b>97.0</b>	—			90.0	

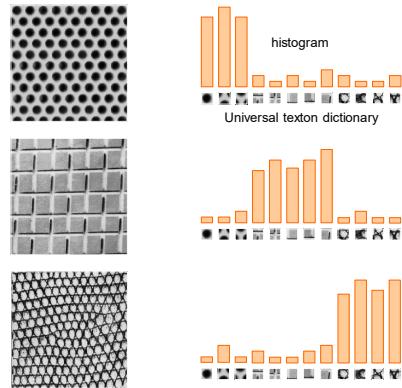
Works pretty well for image-level classification

35



36

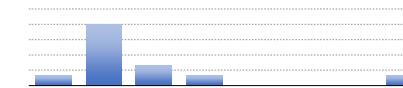
## Texture recognition



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

37

## Vector Space Model



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

38

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

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

just a histogram over words

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

What is the similarity between two documents?



Use any distance you want but the cosine distance is fast.

$$d(\mathbf{v}_i, \mathbf{v}_j) = \cos \theta$$

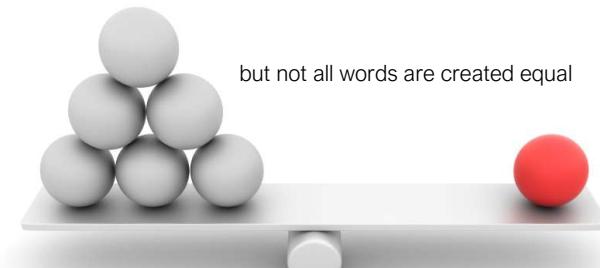
$$= \frac{\mathbf{v}_i \cdot \mathbf{v}_j}{\|\mathbf{v}_i\| \|\mathbf{v}_j\|}$$



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

39

but not all words are created equal



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

40

40

## TF-IDF

Term Frequency Inverse Document Frequency

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

weigh each word by a heuristic

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

$$n(w_{i,d})\alpha_i = n(w_{i,d}) \log \left\{ \frac{\text{term frequency}}{\sum_{d'} \mathbf{1}[w_i \in d']} \right\}$$

(down-weights common terms)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

41

## Standard BOW pipeline (for image classification)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

42

**Dictionary Learning:**  
Learn Visual Words using clustering

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

**Classify:**  
Train and test data using BOWs

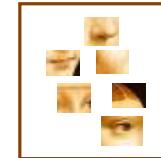


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

43

**Dictionary Learning:**  
Learn Visual Words using clustering

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



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

44

**Dictionary Learning:**  
Learn Visual Words using clustering

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

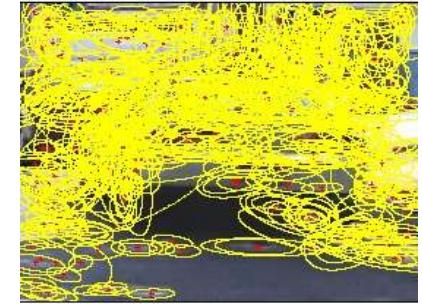


25 SOICT SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

45

*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)



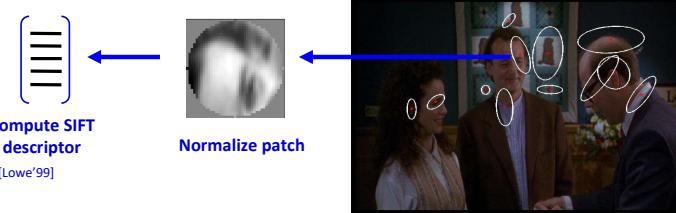
25 SOICT SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

46

**Detect patches**

Compute SIFT descriptor  
[Lowe'99]

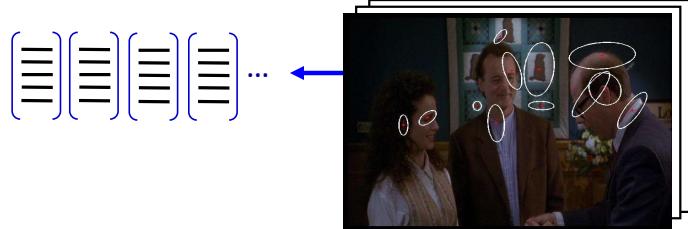
Normalize patch



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

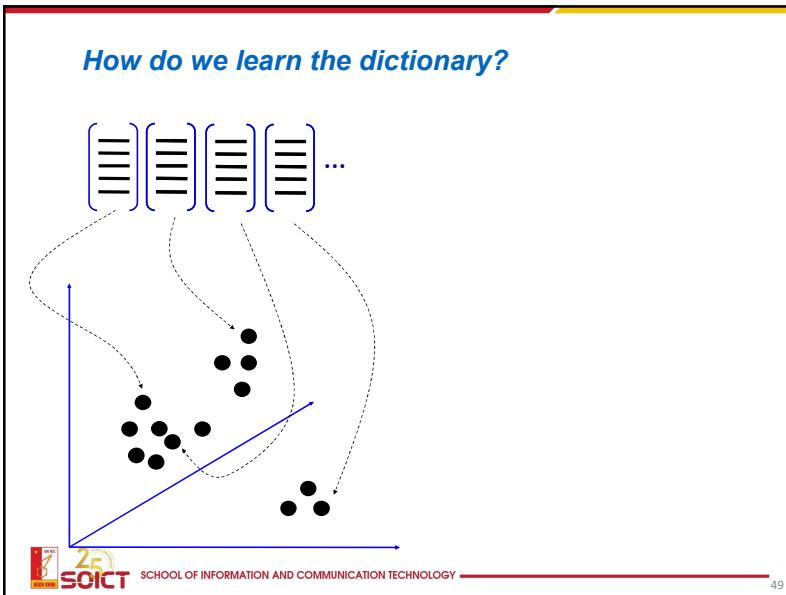
25 SOICT SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

47

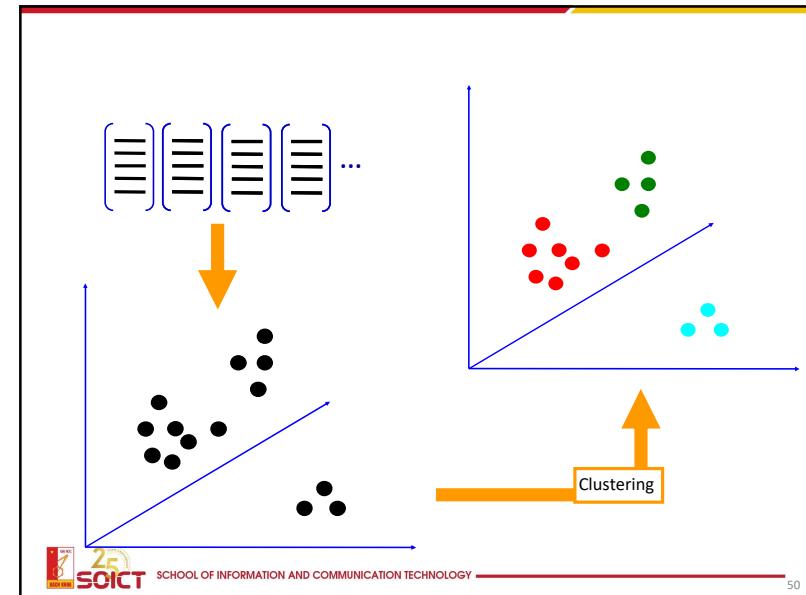


25 SOICT SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

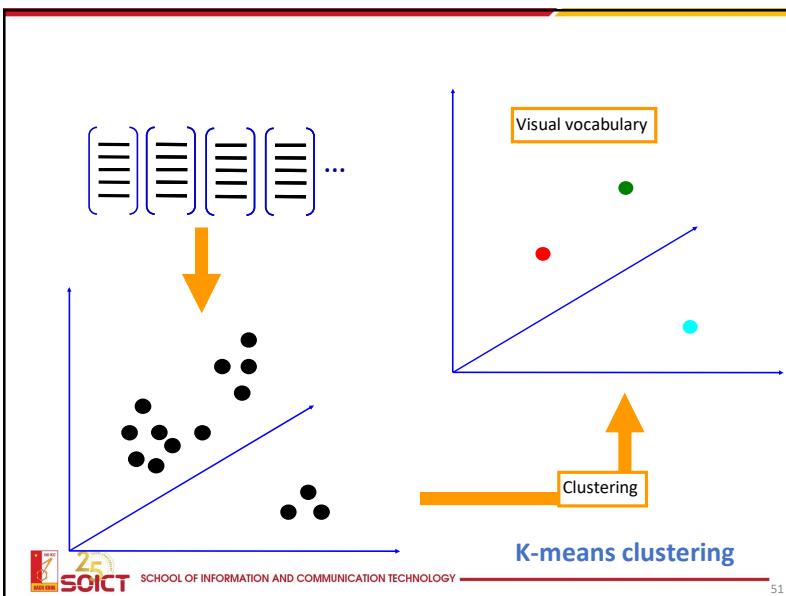
48



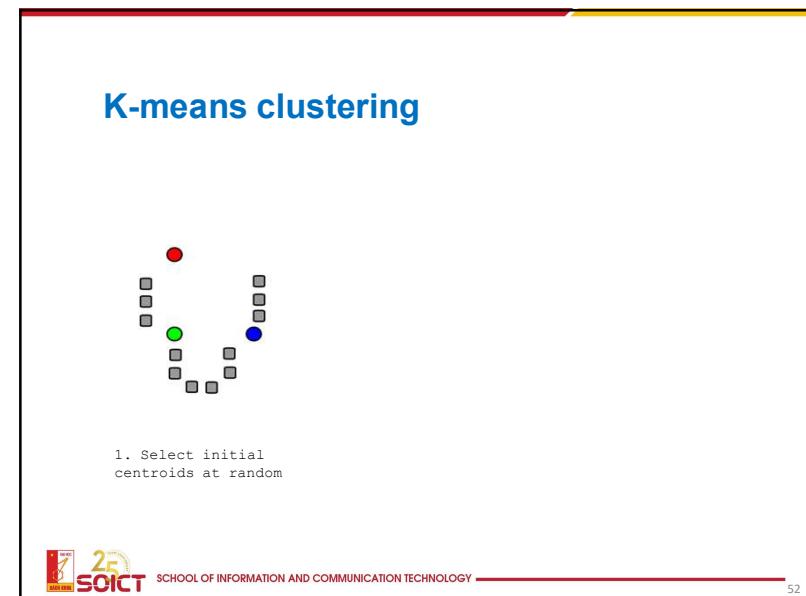
49



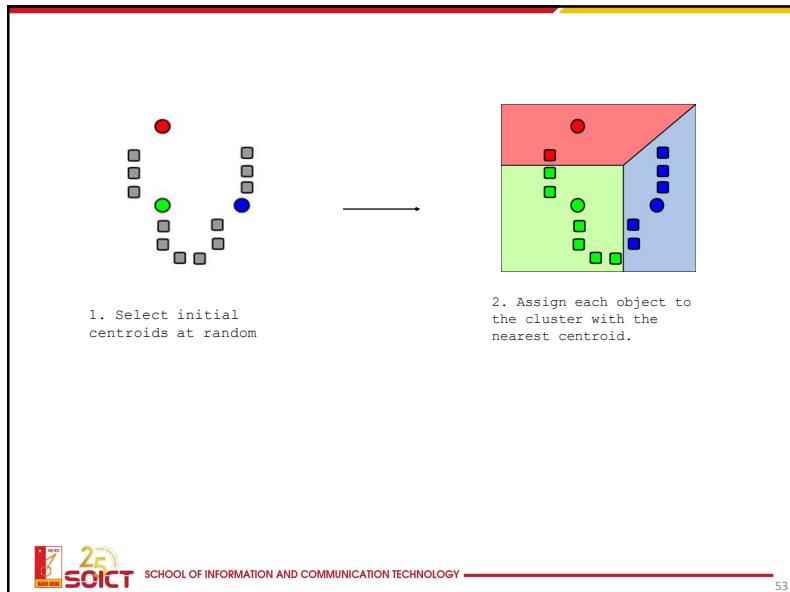
50



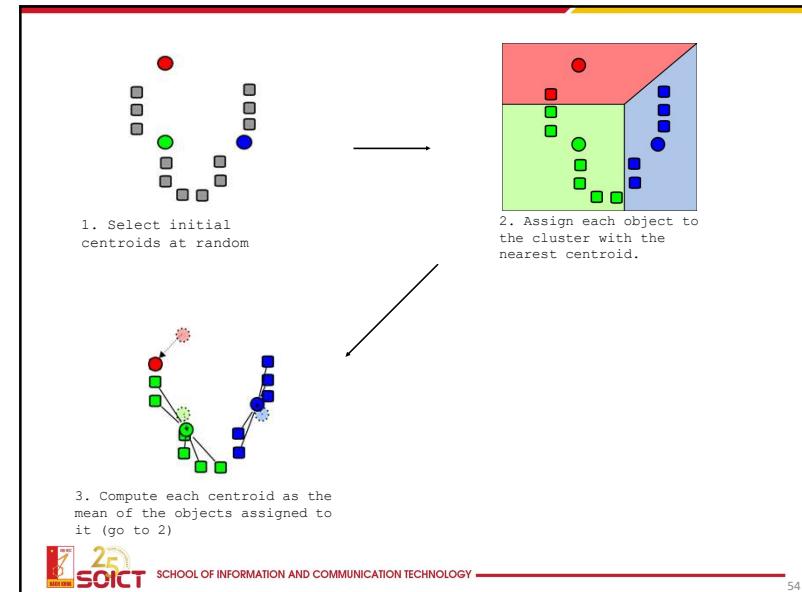
51



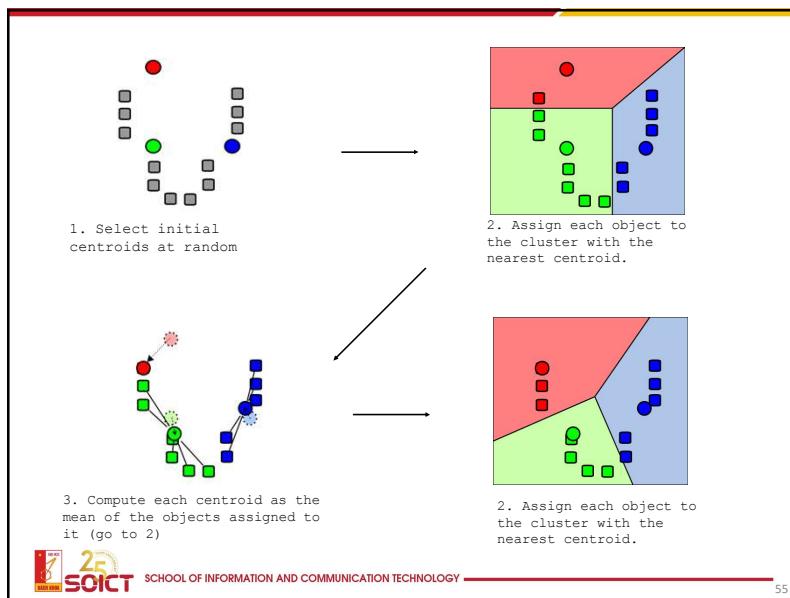
52



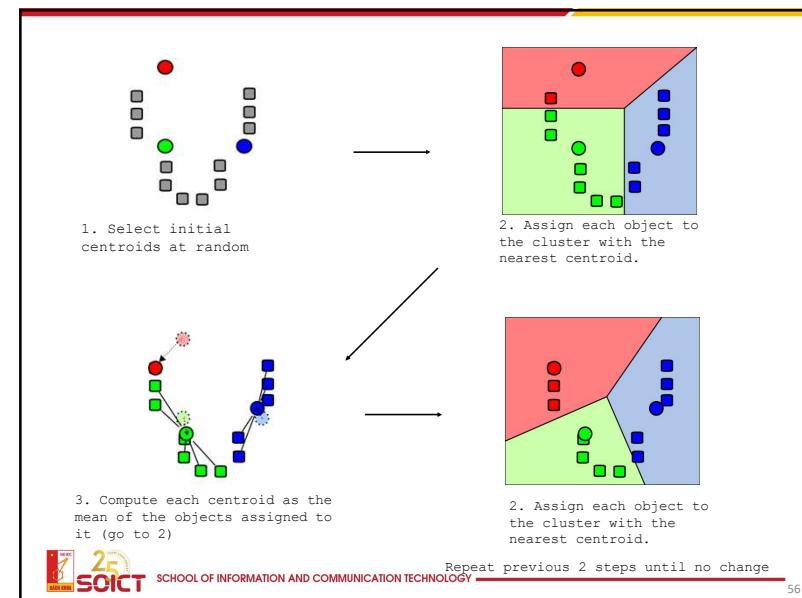
53



54



55



56

## K-means Clustering Algorithm

Given  $k$ :

1. Select initial centroids at random.
2. Assign each object to the cluster with the nearest centroid.
3. Compute each centroid as the mean of the objects assigned to it.
4. Repeat previous 2 steps until no change.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

57

## *From what data should I learn the dictionary?*

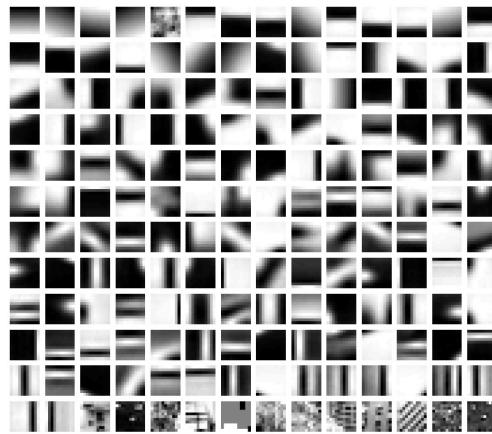
- Dictionary can be learned on separate training set
- Provided the training set is sufficiently representative, the dictionary will be “universal”



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

58

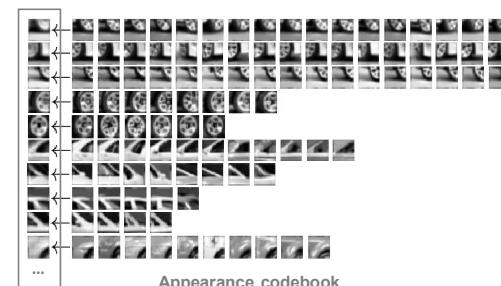
## Example visual dictionary



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

59

## Example dictionary

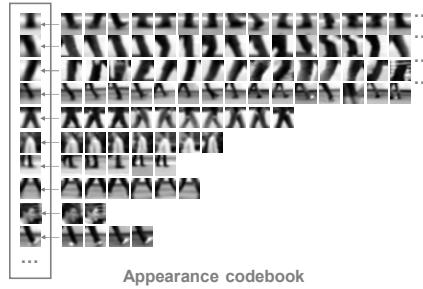
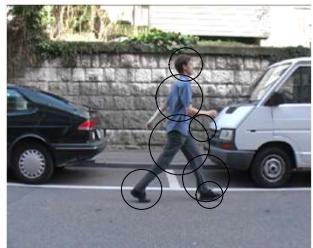


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Source: B. Leibe

60

## Another dictionary



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Source: B. Leibe

61

**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

62



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

### Encode:

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



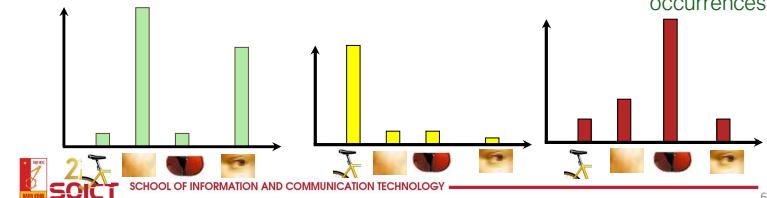
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

63

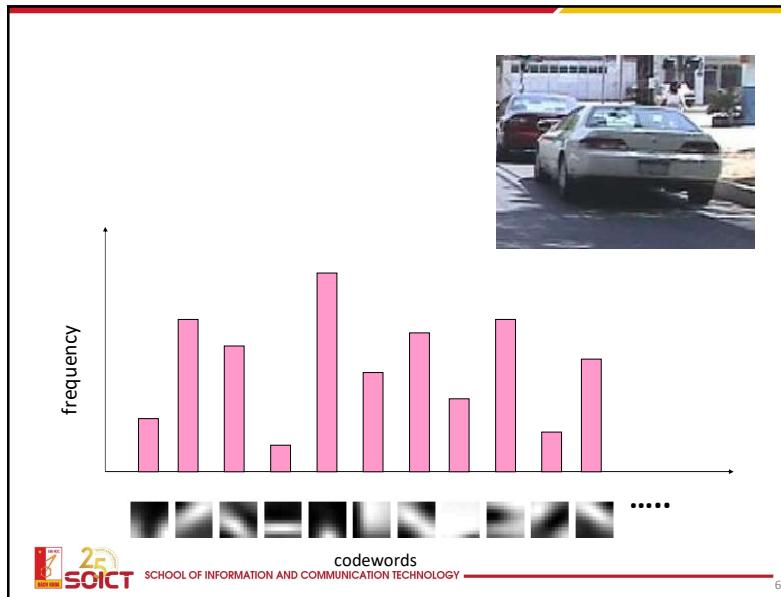
### Encode:

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

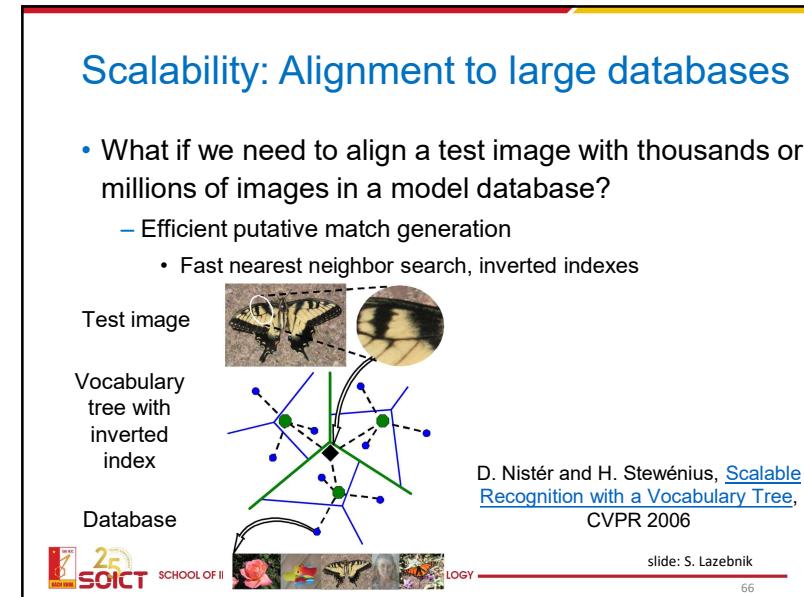
2. Histogram: count the  
number of visual word  
occurrences



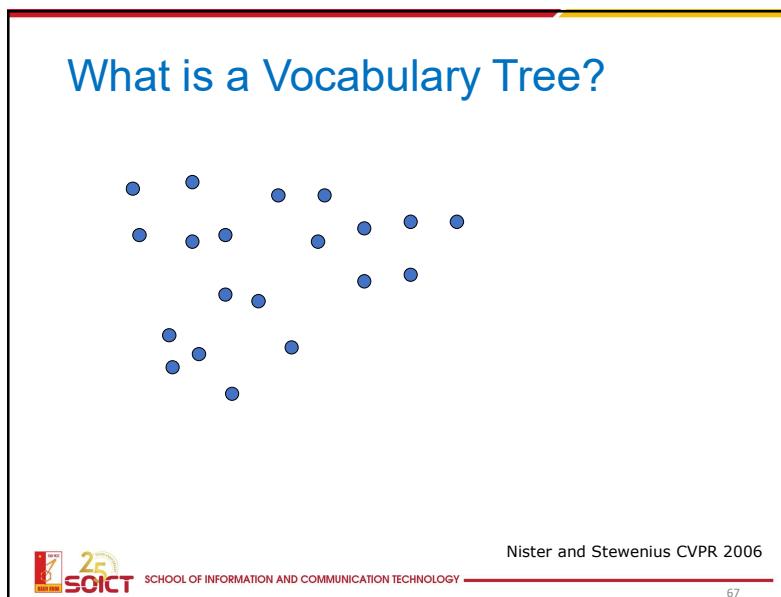
64



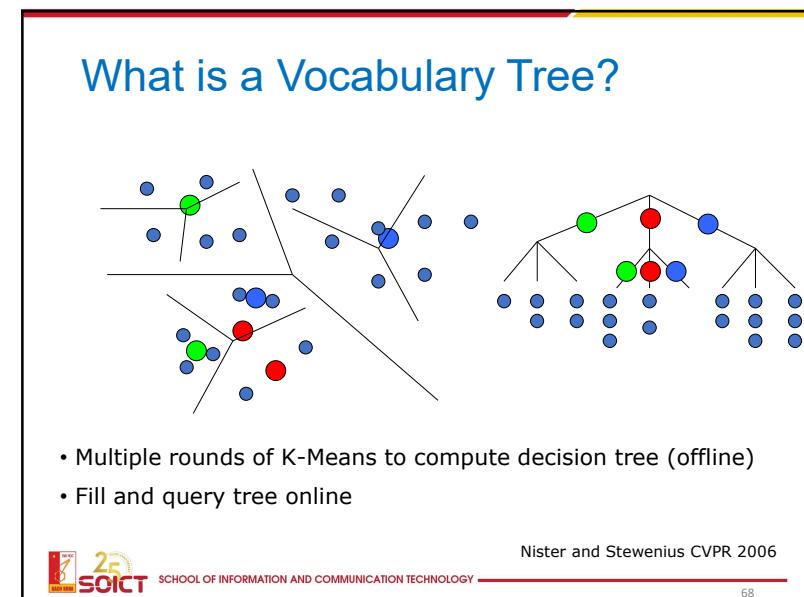
65



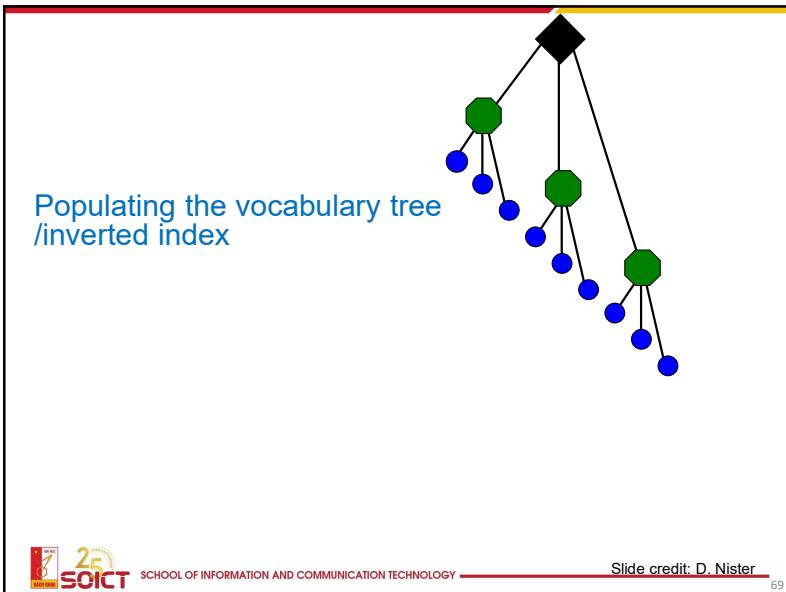
66



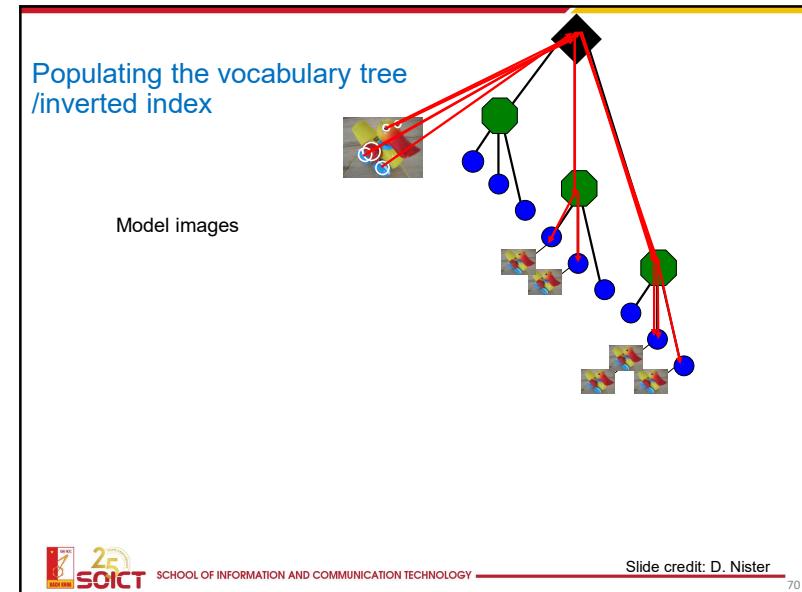
67



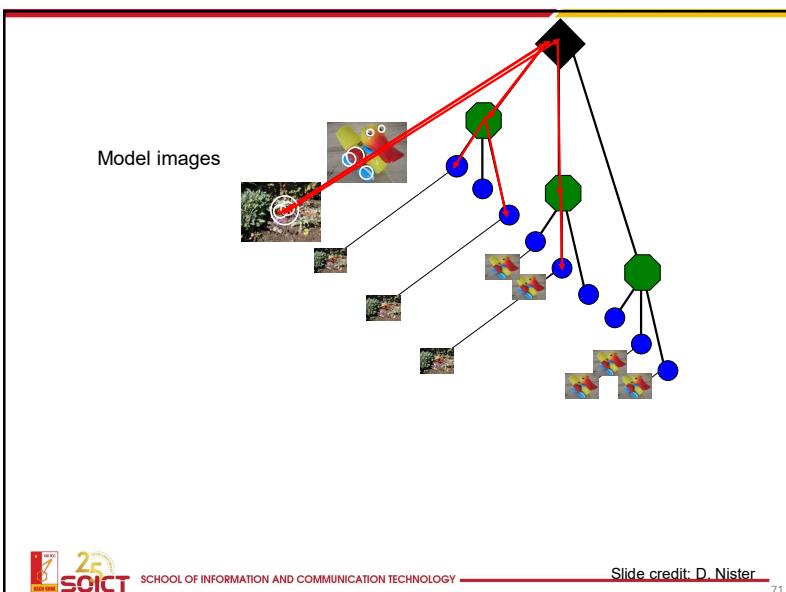
68



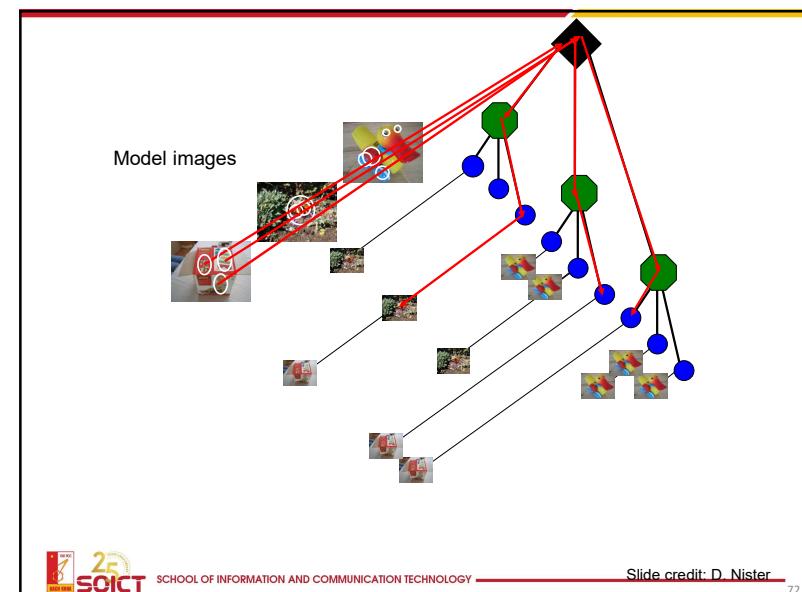
69



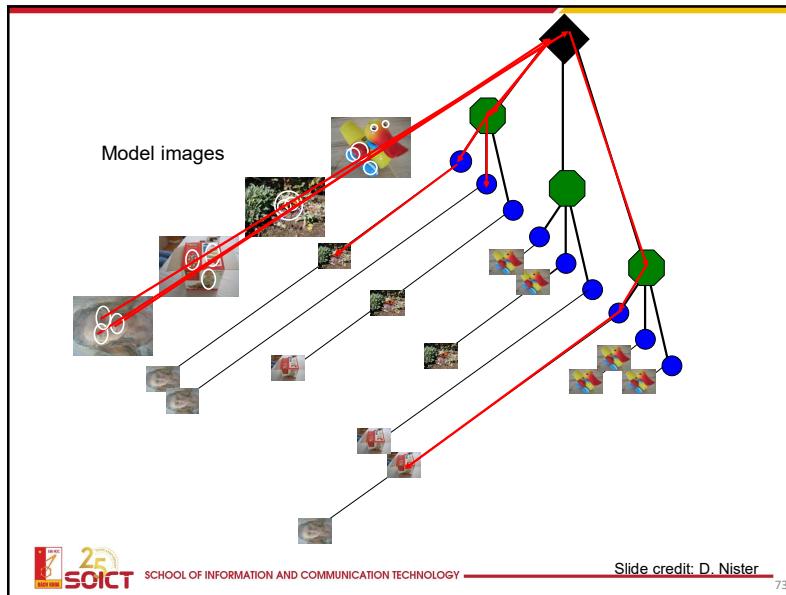
70



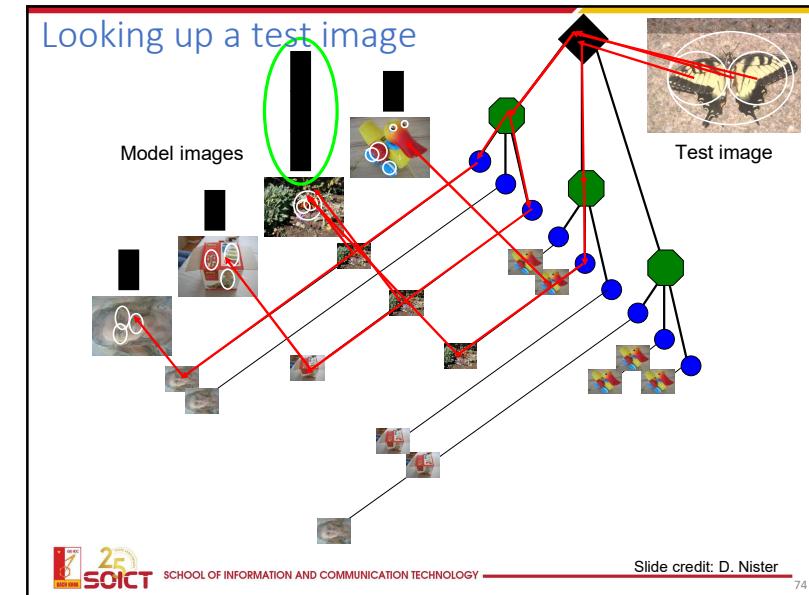
71



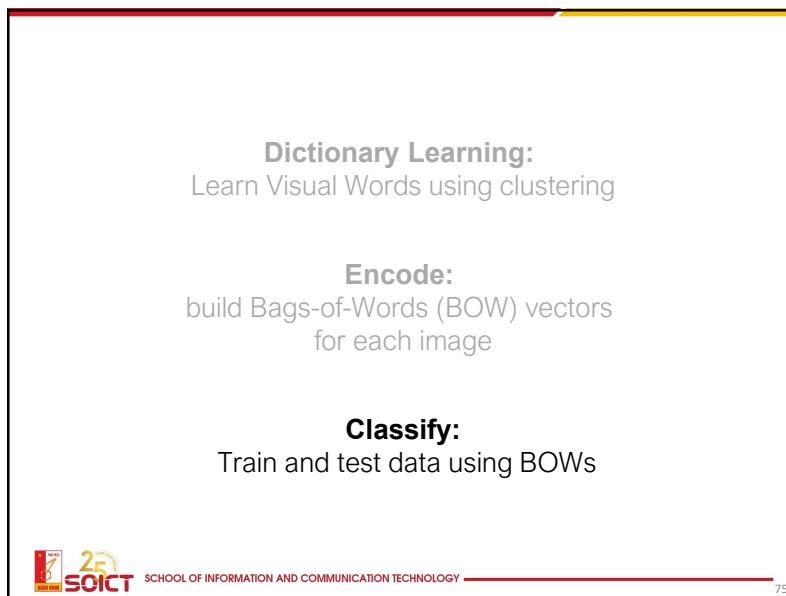
72



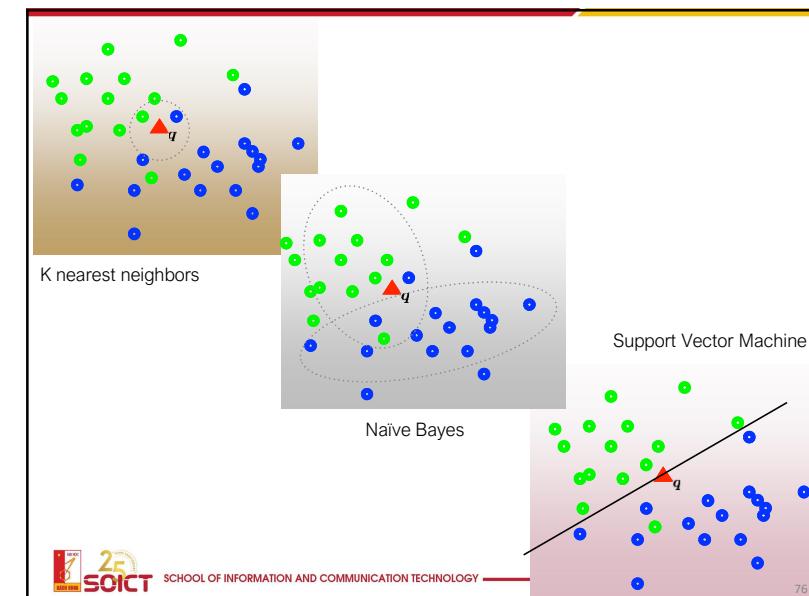
73



74



75



76

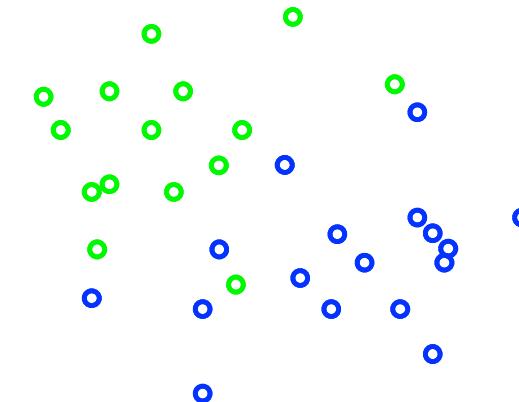
## Classification K nearest neighbors



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

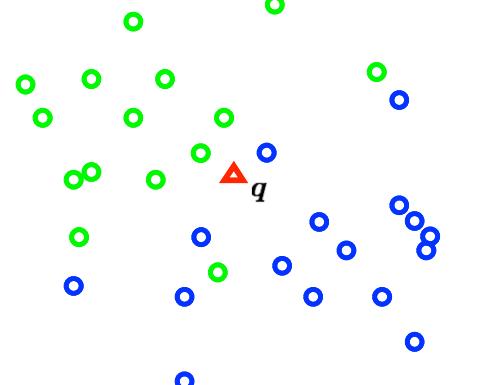
77

Distribution of data from two classes



78

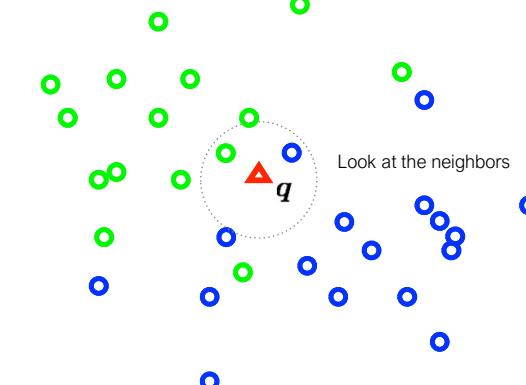
Distribution of data from two classes

Which class does  $q$  belong too?

SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

79

Distribution of data from two classes



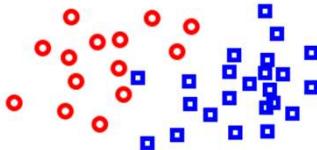
80



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

80

## K-Nearest Neighbor (KNN) Classifier

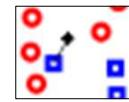


**Non-parametric** pattern classification approach

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

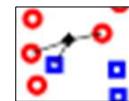
For a given query point q, assign the class of **the nearest neighbor**

k = 1



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

k = 3



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

81

## 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
- 28x28 pixel images: d = 784
- 60,000 training samples
- 10,000 test samples



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

	Test Error Rate (%)
Linear classifier (1-layer NN)	12.0
K-nearest-neighbors, Euclidean	5.0
K-nearest-neighbors, Euclidean, deskewed	2.4
1000 RBF + linear classifier	3.6
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

82

## 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

83

## Distance metrics

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

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

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



84

## Distance metrics

L1 (Manhattan) distance

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

L2 (Euclidean) distance

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

- Two most commonly used special cases of p-norm

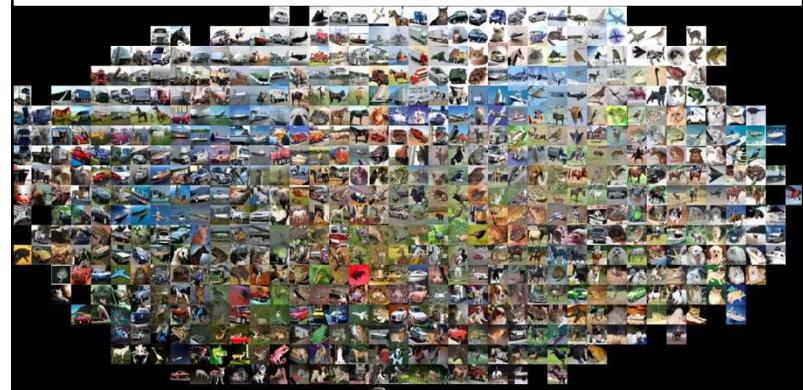
$$\|x\|_p = \left( |x_1|^p + \dots + |x_n|^p \right)^{\frac{1}{p}} \quad p \geq 1, x \in \mathbb{R}^n$$



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

85

## Visualization: L2 distance



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

86

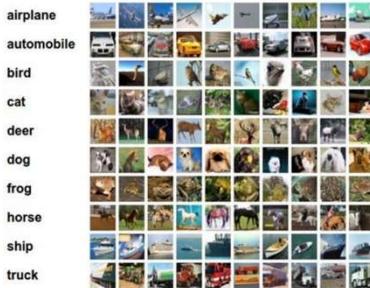
## CIFAR-10 and NN results

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

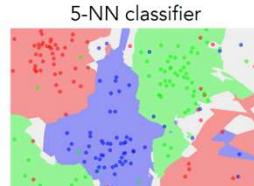
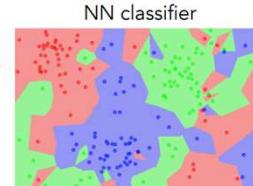
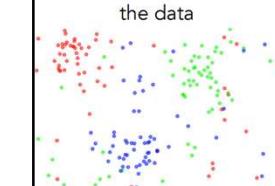


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

87

## 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

88

## 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



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

89

Try out what hyperparameters work best on test set.



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

90

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.



train data      test data



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

91

## Validation

train data      test data

fold 1    fold 2    fold 3    fold 4    fold 5    test data

Validation data

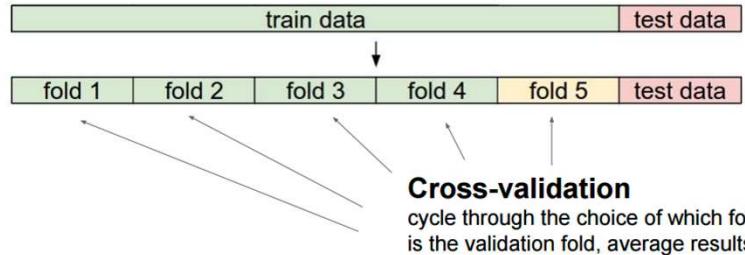
use to tune hyperparameters  
evaluate on test set ONCE at the end



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

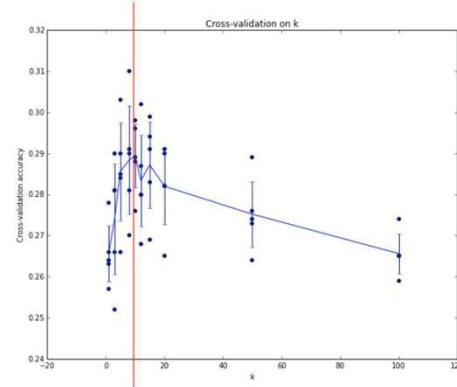
92

## Cross-validation



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

93



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

94

## How to pick hyper-parameters?

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



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

95

## kNN

### Pros

- simple yet effective

### Cons

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



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

96

95

96

## 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)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

97

## Classification Naïve Bayes

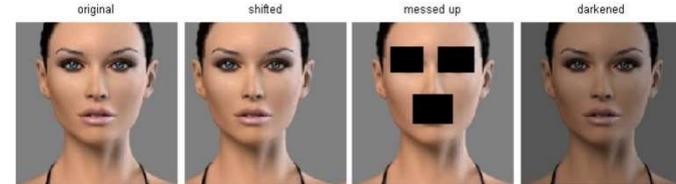


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

99

k-Nearest Neighbor on images **never used**.

- terrible performance at test time
- distance metrics on level of whole images can be very unintuitive



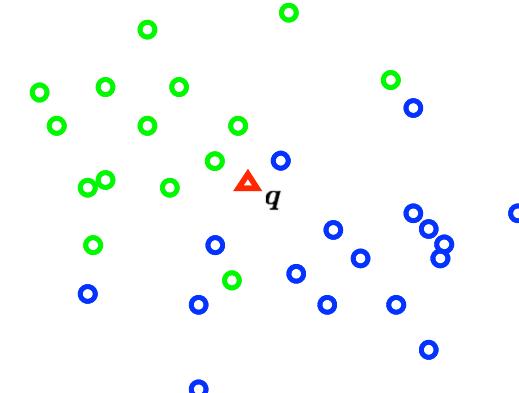
(all 3 images have same L2 distance to the one on the left)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

98

Distribution of data from two classes



*Which class does q belong to?*

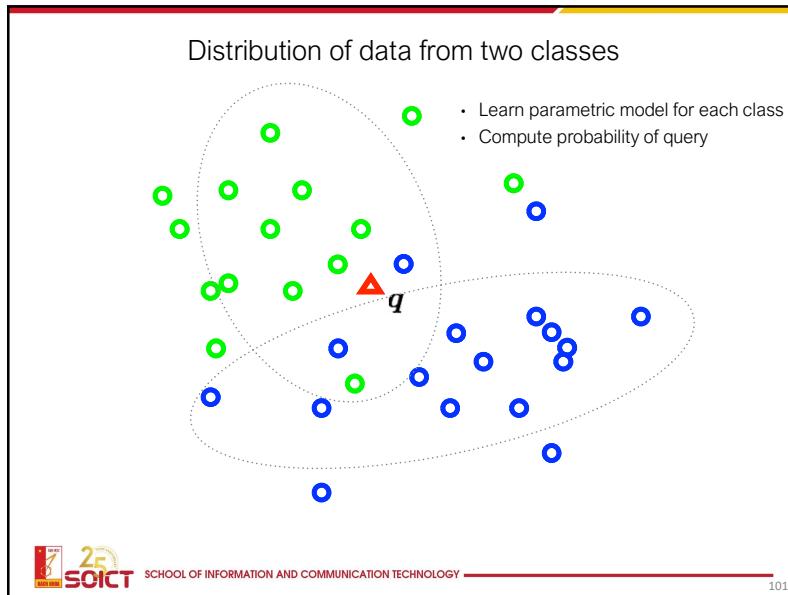


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

100

99

100



101

This is called the posterior:  
the probability of a class  $z$  given the observed features  $X$

$$p(z|X)$$

For classification,  $z$  is a discrete random variable (e.g., car, person, building)

$X$  is a set of observed features (e.g., features from a single image)

(it's a function that returns a single probability value)

**SOICT 25 Years** SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

102

102

This is called the posterior:  
the probability of a class  $z$  given the observed features  $X$

$$p(z|x_1, \dots, x_N)$$

For classification,  $z$  is a discrete random variable (e.g., car, person, building)

Each  $x$  is an observed feature (e.g., visual words)

(it's a function that returns a single probability value)

**SOICT 25 Years** SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

103

103

**Recall:**

The posterior can be decomposed according to **Bayes' Rule**

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$

In our context...

$$p(z|x_1, \dots, x_N) = \frac{p(x_1, \dots, x_N|z)p(z)}{p(x_1, \dots, x_N)}$$

**SOICT 25 Years** SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

104

104

The naive Bayes' classifier is solving this optimization

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(z | \mathbf{X})$$

MAP (maximum a posteriori) estimate

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} \frac{p(\mathbf{X}|z)p(z)}{p(\mathbf{X})}$$

Bayes' Rule

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(\mathbf{X}|z)p(z)$$

Remove constants

To optimize this...we need to compute this

 Compute the likelihood...



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

105

A naive Bayes' classifier assumes all features are  
**conditionally independent**

$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_N | \mathbf{z}) &= p(\mathbf{x}_1 | \mathbf{z})p(\mathbf{x}_2, \dots, \mathbf{x}_N | \mathbf{z}) \\ &= p(\mathbf{x}_1 | \mathbf{z})p(\mathbf{x}_2 | \mathbf{z})p(\mathbf{x}_3, \dots, \mathbf{x}_N | \mathbf{z}) \\ &= p(\mathbf{x}_1 | \mathbf{z})p(\mathbf{x}_2 | \mathbf{z}) \cdots p(\mathbf{x}_N | \mathbf{z}) \end{aligned}$$

Recall:



$$p(x, y) = p(x|y)p(y)$$



$$p(x, y) = p(x)p(y)$$



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

106

To compute the MAP estimate

Given (1) a set of known parameters

$$p(\mathbf{z}) \quad p(\mathbf{x}|\mathbf{z})$$

(2) observations

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

Compute which  $\mathbf{z}$  has the largest probability

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(z) \prod_n p(x_n | z)$$



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

107



DARPA Selects Carnegie Mellon University's Tartan Arm

State's first robot, CHIMP, has been selected by DARPA to compete in the DARPA Robotics Challenge.

The competition will test robots' ability to perform tasks such as driving a vehicle, climbing stairs, opening doors, cutting a cable, and closing valves.

The competition is set for July 26.

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

$$\begin{aligned} p(X|z) &= \prod_v p(x_v|z)^{c(w_v)} \\ &= (0.09)^1 (0.55)^6 \cdots (0.09)^1 \end{aligned}$$

Numbers get really small so use log probabilities

$$\log p(X|z = \text{'grandchallenge'}) = -2.42 - 3.68 - 3.43 - 2.42 - 0.07 - 0.07 - 0.07 - 2.42 = -14.58$$

$$\log p(X|z = \text{'softrobot'}) = -7.63 - 9.37 - 15.18 - 2.97 - 0.02 - 0.01 - 0.02 - 2.27 = -37.48$$

\* typically add pseudo-counts (0.001)

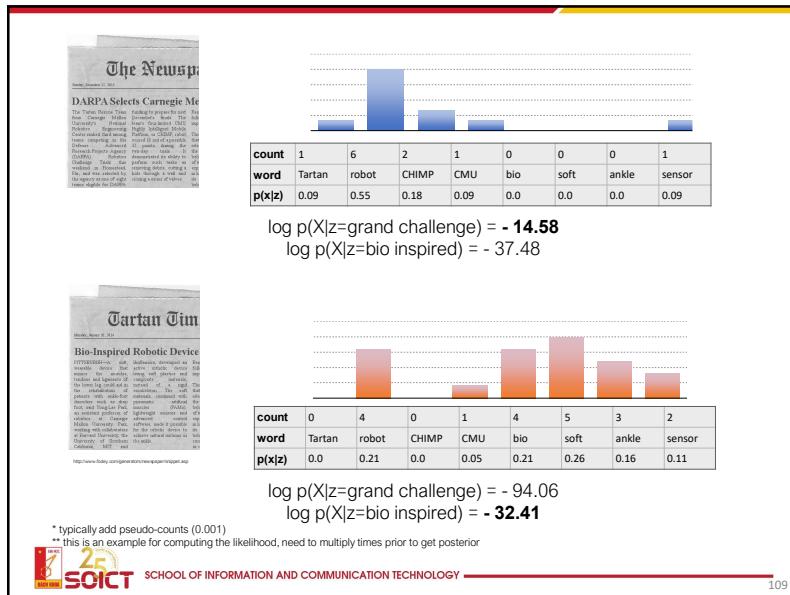
\*\* this is an example for computing the likelihood, need to multiply times **prior** to get posterior



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

108

108



109

## Multinomial Naïve Bayes: Learning

- From training corpus, extract *Vocabulary*

- Calculate  $P(z)$  terms

- For each  $z$  in  $C$  do

$$p(z) = \frac{|docs_z|}{\text{total number of documents}}$$

$docs_z \leftarrow$  all docs with class  $= z$

- Calculate  $P(x_i | z)$  terms

- $Text_z \leftarrow$  single doc containing all  $docs_z$

- For each word  $x_i$  in *Vocabulary*

$$n_i \leftarrow \# \text{ of occurrences of } x_i \text{ in } Text_z$$

$$p(x_i | z) = (n_i + 1) / (\sum_{x_i \in \text{vocabulary}} (n_i + 1))$$



110

	Doc	Words	Class
Training	1	Chinese Beijing Chinese	c
	2	Chinese Chinese Shanghai	c
	3	Chinese Macao	c
	4	Tokyo Japan Chinese	j
Test	5	Chinese Chinese Chinese Tokyo Japan	?

**Priors:**

$$P(z=c) = 3/4$$

$$P(z=j) = 1/4$$

$\hat{P}(w | c) = \frac{\text{count}(w, c) + 1}{\text{count}(c) + |V|}$

**Conditional Probabilities:**

$$P(\text{Chinese} | c) = (5+1) / (8+6) = 6/14 = 3/7$$

$$P(\text{Tokyo} | c) = (0+1) / (8+6) = 1/14$$

$$P(\text{Japan} | c) = (0+1) / (8+6) = 1/14$$

$$P(\text{Chinese} | j) = (1+1) / (3+6) = 2/9$$

$$P(\text{Tokyo} | j) = (1+1) / (3+6) = 2/9$$

$$P(\text{Japan} | j) = (1+1) / (3+6) = 2/9$$

**Choosing a class:**

$$P(z=c | d5) \propto 3/4 * (3/7)^3 * 1/14 * 1/14 \approx 0.0003$$

$$P(z=j | d5) \propto 1/4 * (2/9)^3 * 2/9 * 2/9 \approx 0.0001$$

SOICT SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

111

## Support Vector Machine (SVM)



112

## Image Classification



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

→ cat



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

113

113

## Score function



→ class scores



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

114

114

## Linear Classifier

define a **score function**

**data (histogram)**

$$f(x_i, W, b) = Wx_i + b$$

↑

“weights”      “bias vector”

class scores      “parameters”



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

115

115

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

Convert image to histogram representation



0.2	-0.5	0.1	2.0
1.5	1.3	2.1	0.0
0	0.25	0.2	-0.3

W

56	1.1	-96.8
231	3.2	437.9
24	-1.2	61.95
2		f(x_i; W, b)

x<sub>i</sub>

b

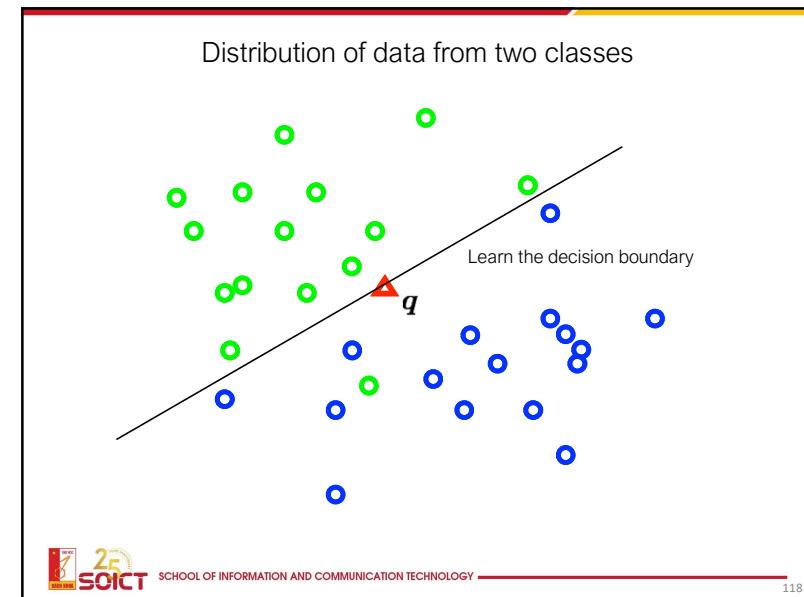
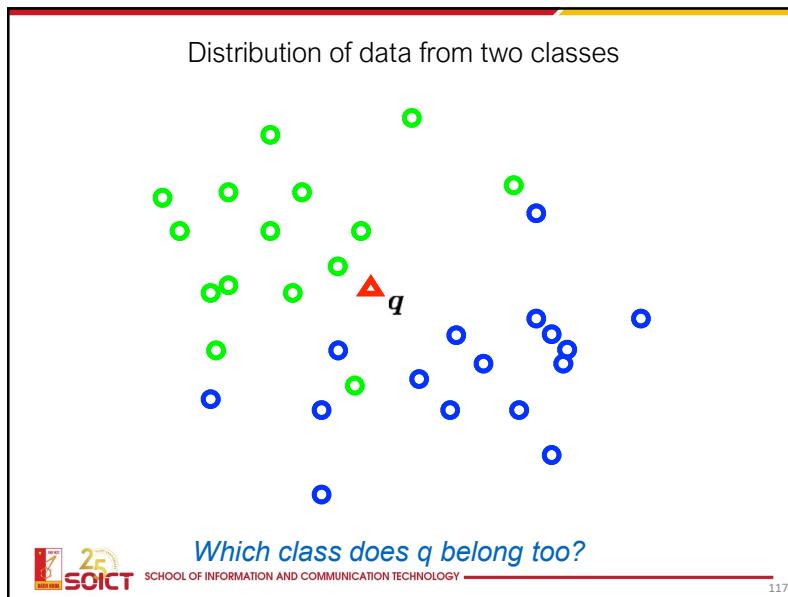
x<sub>i</sub>



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

116

116



First we need to understand hyperplanes...

SOICT SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

119

Hyperplanes (lines) in 2D

$$w_1x_1 + w_2x_2 + b = 0$$

a line can be written as dot product plus a bias

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

$$\mathbf{w} \in \mathbb{R}^2$$

another version, add a weight 1 and push the bias inside

$$\mathbf{w} \cdot \mathbf{x} = 0$$

$$\mathbf{w} \in \mathbb{R}^3$$

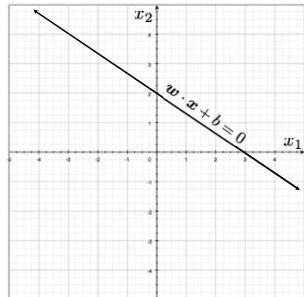
SOICT SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

120

## Hyperplanes (lines) in 2D

$$\mathbf{w} \cdot \mathbf{x} + b = 0 \quad (\text{offset/bias outside}) \quad \mathbf{w} \cdot \mathbf{x} = 0 \quad (\text{offset/bias inside})$$

$$w_1x_1 + w_2x_2 + b = 0$$



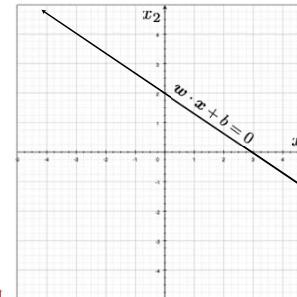
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

121

## Hyperplanes (lines) in 2D

$$\mathbf{w} \cdot \mathbf{x} + b = 0 \quad (\text{offset/bias outside}) \quad \mathbf{w} \cdot \mathbf{x} = 0 \quad (\text{offset/bias inside})$$

$$w_1x_1 + w_2x_2 + b = 0$$



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

122

**Important property:**  
Free to choose any normalization of  $w$

The line

$$w_1x_1 + w_2x_2 + b = 0$$

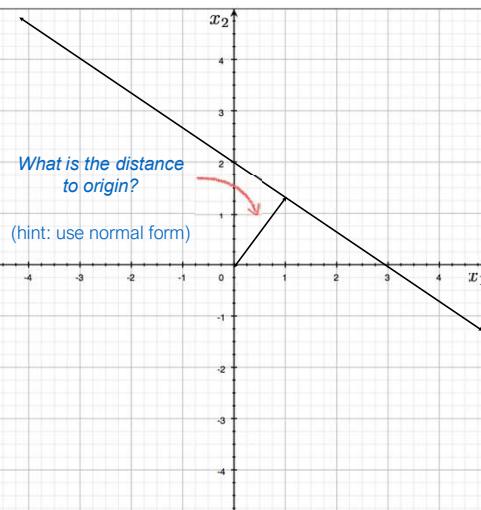
and the line

$$\lambda(w_1x_1 + w_2x_2 + b) = 0$$

define the same line

121

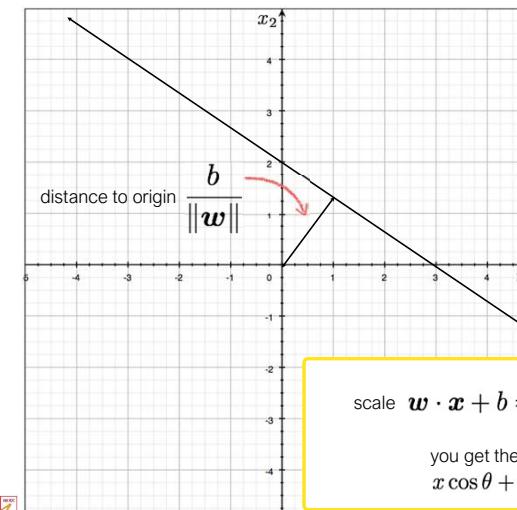
122



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

123

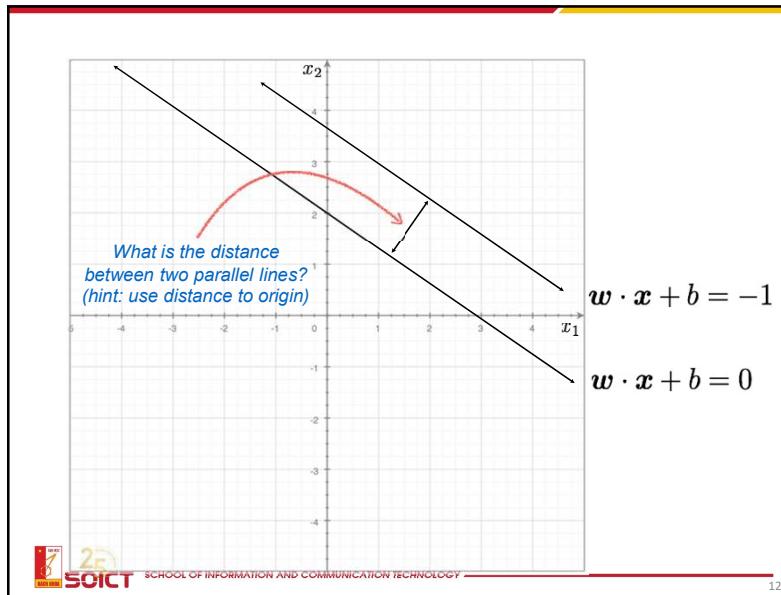


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

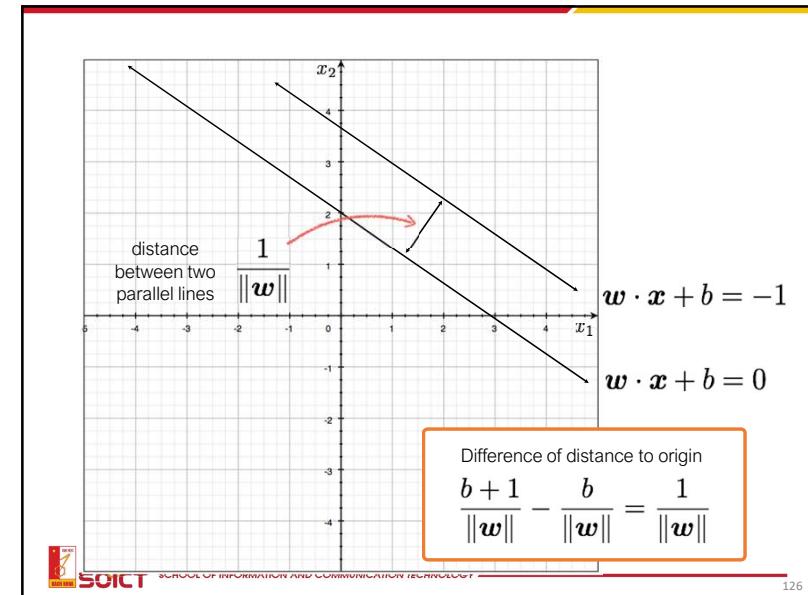
124

123

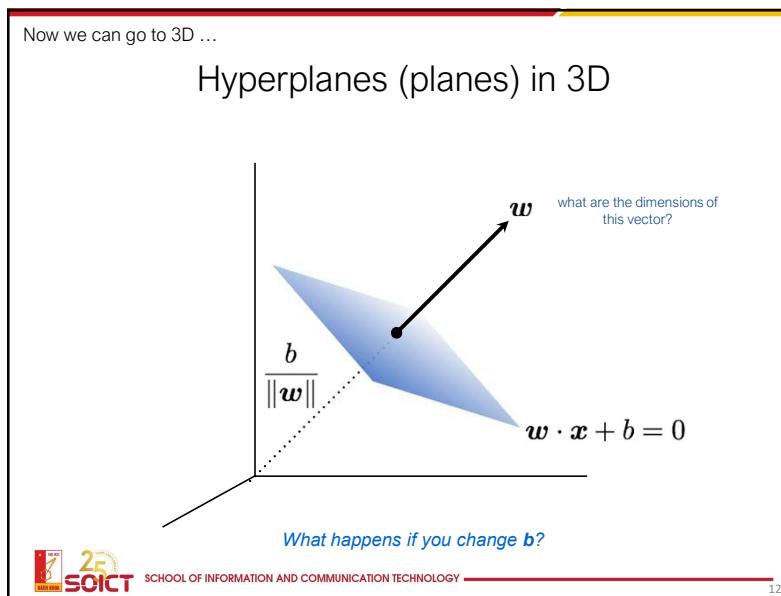
124



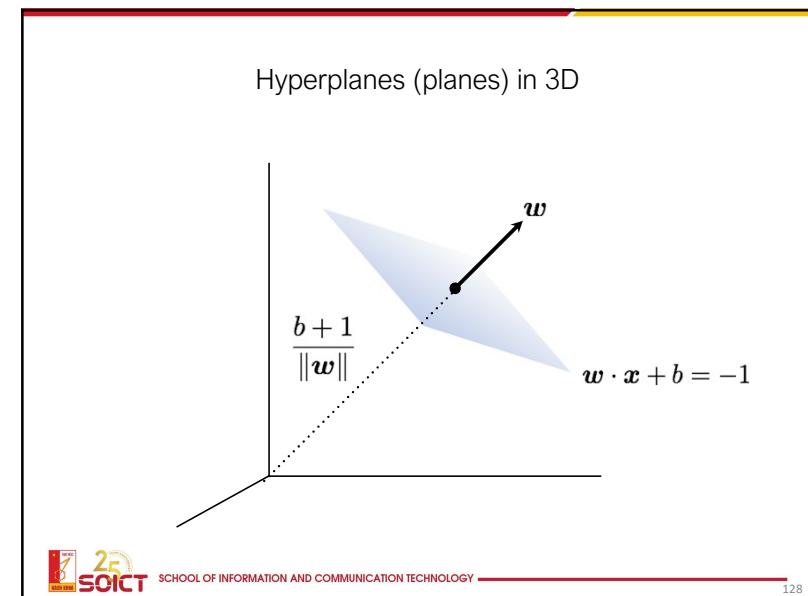
125



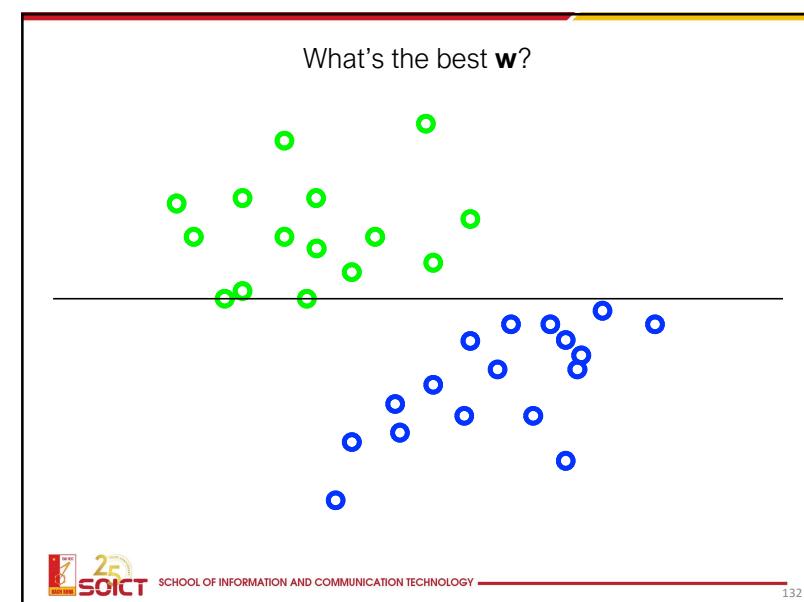
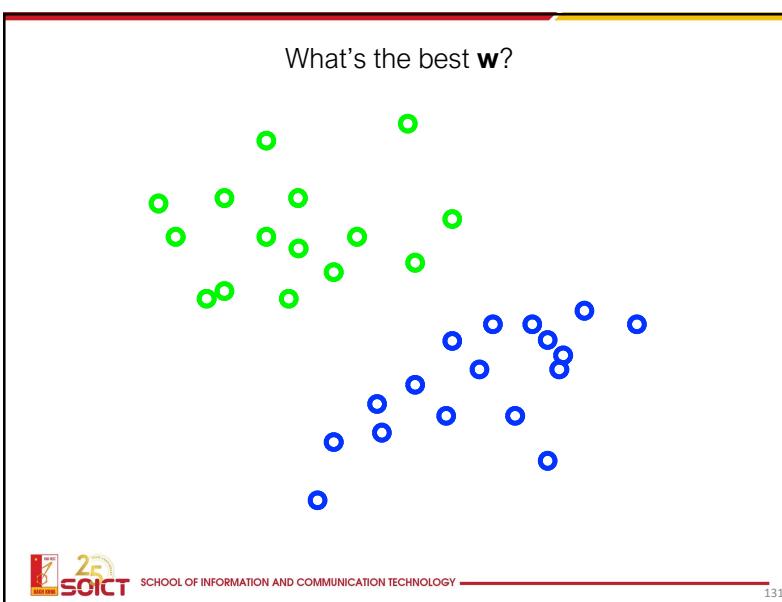
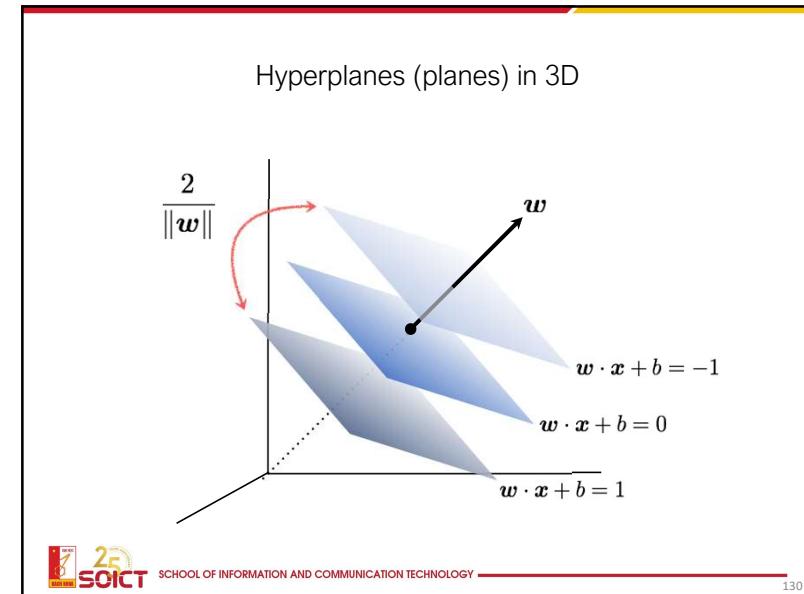
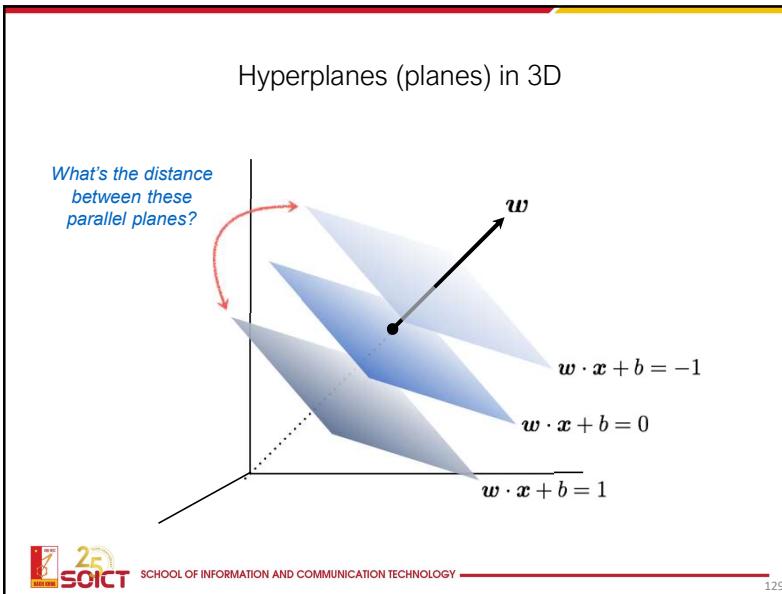
126

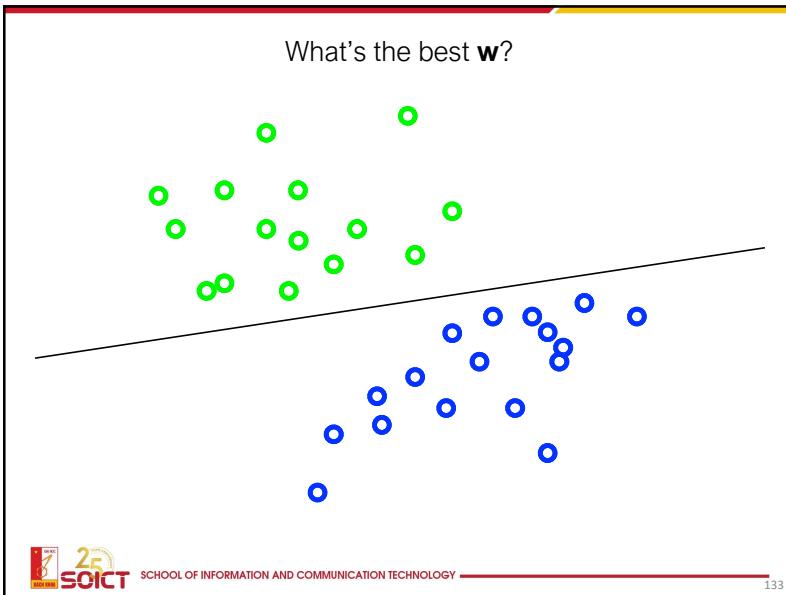


127

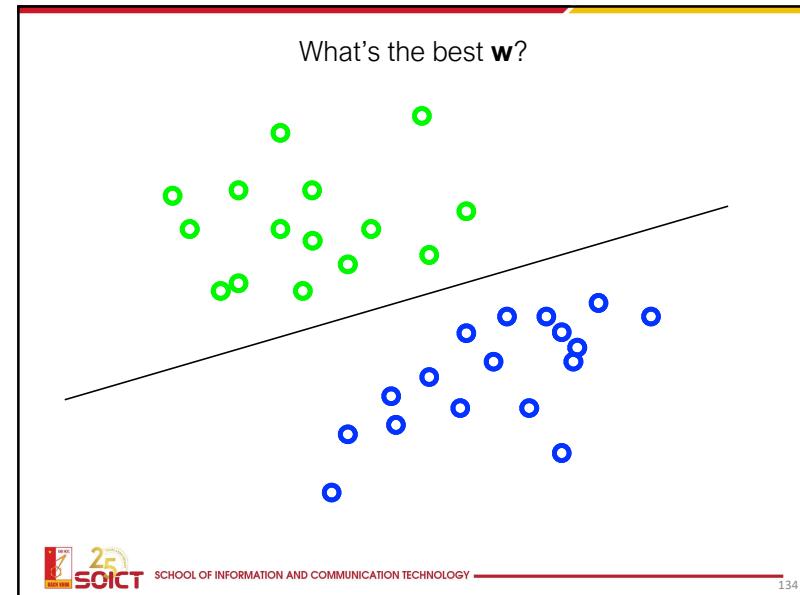


128

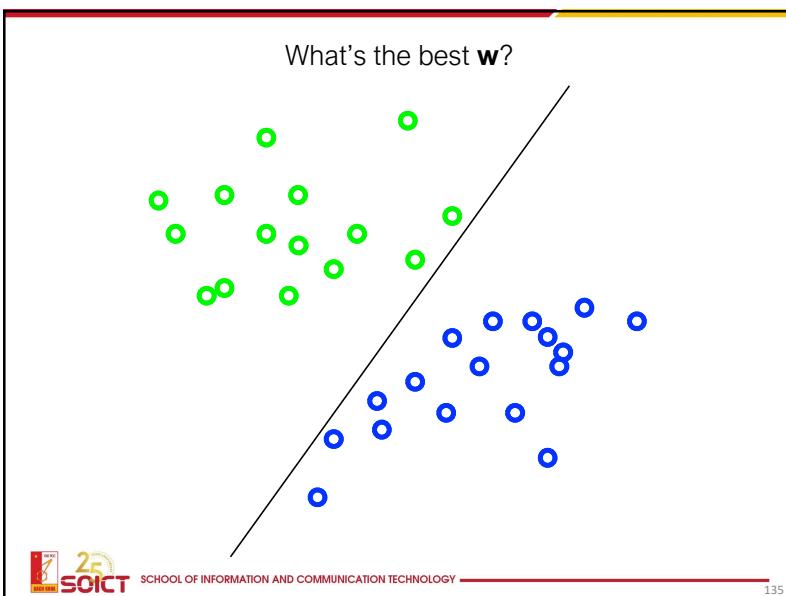




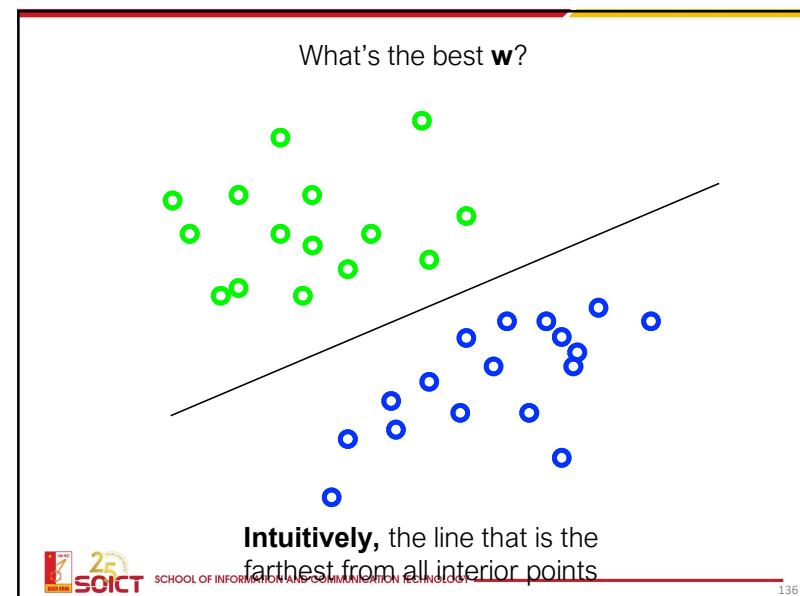
133



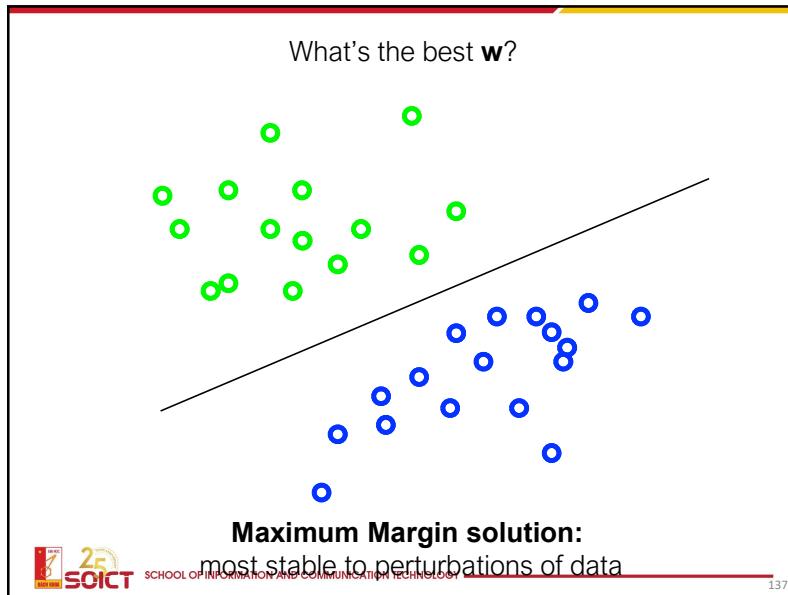
134



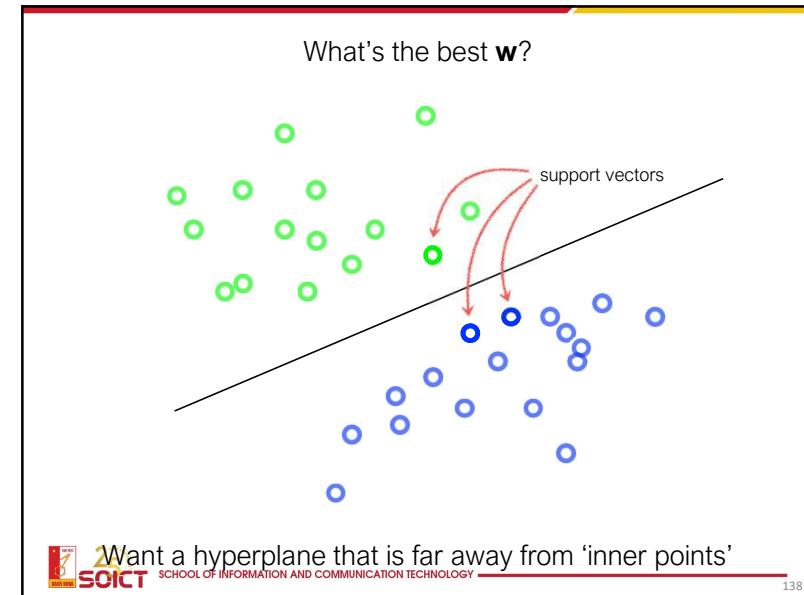
135



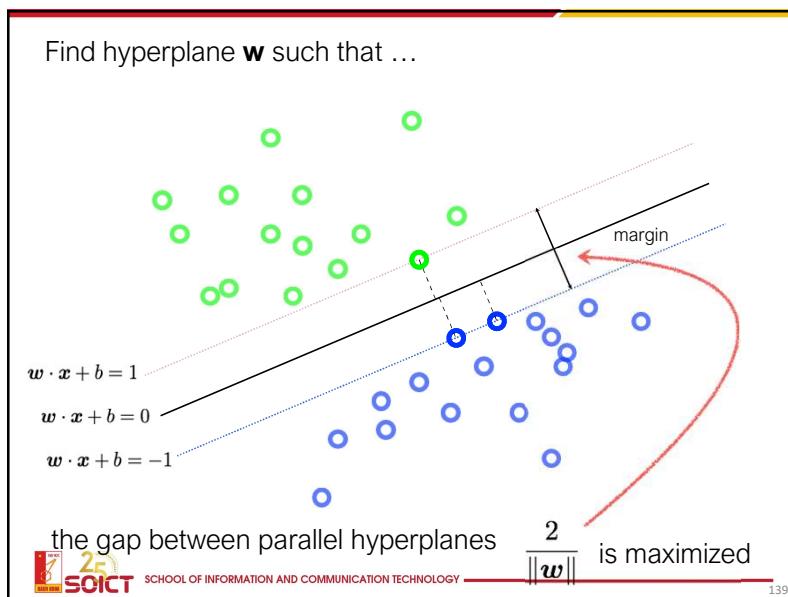
136



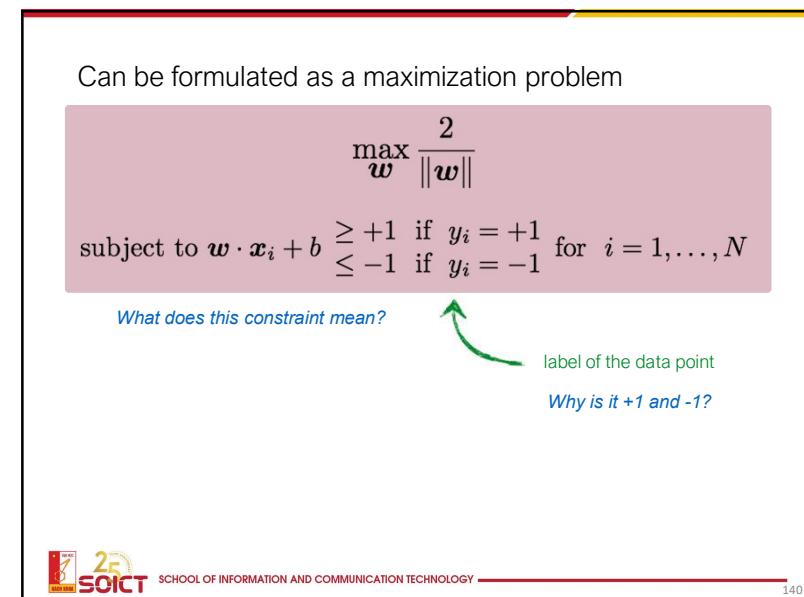
137



138



139



140

Can be formulated as a maximization problem

$$\max_{\mathbf{w}} \frac{2}{\|\mathbf{w}\|}$$

subject to  $\mathbf{w} \cdot \mathbf{x}_i + b \begin{cases} \geq +1 & \text{if } y_i = +1 \\ \leq -1 & \text{if } y_i = -1 \end{cases} \text{ for } i = 1, \dots, N$

Equivalently,

*Where did the 2 go?*

$$\min_{\mathbf{w}} \|\mathbf{w}\|$$

subject to  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \text{ for } i = 1, \dots, N$

*What happened to the labels?*



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

141

'Primal formulation' of a linear SVM

$$\min_{\mathbf{w}} \|\mathbf{w}\|$$

Objective Function

subject to  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \text{ for } i = 1, \dots, N$

Constraints

This is a convex quadratic programming (QP) problem  
(a unique solution exists)



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

142

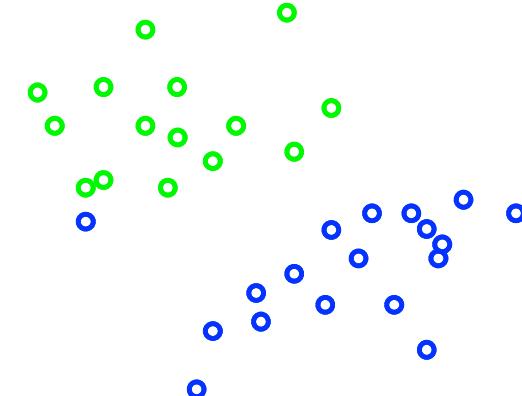
SVM: 'soft' margin



SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

143

What's the best  $\mathbf{w}$ ?

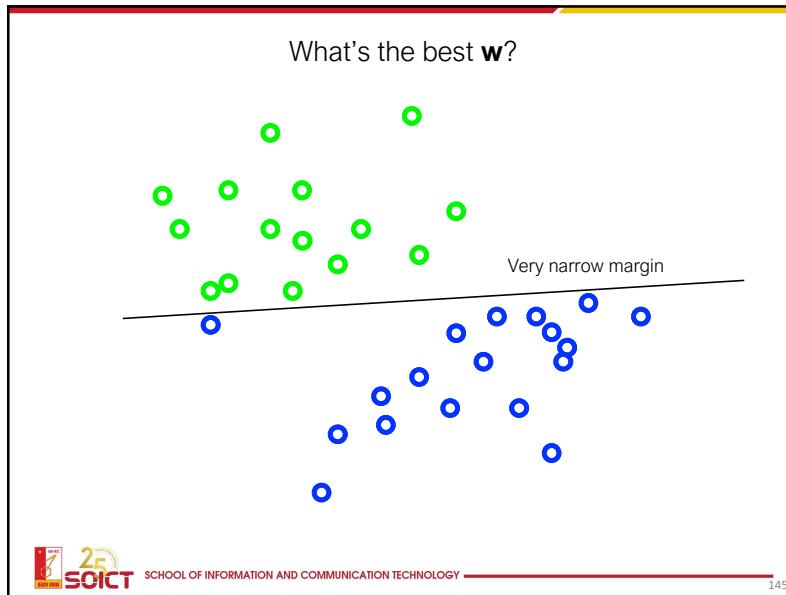


SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

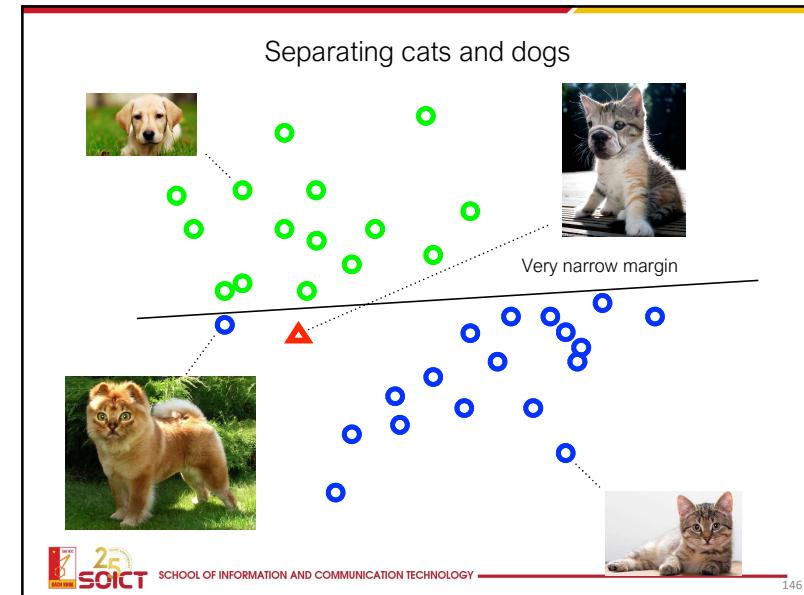
144

143

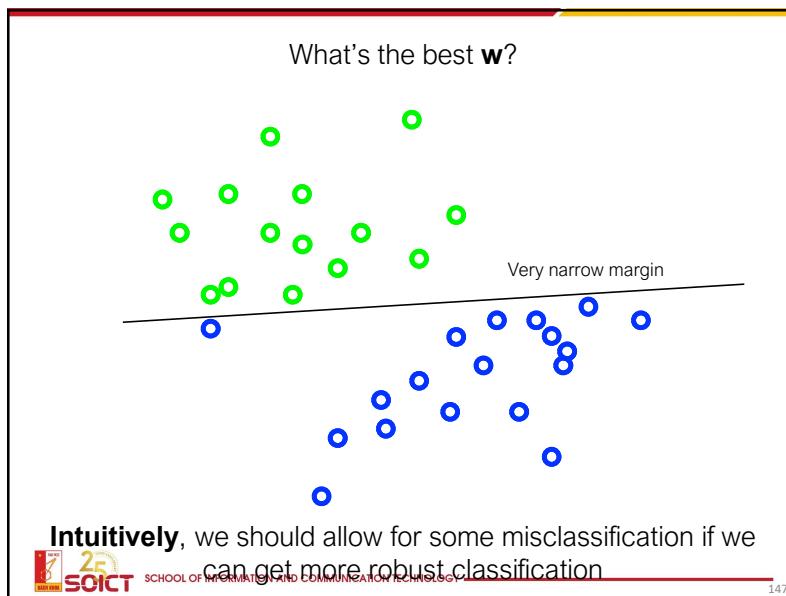
144



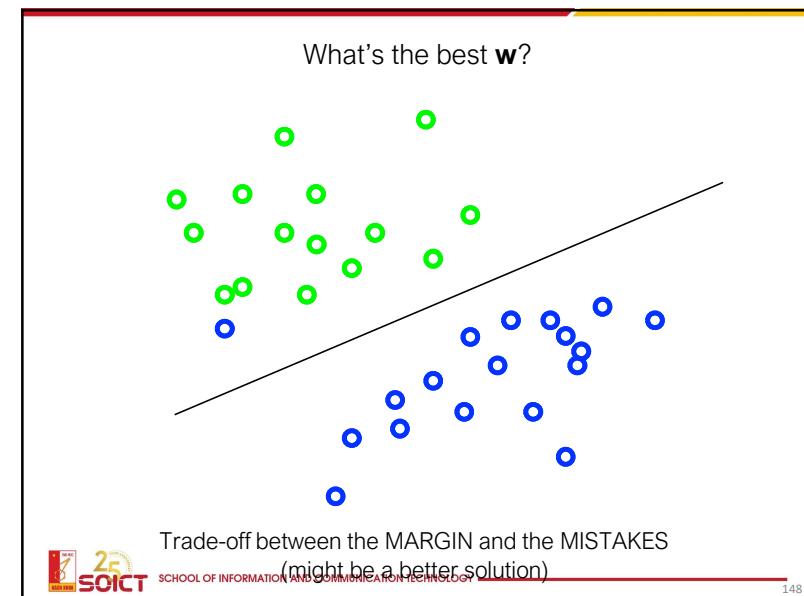
145



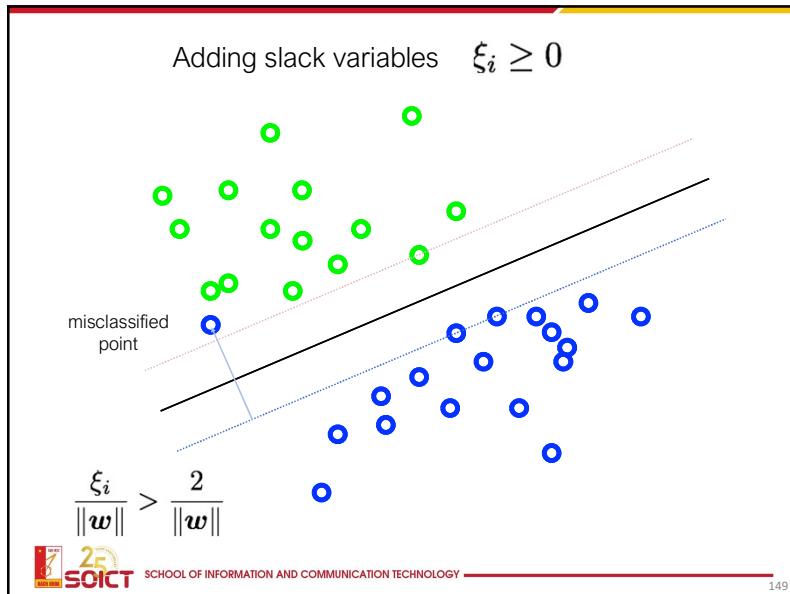
146



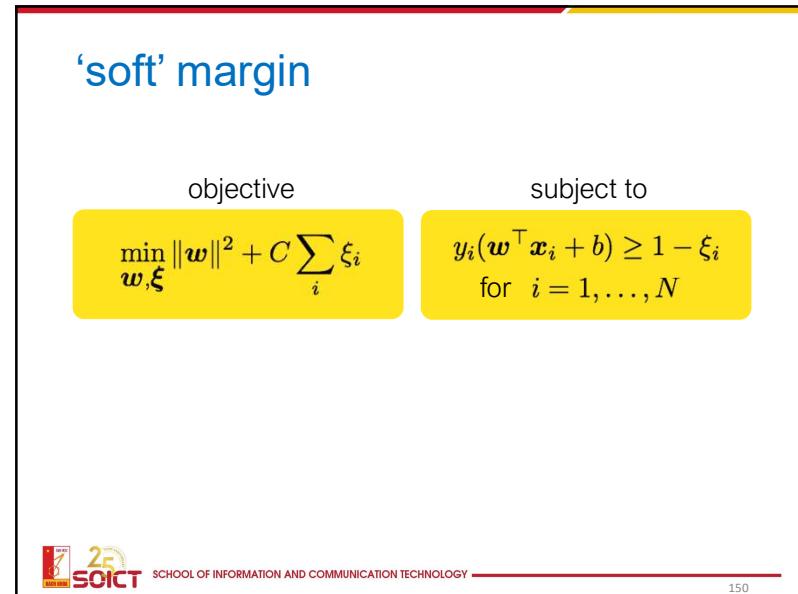
147



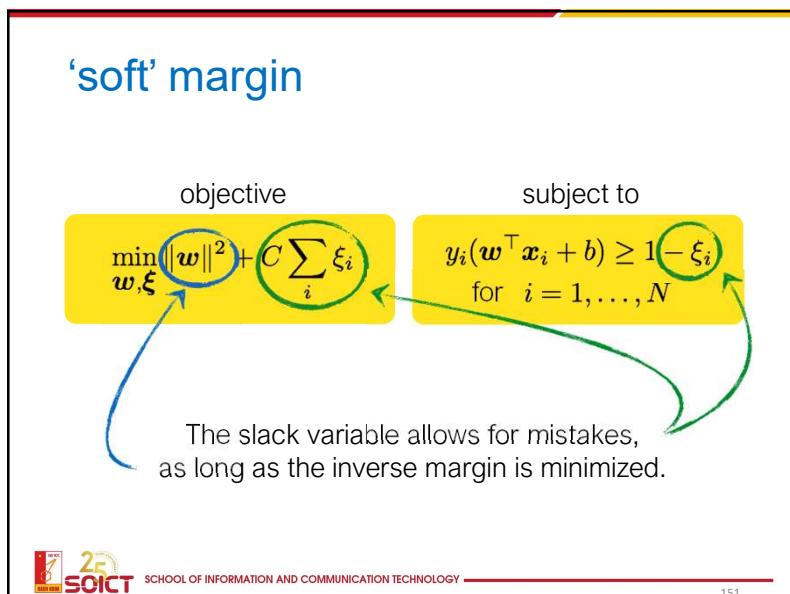
148



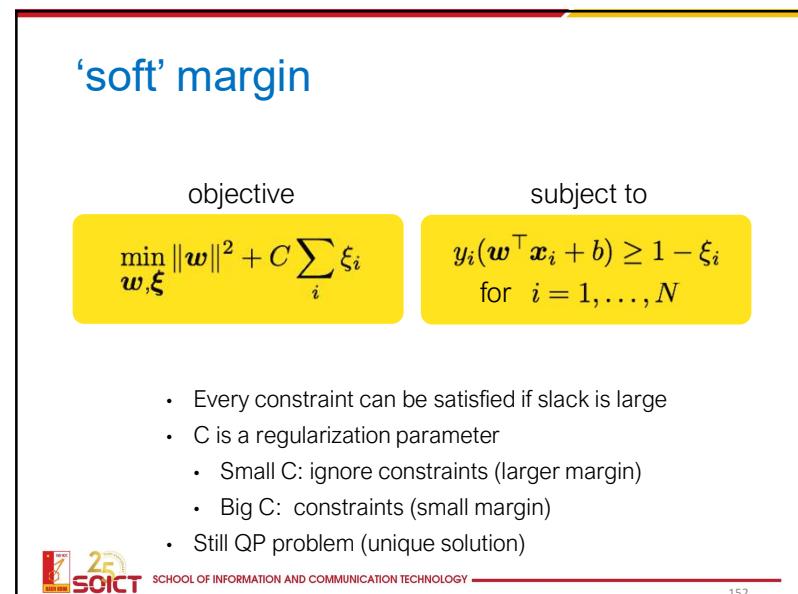
149



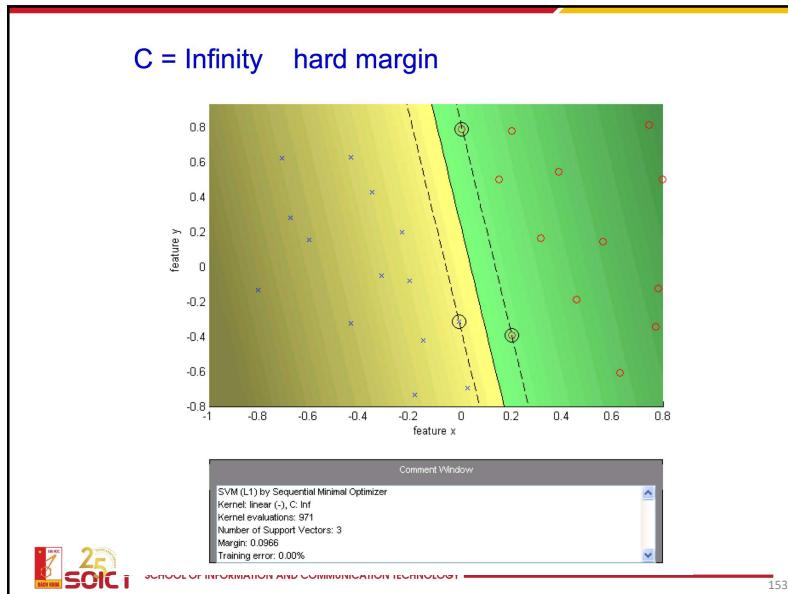
150



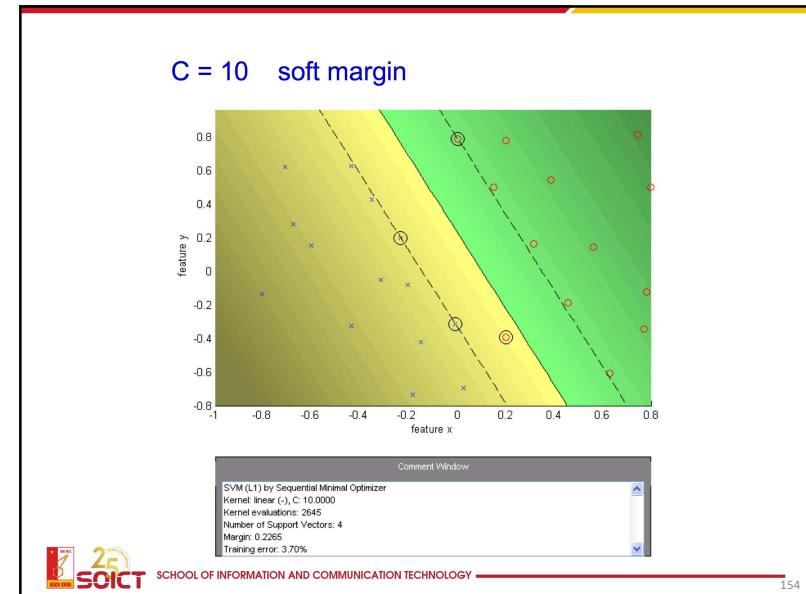
151



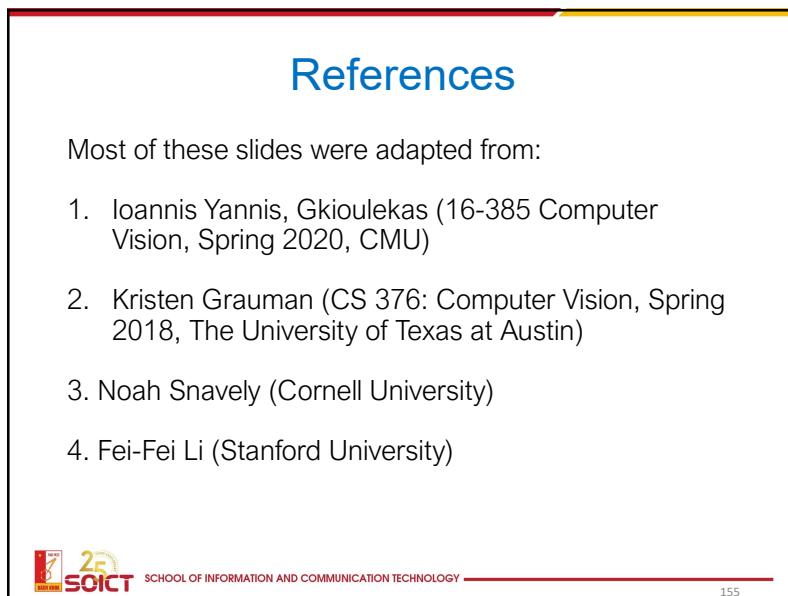
152



153



154



155



156