



Computer Vision I:

Jingya Wang

Email: wangjingya@shanghaitech.edu.cn



Introduction to object recognition



Overview

- Basic recognition tasks
- A statistical learning approach
- Traditional or “shallow” recognition pipeline
 - Bags of features
 - Classifiers
- Next time: neural networks and “deep” recognition pipeline

Common recognition tasks

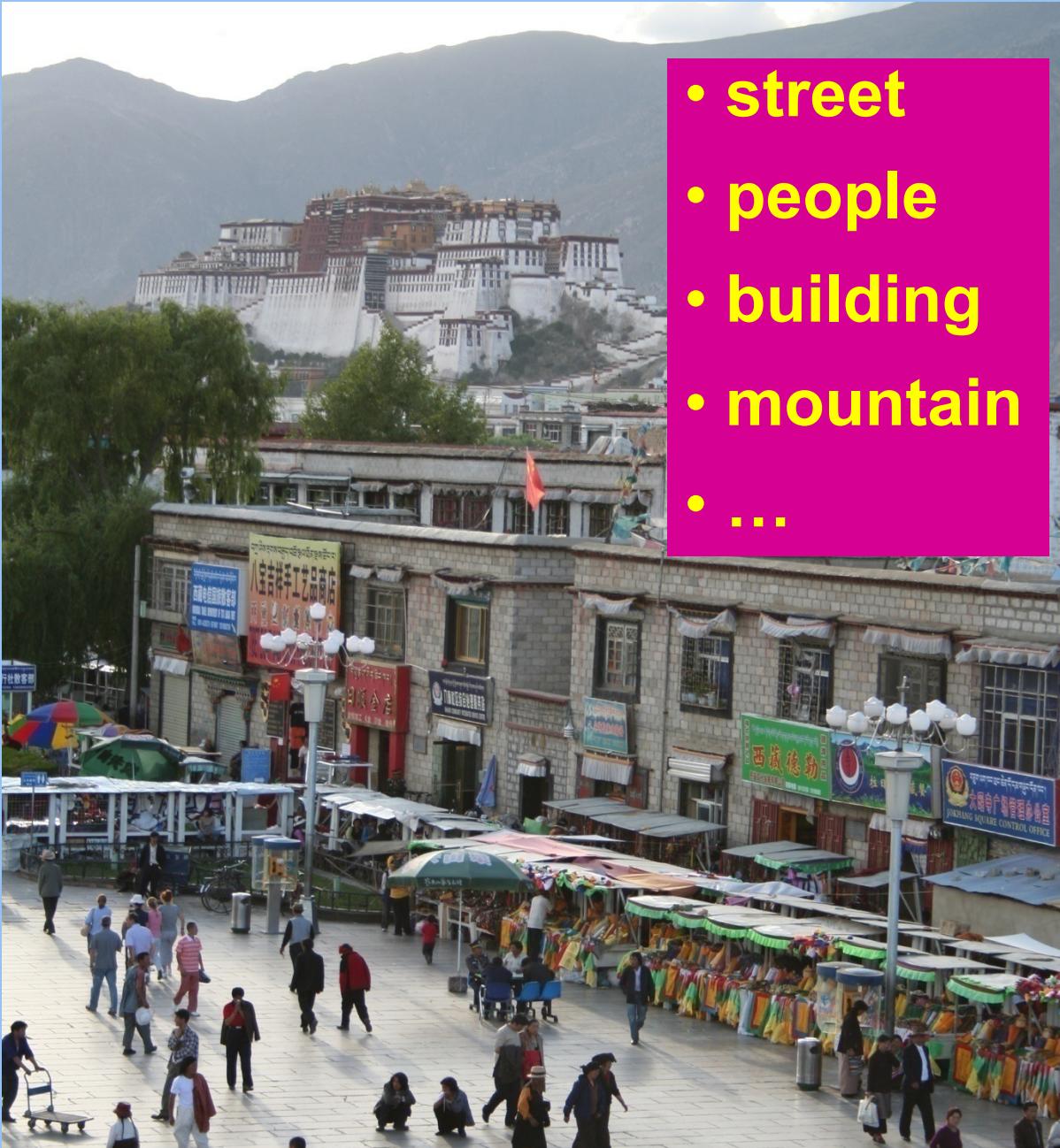


Image classification



- outdoor/indoor
- city/forest/factory/etc.

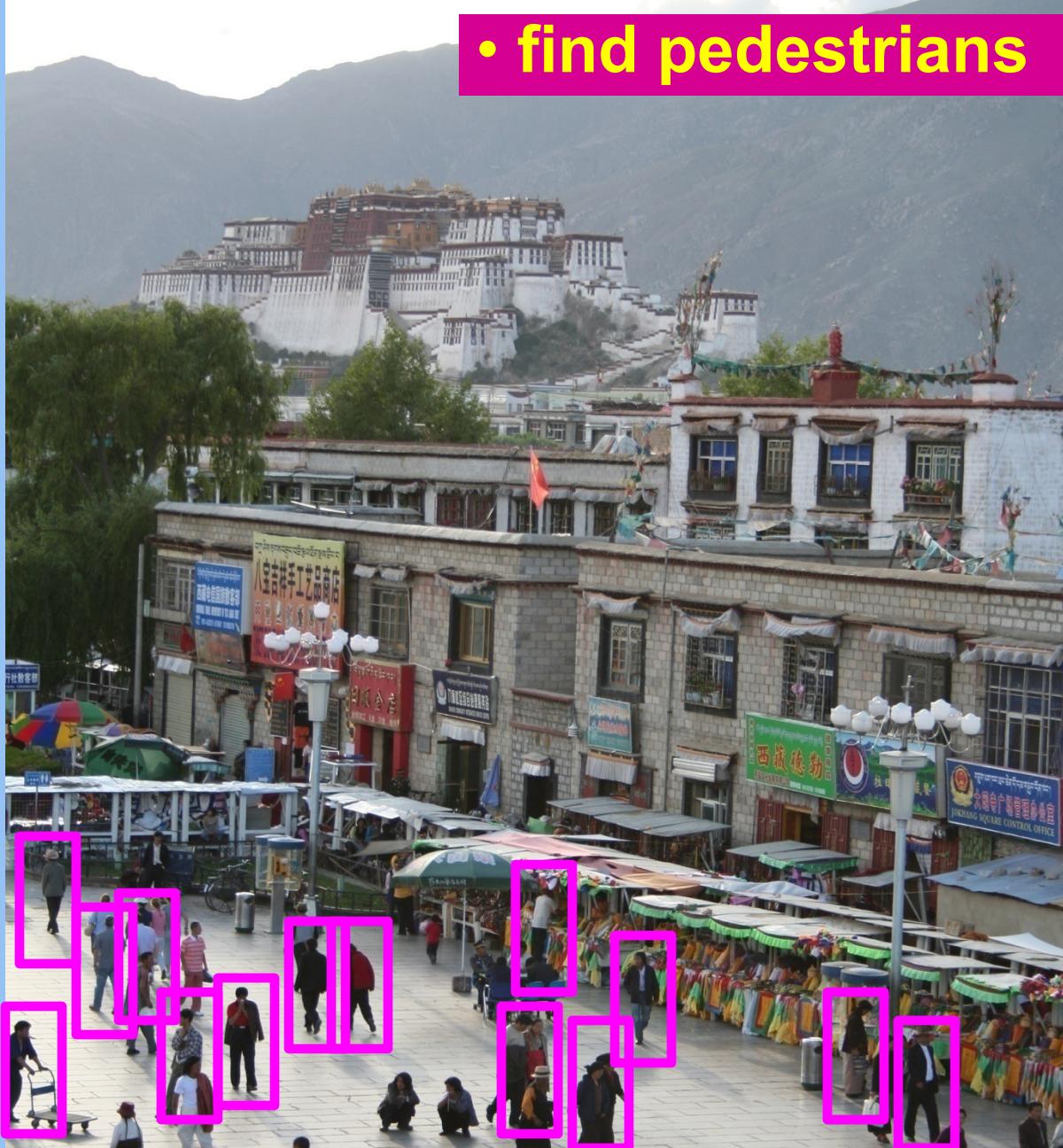
Image tagging



- street
- people
- building
- mountain
- ...

Object detection

- find pedestrians



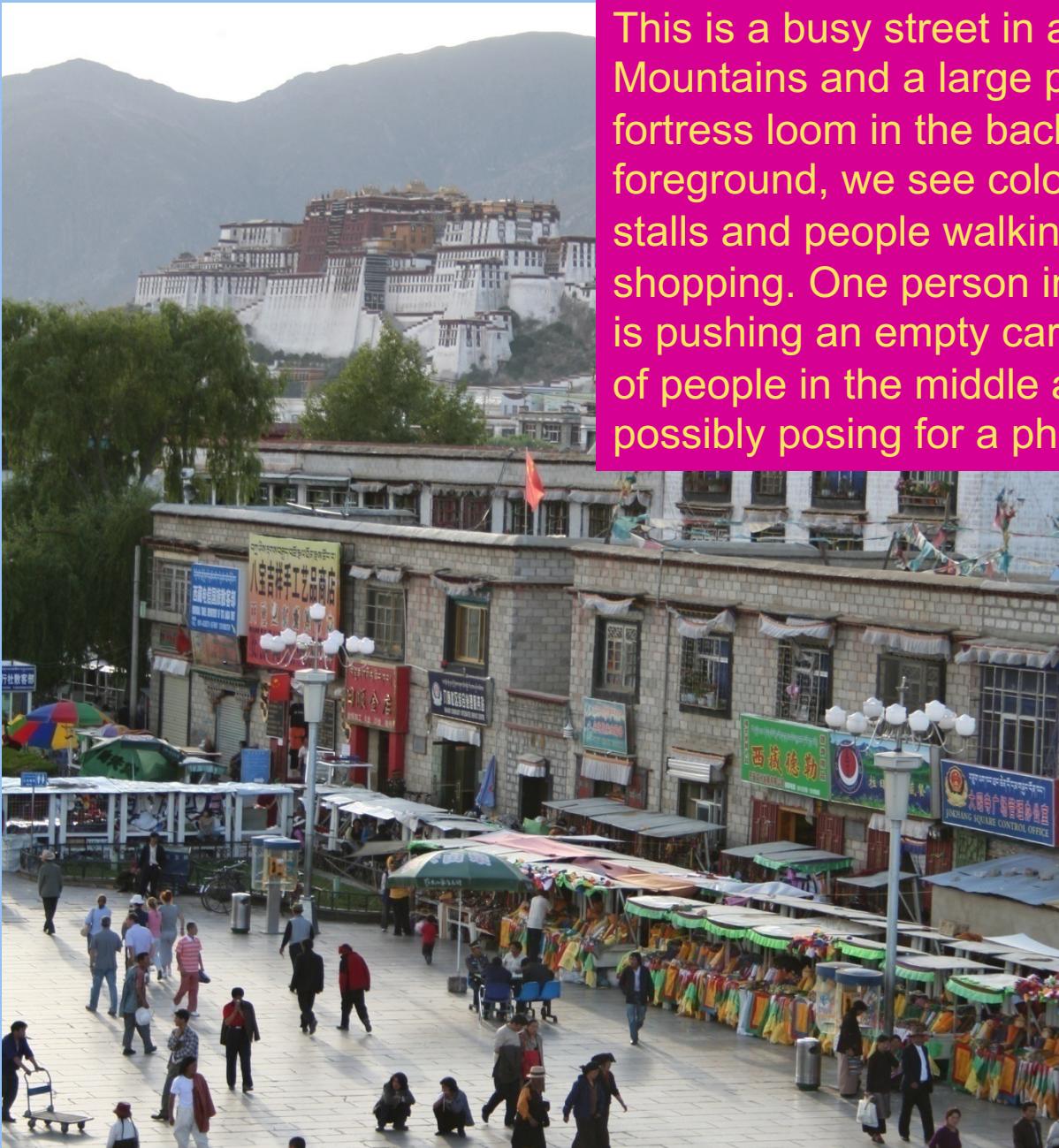
Activity recognition



Image parsing



Image description



This is a busy street in an Asian city. Mountains and a large palace or fortress loom in the background. In the foreground, we see colorful souvenir stalls and people walking around and shopping. One person in the lower left is pushing an empty cart, and a couple of people in the middle are sitting, possibly posing for a photograph.

Image Classification: A core task in Computer Vision

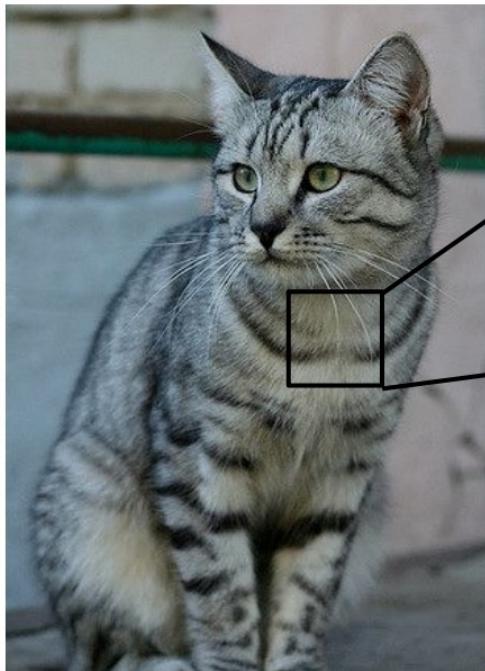


This image by [Nikita](#) is
licensed under [CC-BY 2.0](#)



cat

The Problem: Semantic Gap



```
[ [105 112 108 111 104 99 106 99 96 103 112 119 104 97 93 87]
[ 91 98 102 106 104 79 98 103 99 105 123 136 110 105 94 85]
[ 76 85 90 105 128 105 87 96 95 99 115 112 106 103 99 85]
[ 99 81 81 93 120 131 127 100 95 98 102 99 96 93 101 94]
[114 108 85 55 55 69 64 54 64 87 112 129 98 74 84 91]
[133 137 147 103 65 81 80 65 52 54 74 84 102 93 85 82]
[128 137 144 140 109 95 86 70 62 65 63 63 60 73 86 101]
[125 133 148 137 119 121 117 94 65 79 80 65 54 64 72 98]
[127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]
[115 114 109 123 150 148 131 118 113 109 100 92 74 65 72 78]
[ 89 93 90 97 108 147 131 118 113 114 113 109 106 95 77 80]
[ 63 77 86 81 77 79 102 123 117 115 117 125 125 130 115 87]
[ 62 65 82 89 78 71 80 101 124 126 119 101 107 114 131 119]
[ 63 65 75 88 89 71 62 81 120 138 135 105 81 98 110 118]
[ 87 65 71 87 106 95 69 45 76 130 126 107 92 94 105 112]
[118 97 82 86 117 123 116 66 41 51 95 93 89 95 102 107]
[164 146 112 80 82 120 124 104 76 48 45 66 88 101 102 109]
[157 170 157 120 93 86 114 132 112 97 69 55 70 82 99 94]
[130 128 134 161 139 100 109 118 121 134 114 87 65 53 69 86]
[128 112 96 117 150 144 126 115 104 107 102 93 87 81 72 79]
[123 107 96 86 83 112 153 149 122 109 104 75 80 107 112 99]
[122 121 102 80 82 86 94 117 145 148 153 102 58 78 92 107]
[122 164 148 103 71 56 78 83 93 103 119 139 102 61 69 84]]
```

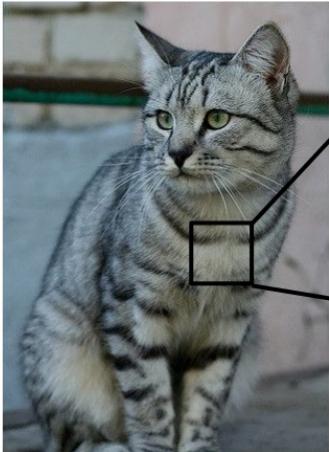
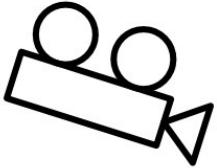
What the computer sees

An image is a tensor of integers between [0, 255]:

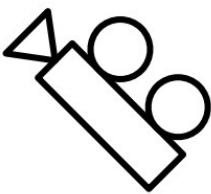
e.g. 800 x 600 x 3
(3 channels RGB)

This image by [Nikita](#) is
licensed under [CC-BY 2.0](#)

Challenges: Viewpoint variation



```
[1185 112 100 111 184 99 186 99 96 183 112 110 184 97 93 97]
[ 91 98 102 106 104 79 98 103 99 105 123 136 110 185 94 95]
[ 76 85 90 105 128 105 87 96 95 99 115 112 106 103 94 95]
[ 99 81 89 128 131 127 180 95 98 102 99 96 93 103 94]
[106 91 61 64 69 91 88 85 101 107 109 98 75 84 96 95]
[114 100 85 55 69 88 54 64 123 126 128 100 70 84 86 95]
[137 117 114 85 91 88 54 62 54 74 84 103 93 96 82]
[128 137 144 148 189 95 86 78 62 65 63 60 73 86 181]
[125 133 141 137 119 121 117 94 65 79 88 65 54 64 72 98]
[127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]
[115 114 100 123 150 140 131 118 113 109 100 92 74 65 72 78]
[ 90 85 80 88 100 100 100 100 100 100 100 100 100 100 100 100]
[ 63 77 86 81 77 79 102 123 117 115 117 125 125 130 115 87]
[ 62 65 82 89 78 71 88 101 124 126 119 101 107 114 131 119]
[ 87 65 73 87 106 95 69 45 76 138 126 107 92 94 105 112]
[ 63 65 75 88 89 71 62 81 120 138 135 103 81 98 118 118]
[ 87 65 73 87 106 95 69 45 76 138 126 107 92 94 105 112]
[118 97 85 86 117 123 116 66 41 51 95 93 89 95 102 107]
[110 104 102 103 100 100 100 100 100 100 100 100 100 100 100 100]
[157 178 157 128 93 86 114 132 112 97 69 55 70 82 99 94]
[130 128 134 161 139 100 109 118 121 134 114 87 65 53 69 86]
[128 112 96 117 158 144 128 115 104 107 102 93 87 81 72 79]
[123 107 96 86 83 112 153 149 122 109 104 75 80 107 112 99]
[122 121 102 88 82 86 94 117 145 148 153 103 58 78 92 107]
[122 164 148 103 71 56 78 83 93 103 119 139 102 61 69 841]
```



All pixels change when
the camera moves!

Challenges: Background Clutter



[This image](#) is [CC0 1.0](#) public domain



[This image](#) is [CC0 1.0](#) public domain

Challenges: Illumination



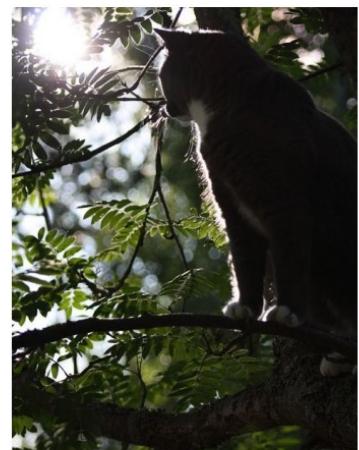
[This image](#) is CC0 1.0 public domain



[This image](#) is CC0 1.0 public domain



[This image](#) is CC0 1.0 public domain



[This image](#) is CC0 1.0 public domain

Challenges: Occlusion



[This image](#) is [CC0 1.0](#) public domain

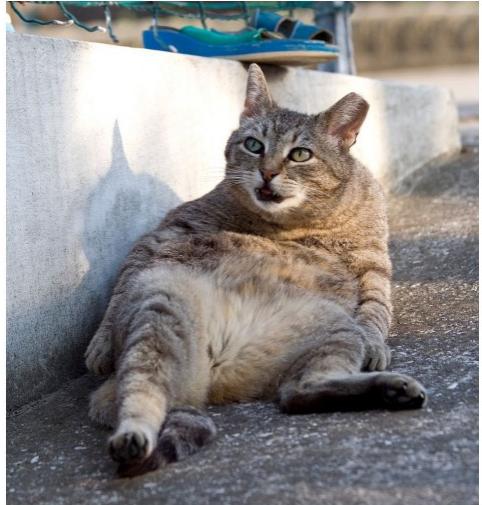


[This image](#) is [CC0 1.0](#) public domain



[This image](#) by [jonsson](#) is licensed
under [CC-BY 2.0](#)

Challenges: Deformation



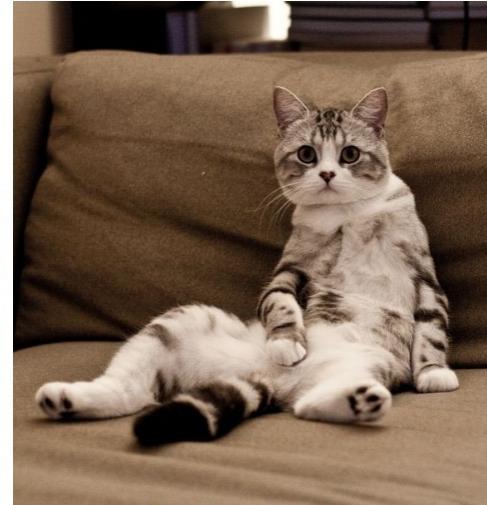
[This image](#) by [Umberto Salvagnin](#)
is licensed under [CC-BY 2.0](#)



[This image](#) by [Umberto Salvagnin](#)
is licensed under [CC-BY 2.0](#)



[This image](#) by [sare bear](#) is
licensed under [CC-BY 2.0](#)



[This image](#) by [Tom Thai](#) is
licensed under [CC-BY 2.0](#)

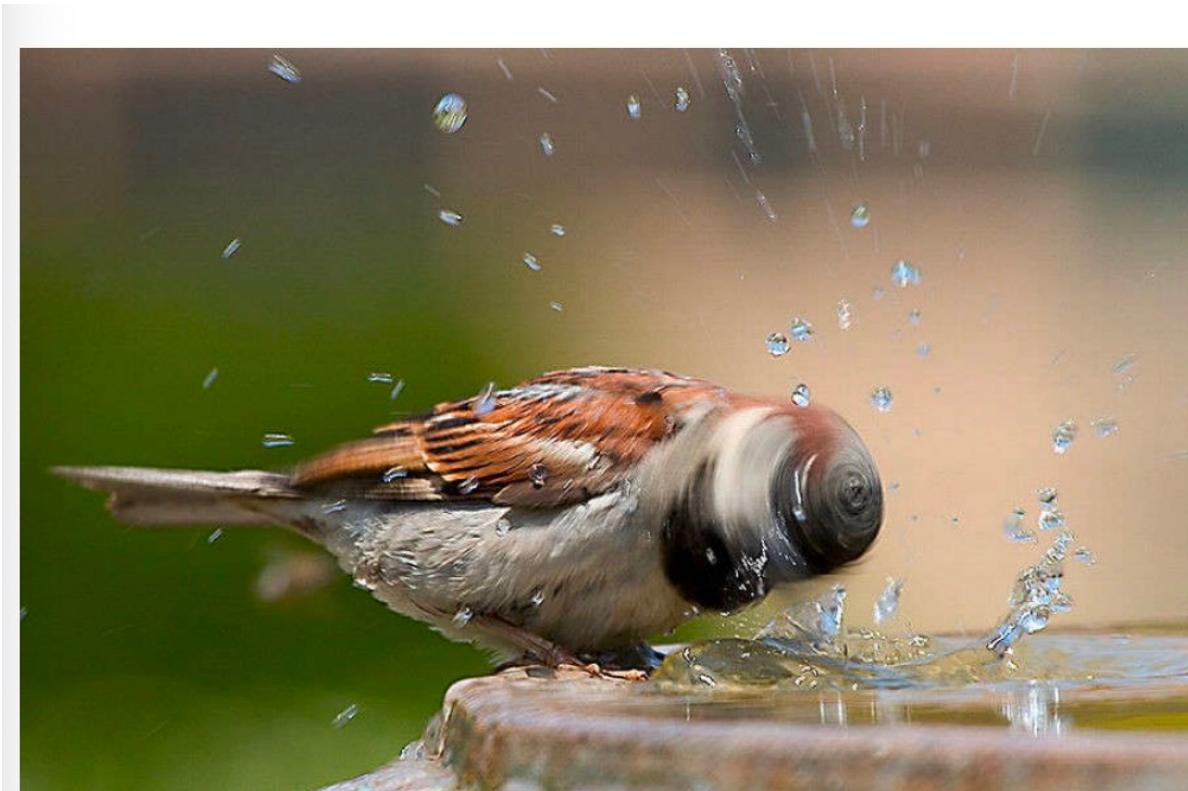
Challenges: scale

and small things
from Apple.
(Actual size)



Source: S. Lazebnik

Challenges: motion blur



Source: S. Lazebnik

Challenges: Intraclass variation



This image is [CC0 1.0](#) public domain

Challenges: Intraclass variation

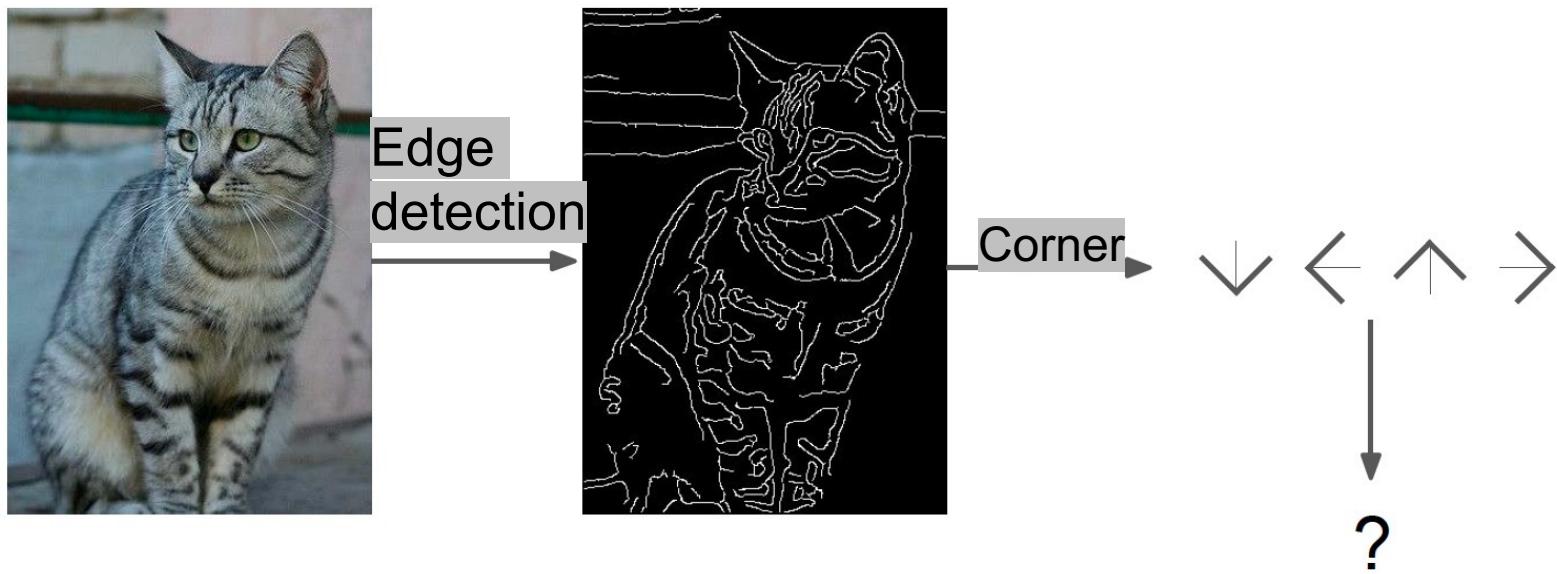


Source: Fei-Fei

Solution

- Rule based methods? Hard-coding

```
def classify_image ( image ):  
    # Do something magical here  
    return class_label
```



Source: Fei-Fei Li

Machine Learning: Data-Driven Approach

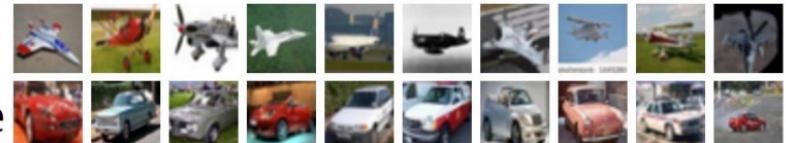
1. Collect a dataset of images and labels
2. Use Machine Learning algorithms to train a classifier
3. Evaluate the classifier on new images

Example training set

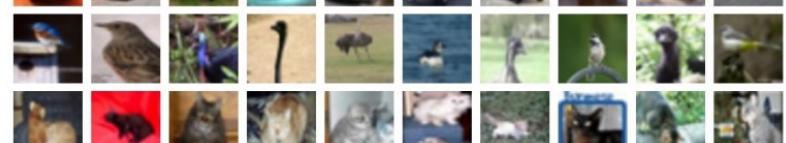
```
def train(images, labels):
    # Machine learning!
    return model
```

```
def predict(model, test_images):
    # Use model to predict labels
    return test_labels
```

airplane



automobile



bird



cat

deer

Image classification



The statistical learning framework

- Apply a prediction function to a feature representation of the image to get the desired output:

$$f(\text{apple}) = \text{"apple"}$$
$$f(\text{tomato}) = \text{"tomato"}$$
$$f(\text{cow}) = \text{"cow"}$$

The statistical learning framework

$$y = f(\mathbf{x})$$

A diagram illustrating the components of the prediction function $y = f(\mathbf{x})$. The equation is written in blue. Three red arrows point upwards from labels to the corresponding parts of the equation: a vertical arrow points to the output variable y , another vertical arrow points to the prediction function f , and a diagonal arrow points to one of the input features \mathbf{x} .

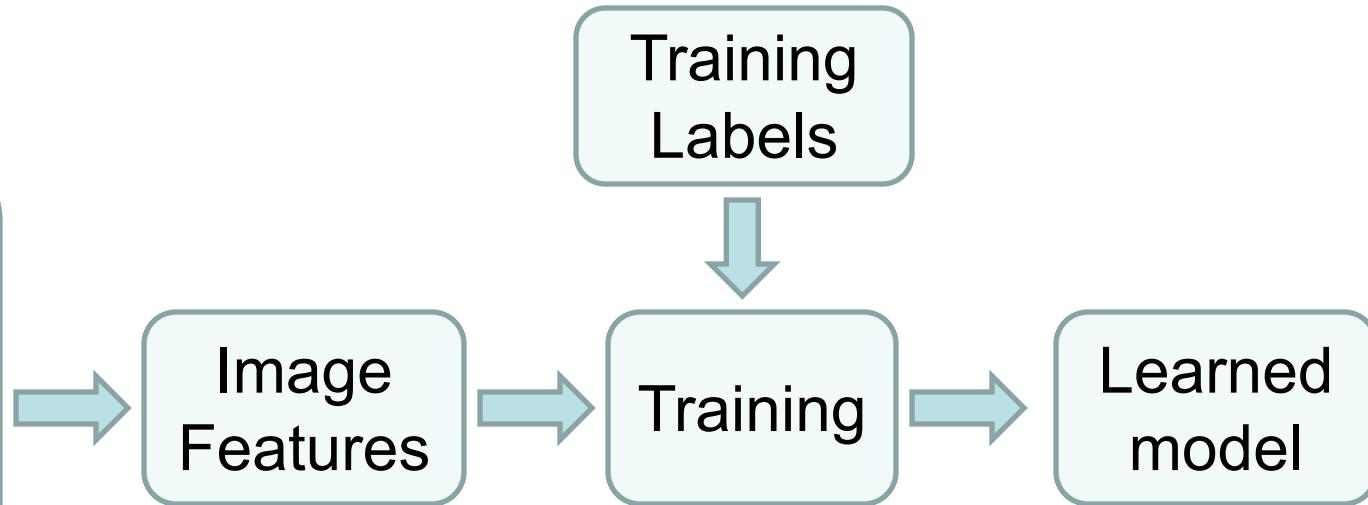
output prediction function Image feature

- **Training:** given a *training set* of labeled examples $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, estimate the prediction function f by minimizing the prediction error on the training set
- **Testing:** apply f to a never before seen *test example* \mathbf{x} and output the predicted value $y = f(\mathbf{x})$

Steps

Training

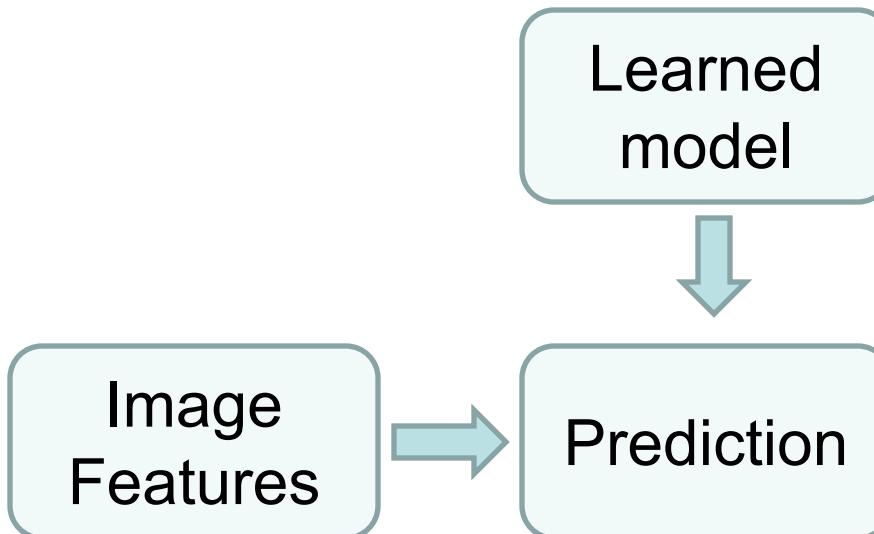
Training Images



Testing

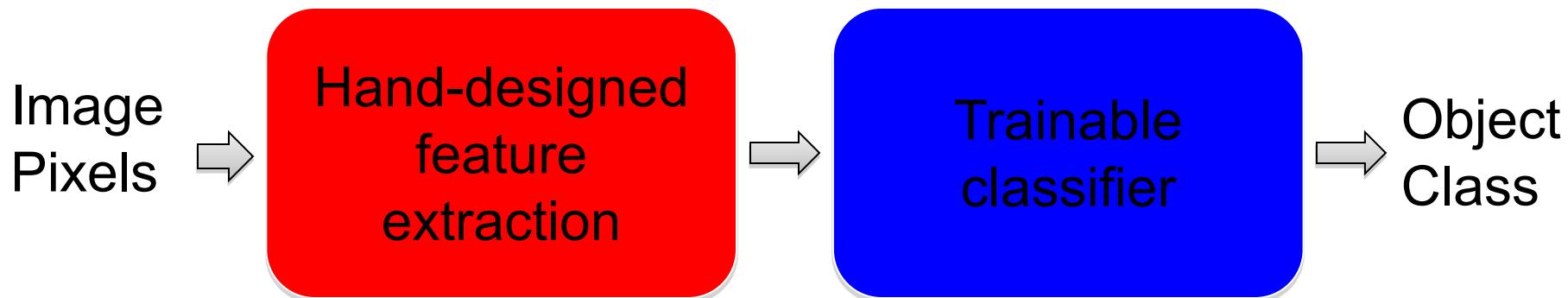


Test Image



Traditional recognition pipeline

feature 不可学习, fixed!!



不是 data-driven

- Features are not learned
- Trainable classifier is often generic (e.g. SVM)

Bags of features: Motivation

- Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)

2007-01-23: State of the Union Address

George W. Bush (2001-)

abandon accountable affordable afghanistan africa aided ally anbar armed army **baghdad** bless **challenges** chamber chaos
choices civilians coalition **commanders** **commitment** confident confront congressman constitution corps debates deduction
deficit deliver **democratic** deploy dikembe diplomacy disruptions earmarks **economy** einstein **elections** eliminates
expand **extremists** failing faithful families **freedom** fuel **funding** god haven ideology immigration impose

insurgents iran **iraq** islam julie lebanon love madam marine math medicare moderation neighborhoods nuclear offensive
palestinian payroll province pursuing **qaeda** radical **regimes** resolve retreat rieman sacrifices science sectarian senate

september **shia** stays strength students succeed sunni **tax** territories **terrorists** threats uphold victory
violence violent **War** washington weapons wesley

Bags of features: Motivation

- Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)



TF-IDF

↑ 重要性↑

词频

- **TF: Term Frequency**, which measures how frequently a term occurs in a document. Since every document is different in length, it is possible that a term would appear much more times in long documents than shorter ones. Thus, the term frequency is often divided by the document length (aka. the total number of terms in the document) as a way of normalization:

$$TF(t) = (\text{Number of times term } t \text{ appears in a document}) / (\text{Total number of terms in the document}).$$

逆文档词频

过滤不重要的高频词

- **IDF: Inverse Document Frequency**, which measures how important a term is. While computing TF, all terms are considered equally important. However it is known that certain terms, such as "is", "of", and "that", may appear a lot of times but have little importance. Thus we need to weigh down the frequent terms while scale up the rare ones, by computing the following:

$$IDF(t) = \log(\text{Total number of documents} / \text{Number of documents with term } t \text{ in it}).$$

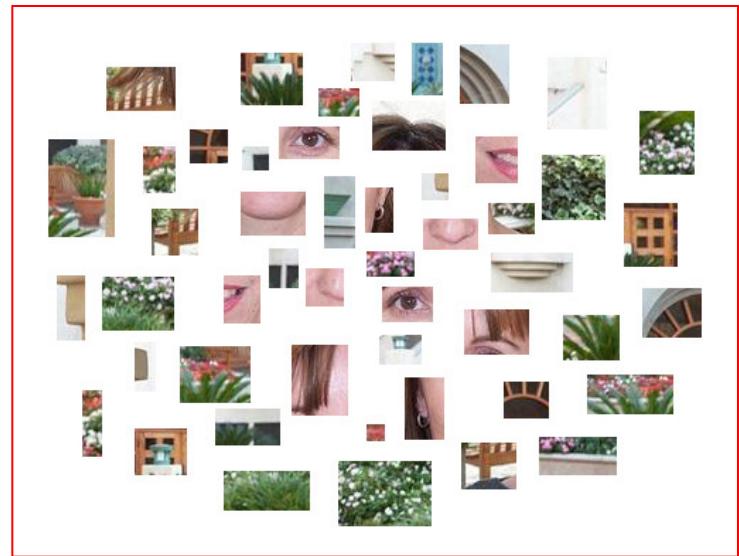
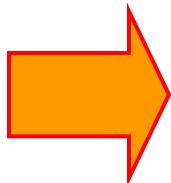
Example

Consider a document containing 100 words wherein the word *cat* appears 3 times. The term frequency (i.e., tf) for *cat* is ? 3/100

Now, assume we have 10 million documents and the word *cat* appears in one thousand of these. Then, the inverse document frequency (i.e., idf) is calculated as ? $\log_{10}\left(\frac{10,000,000}{10^3}\right) = 4$

$$\text{TF-IDF} = 0.03 \times 4 = 0.12$$

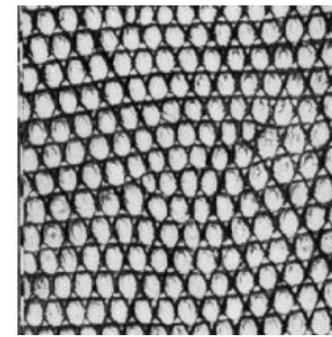
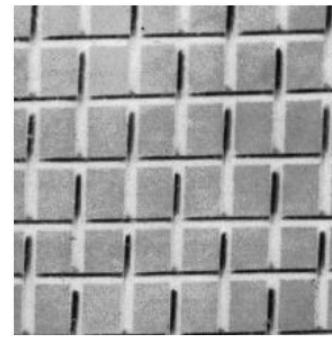
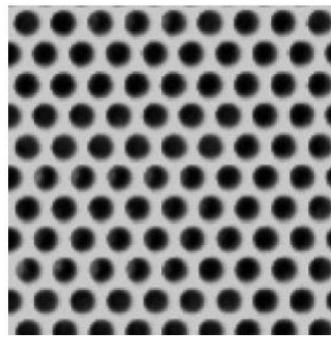
Bags of features



Origin 1: Texture recognition

Texture is characterized by the repetition of basic elements or *textons*

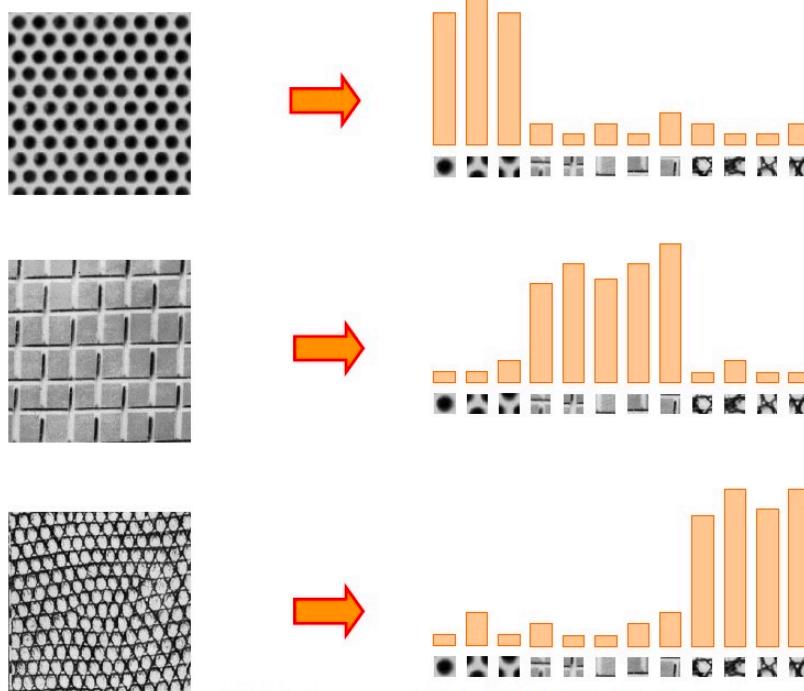
For stochastic textures, it is the identity of the textons, not their spatial arrangement, that matters



Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001;
Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

Source: Lazebnik

Origin 1: Texture recognition

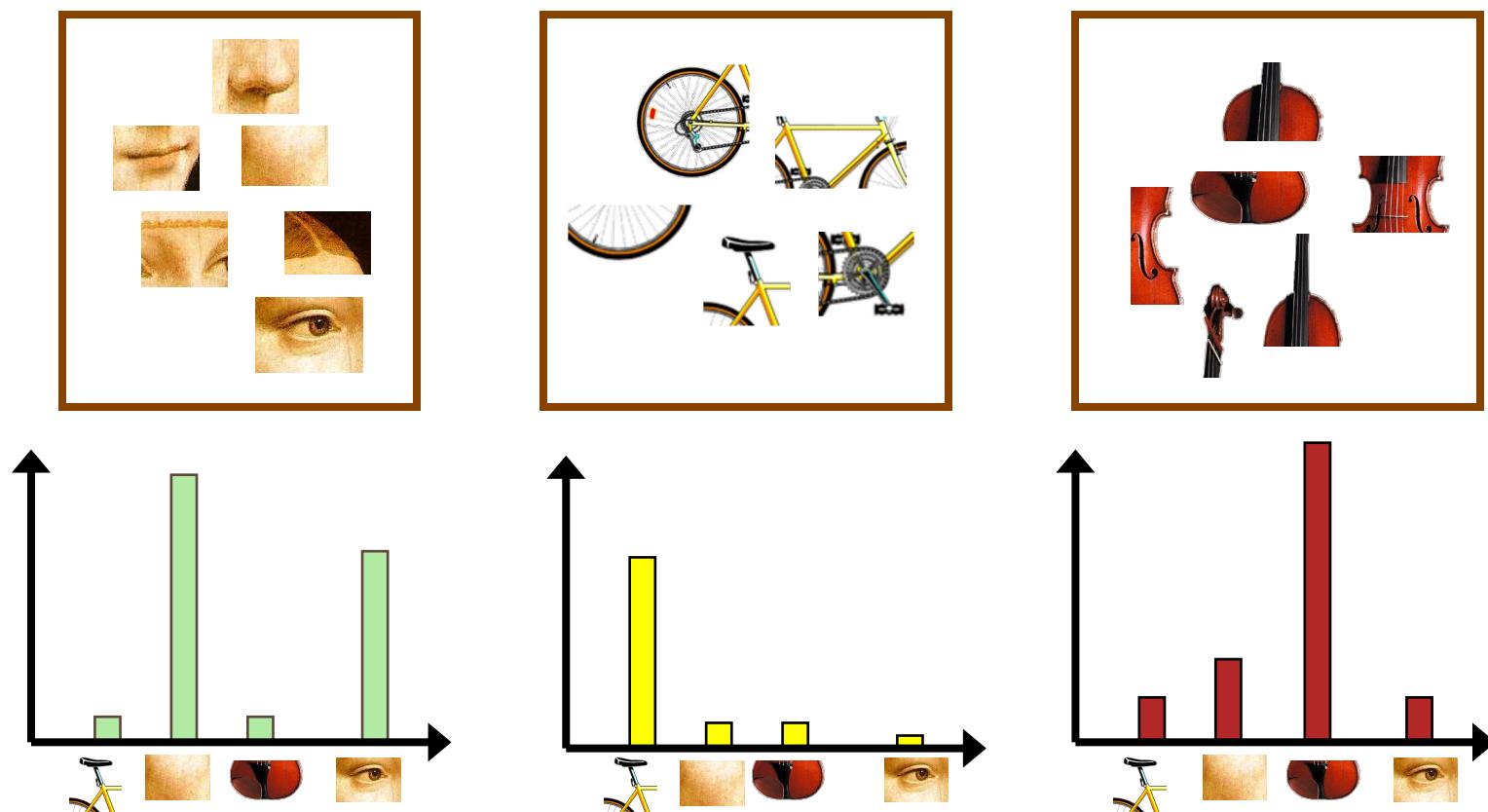


Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

Source: Lazebnik

Traditional features: Bags-of-features

1. Extract local features
2. Learn “visual vocabulary”
3. Quantize local features using visual vocabulary
4. Represent images by frequencies of “visual words”

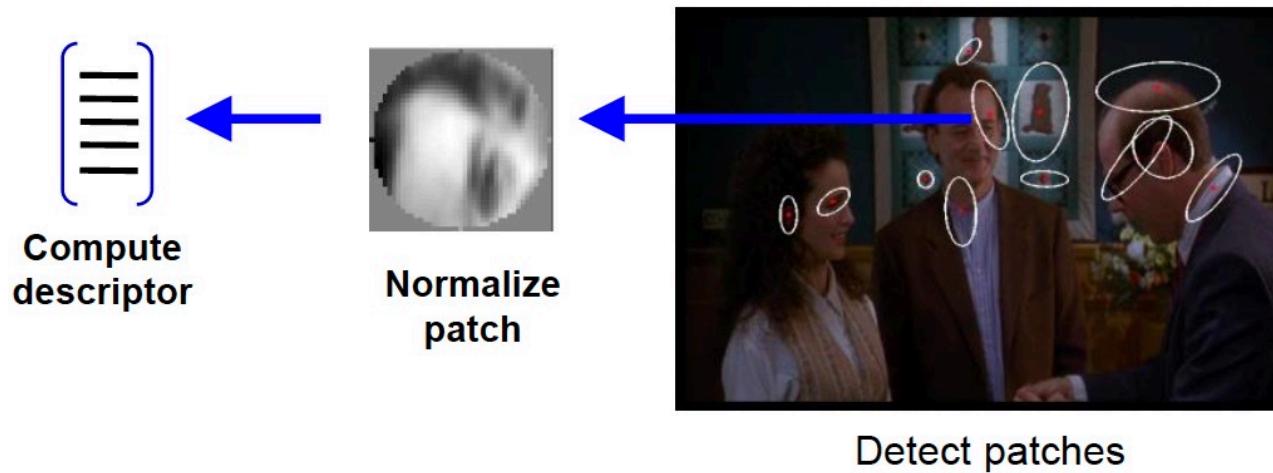


1. Local feature extraction

- Sample patches and extract descriptors

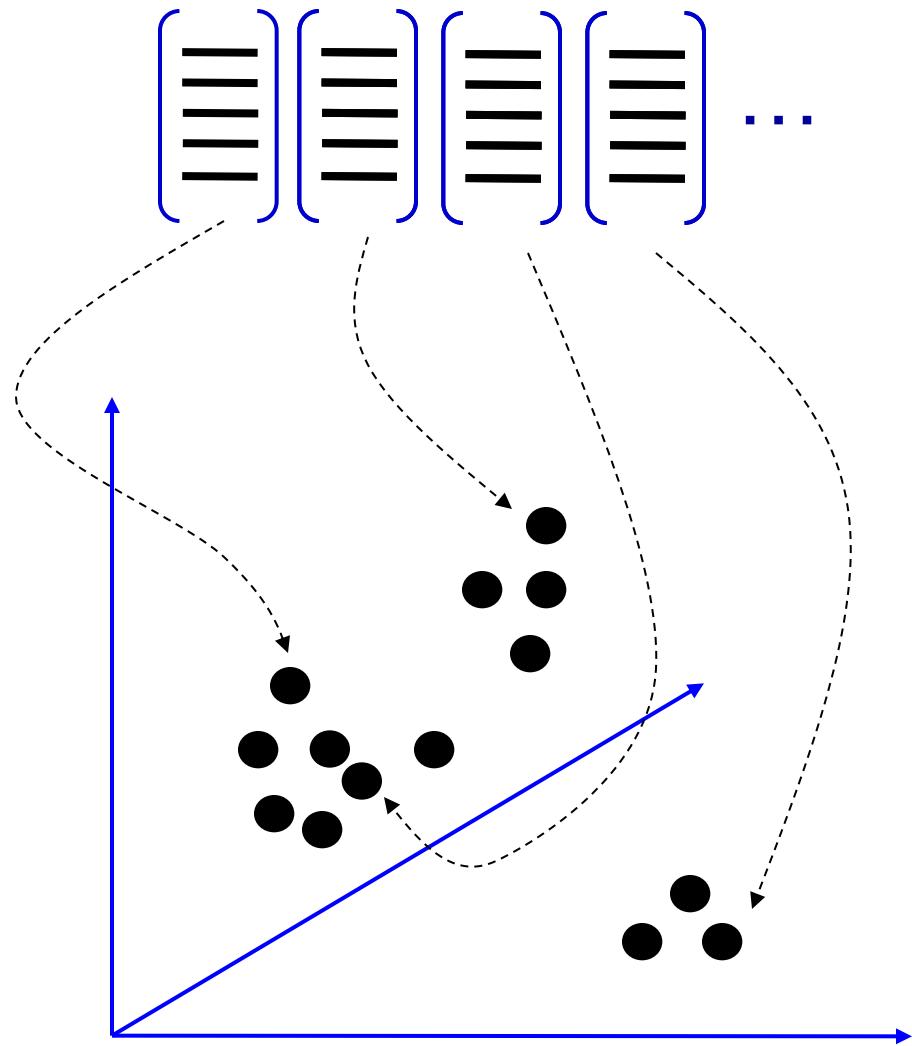


1. Feature extraction



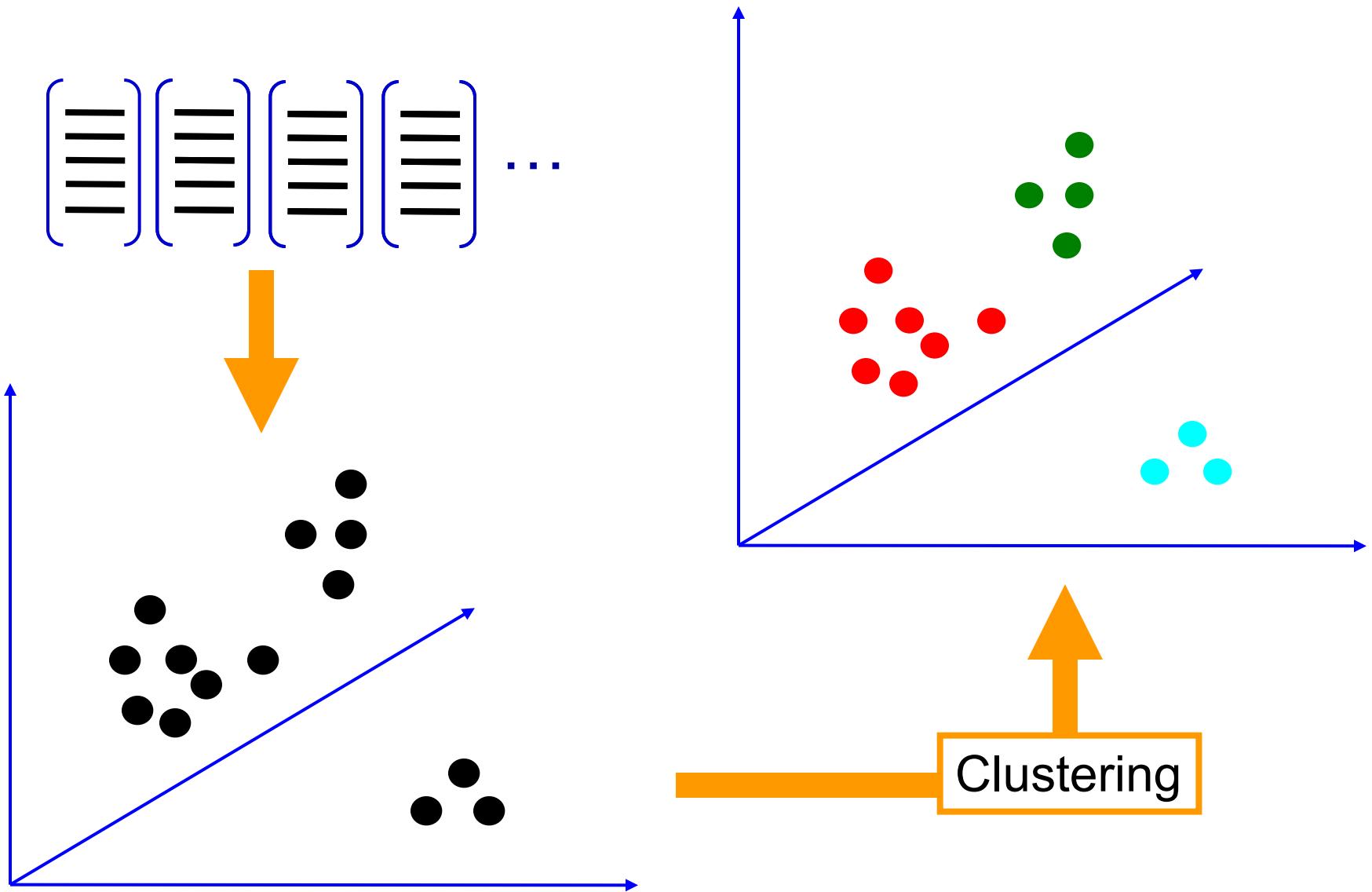
Source: Josef Sivic

2. Learning the visual vocabulary

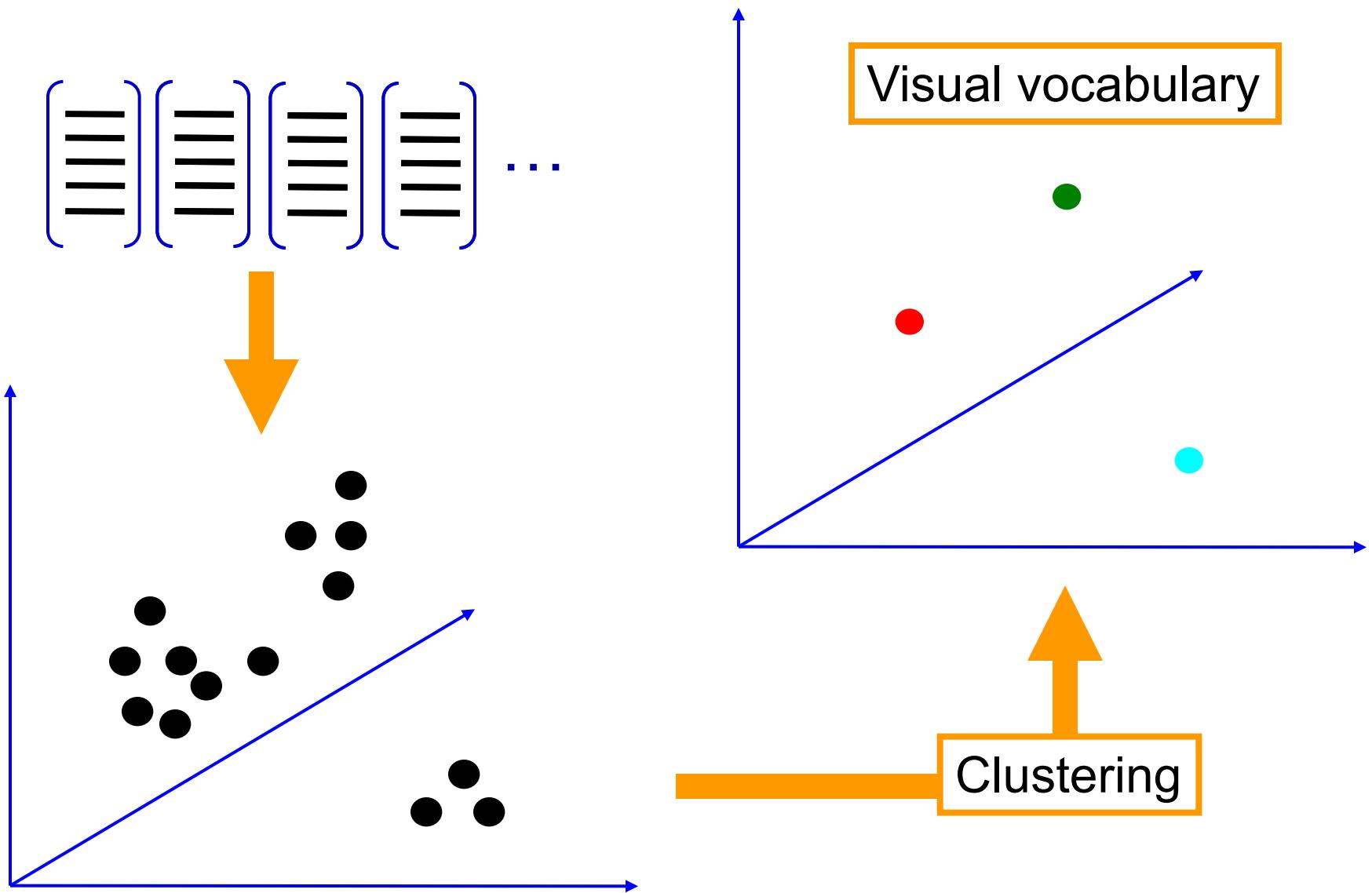


Extracted descriptors
from the training set

2. Learning the visual vocabulary



2. Learning the visual vocabulary



K-means

A clustering algorithm (unsupervised learning method)

Group similar things together.

K-means clustering

- Want to minimize sum of squared Euclidean distances between points \mathbf{x}_i and their nearest cluster centers \mathbf{m}_k

$$D(X, M) = \sum_{\text{cluster } k} \sum_{\substack{\text{point } i \text{ in} \\ \text{cluster } k}} (\mathbf{x}_i - \mathbf{m}_k)^2$$

Algorithm:

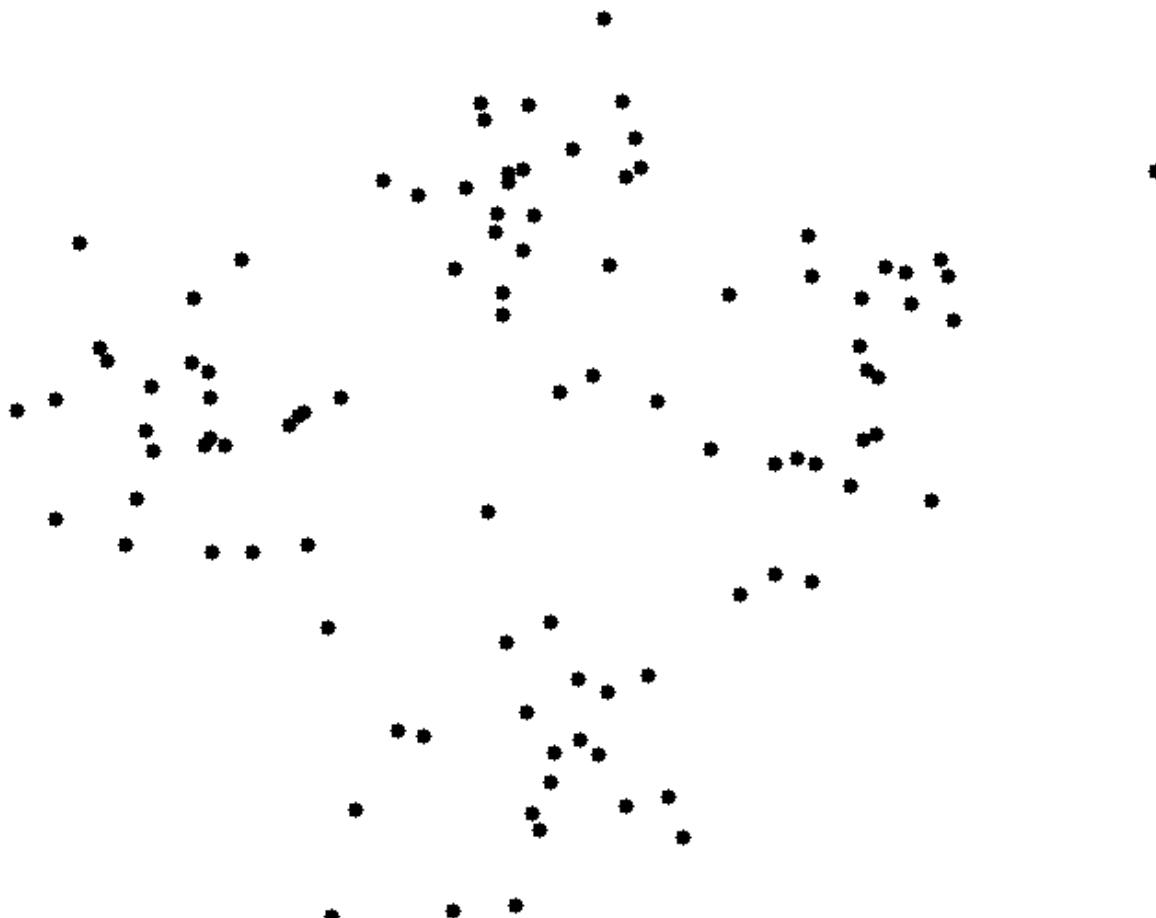
- Randomly initialize K cluster centers
- Iterate until convergence:
 - Assign each data point to the nearest center
 - Recompute each cluster center as the mean of all points assigned to it

K-means demo

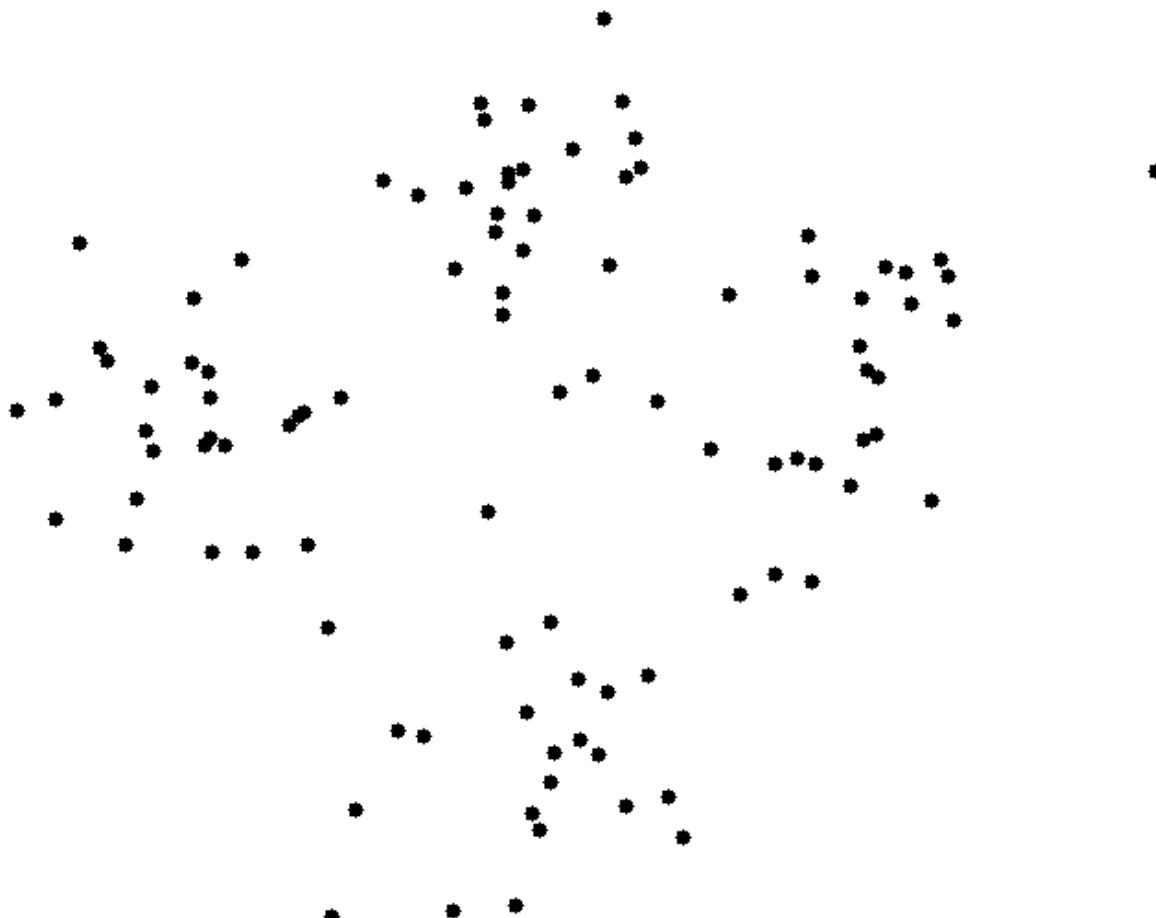


Source: <http://shabal.in/visuals/kmeans/1.html>

Another demo: <http://www.kovan.ceng.metu.edu.tr/~maya/kmeans/>

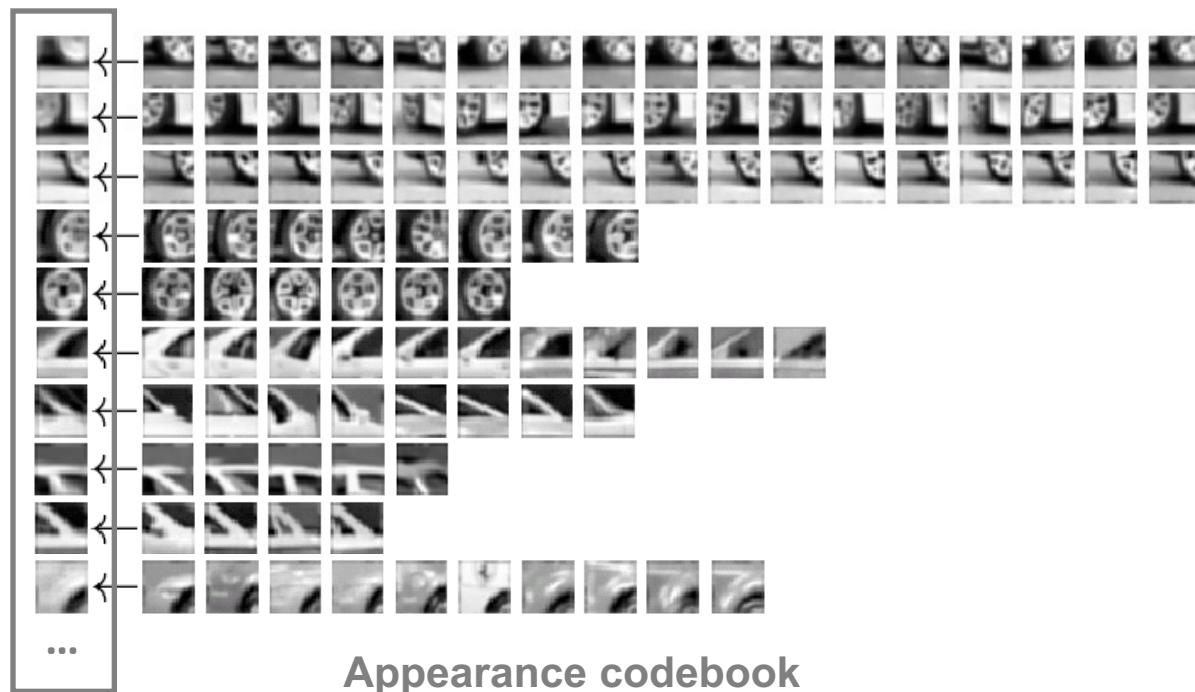
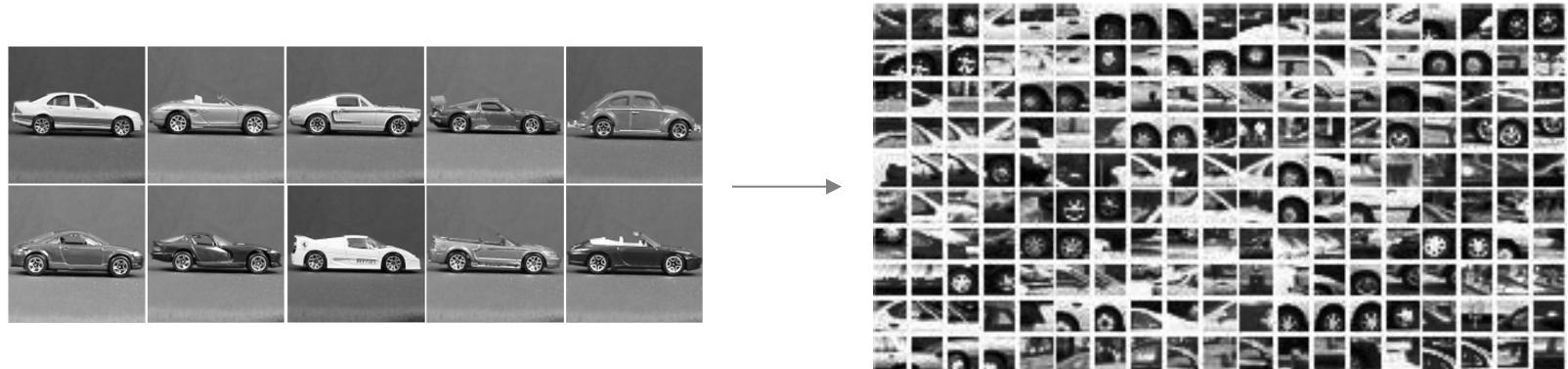


Local Minimum



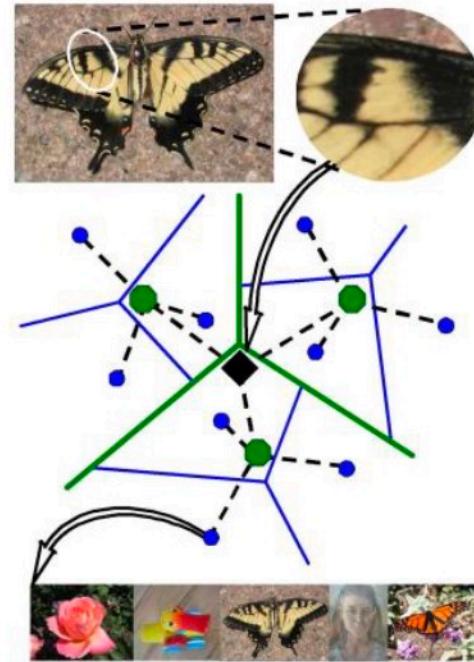
Local Minimum

Example visual vocabulary



Visual vocabularies: Issues

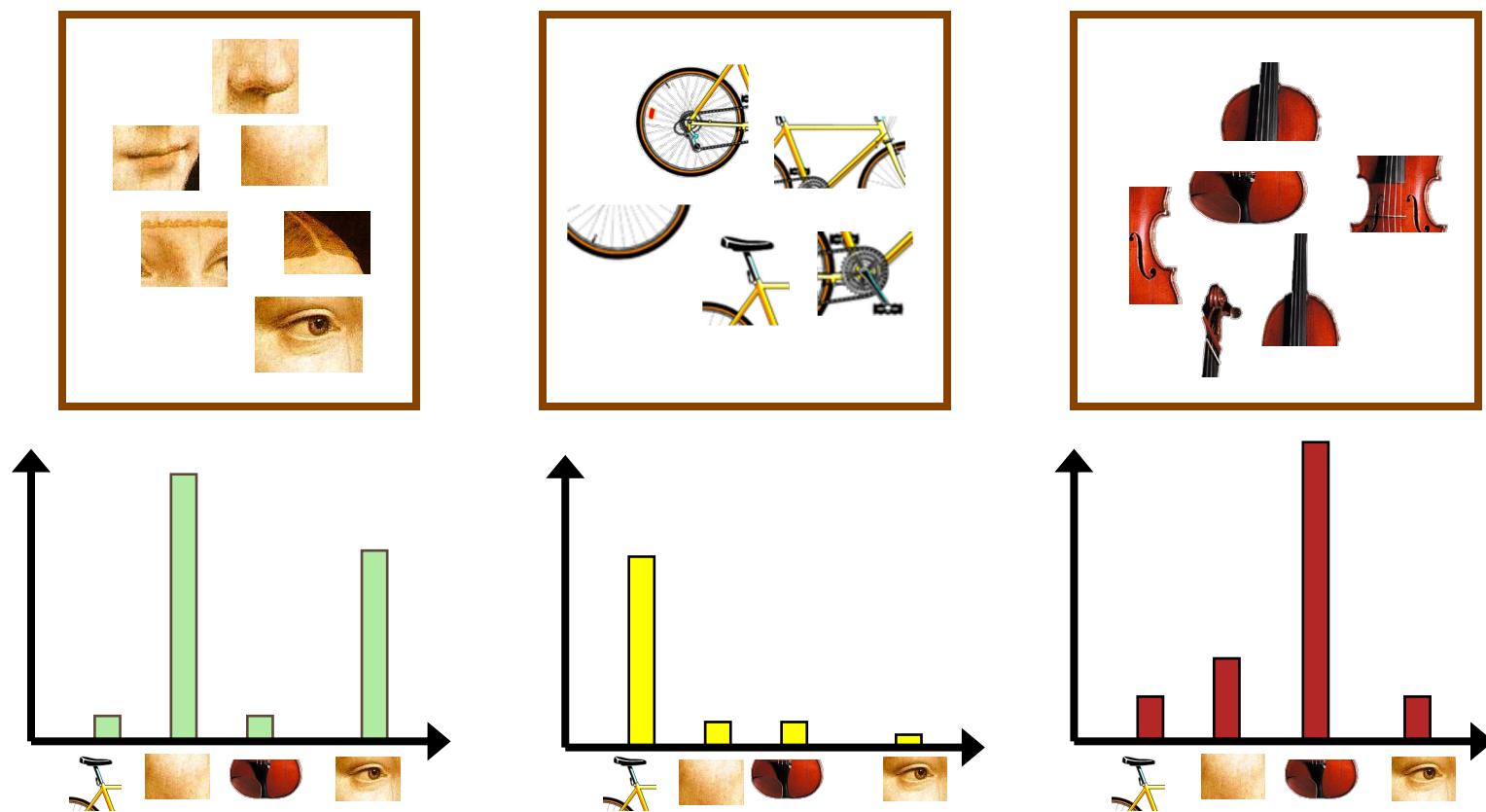
- How to choose vocabulary size?
 - Too small: visual words not representative of all patches
 - Too large: quantization artifacts, overfitting
- Computational efficiency
 - Vocabulary trees
(Nister & Stewenius, 2006)



Source: Lazebnik

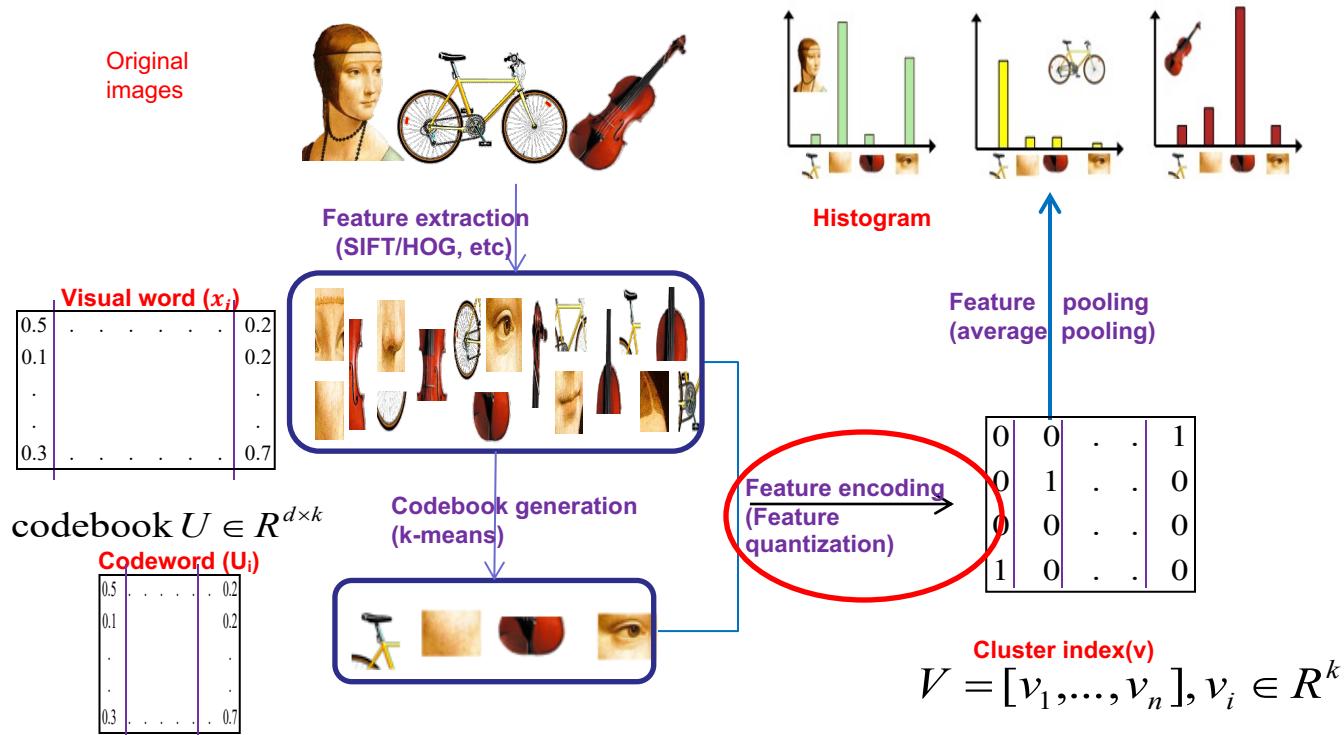
Bag-of-features steps

1. Extract local features
2. Learn “visual vocabulary”
3. **Quantize local features using visual vocabulary**
4. Represent images by frequencies of “visual words”



Bag-of-(Visual)Words (BoW) model:

(i) Feature extraction; (ii) Dictionary learning and Feature encoding; (iii) pooling



(credited to Fei-Fei Li)

Bags of features: Motivation

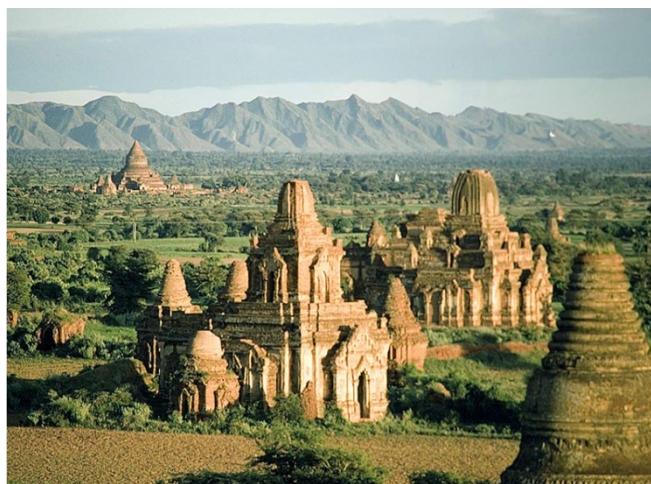
- Orderless document representation: frequencies of words from a dictionary Salton & McGill (1983)

Spatial pyramid representation

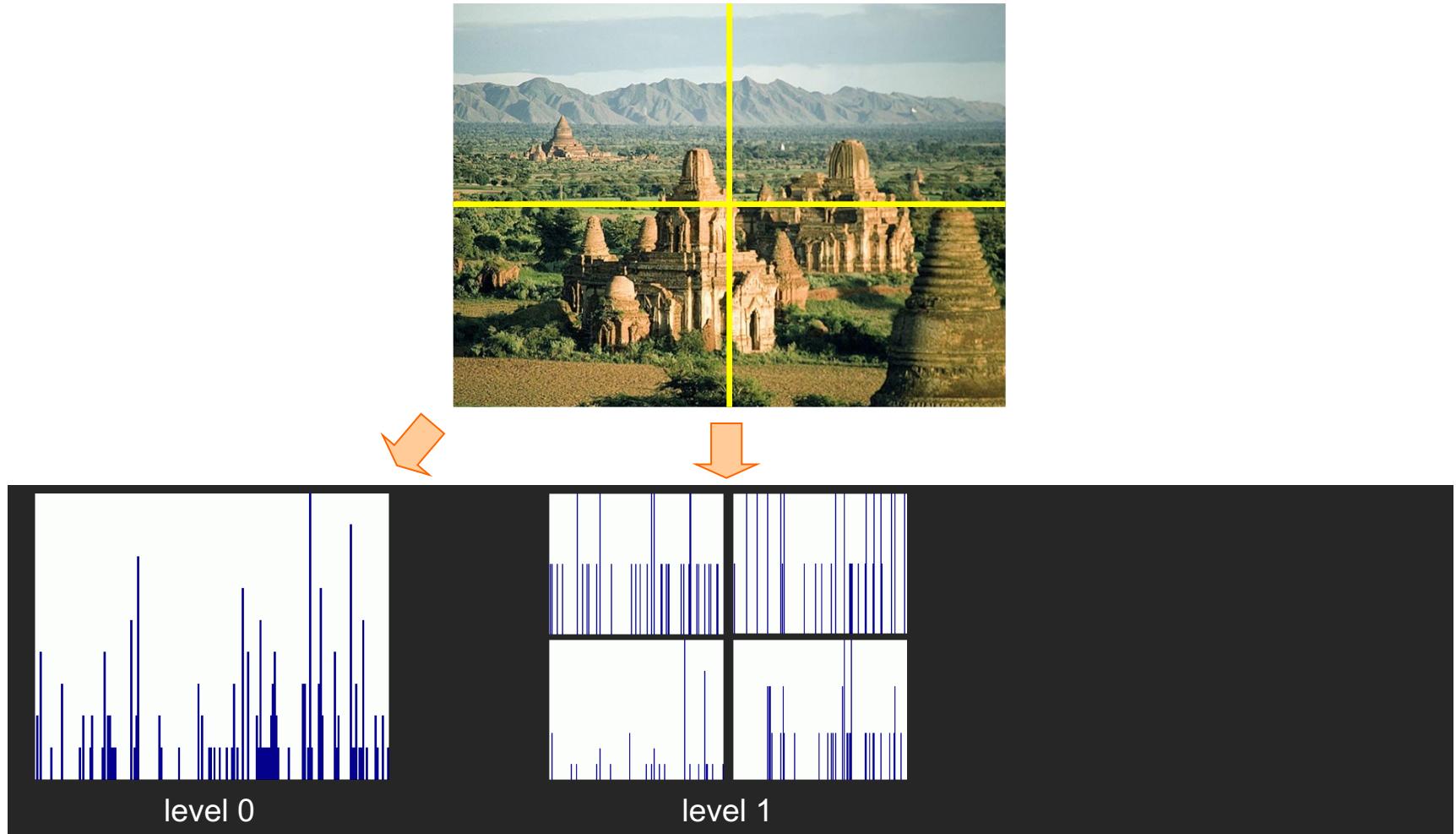
Lazebnik, Schmid & Ponce (CVPR 2006)

Extension of a bag of features

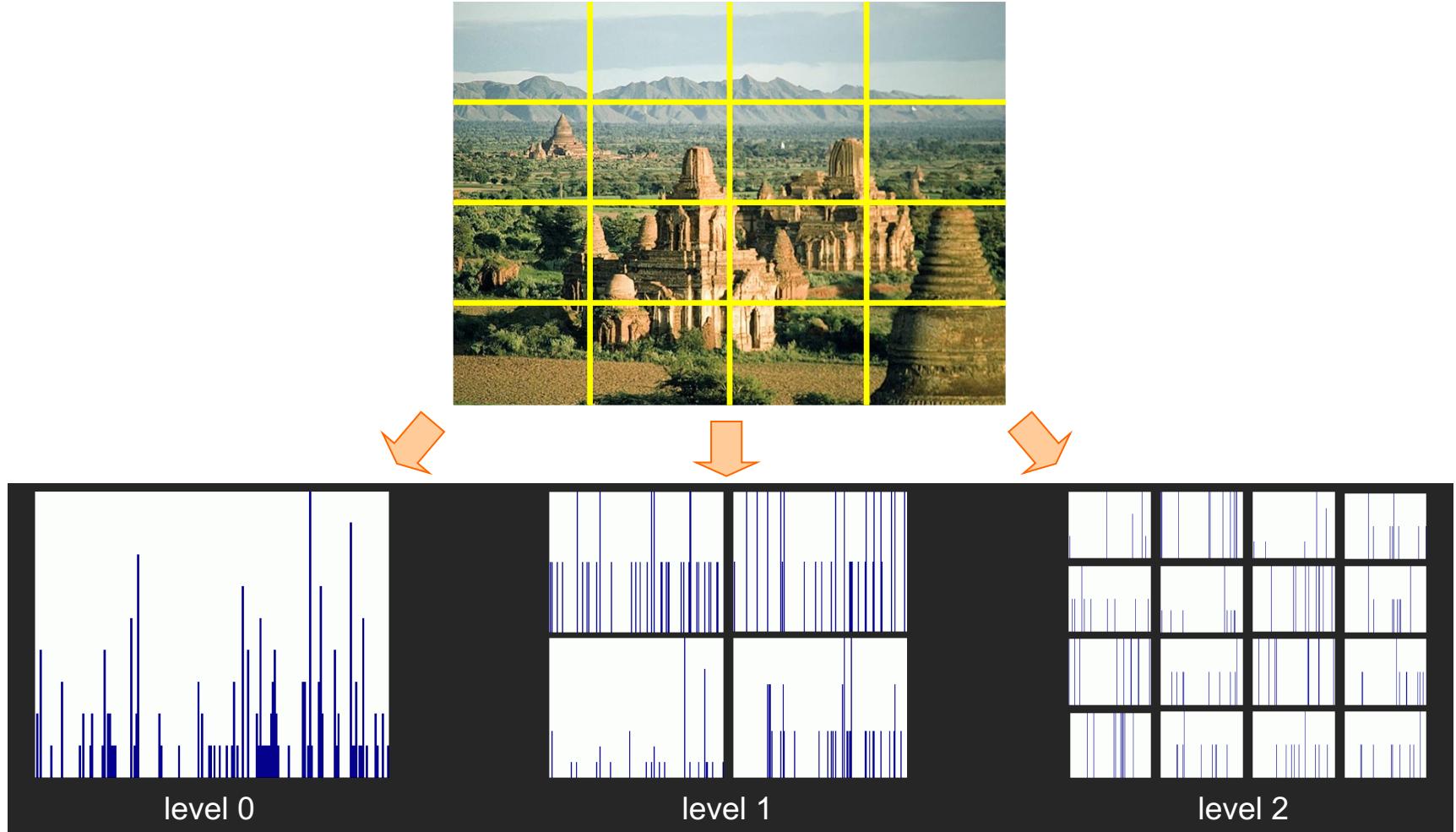
Locally orderless representation at several levels of resolution



Spatial pyramids

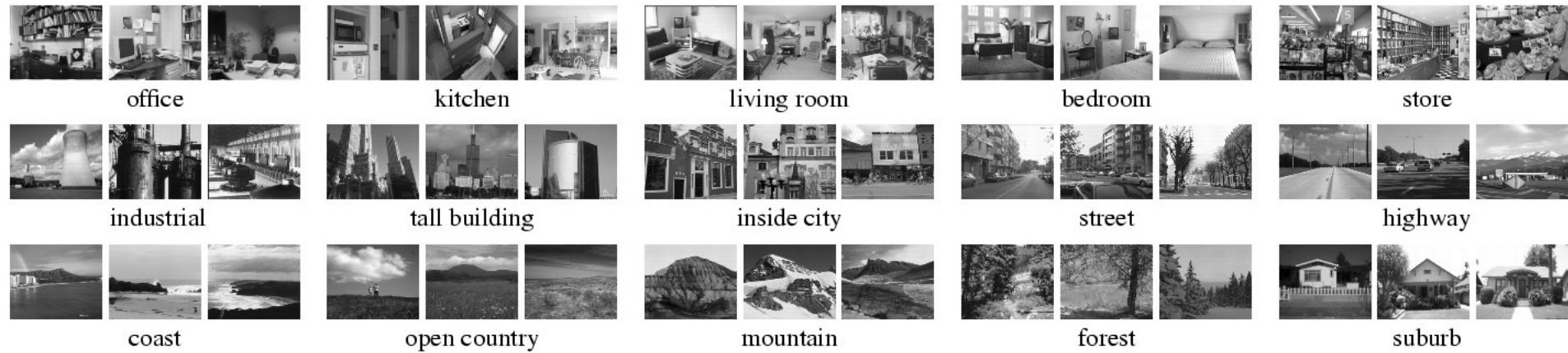


Spatial pyramids



Spatial pyramids

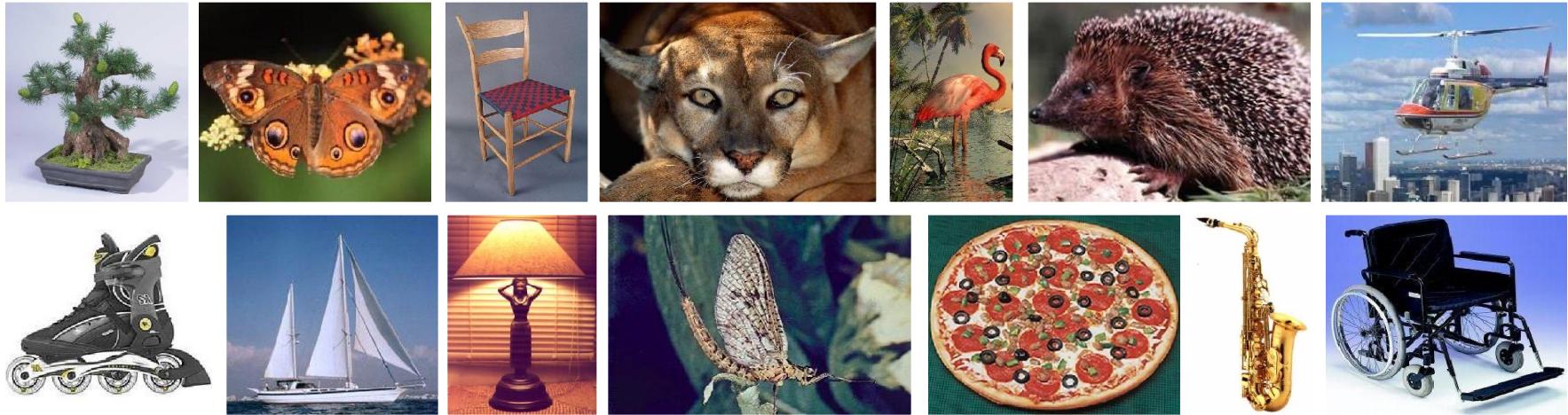
- Scene classification results



	Weak features (vocabulary size: 16)		Strong features (vocabulary size: 200)	
Level	Single-level	Pyramid	Single-level	Pyramid
0 (1×1)	45.3 ± 0.5		72.2 ± 0.6	
1 (2×2)	53.6 ± 0.3	56.2 ± 0.6	77.9 ± 0.6	79.0 ± 0.5
2 (4×4)	61.7 ± 0.6	64.7 ± 0.7	79.4 ± 0.3	81.1 ± 0.3
3 (8×8)	63.3 ± 0.8	66.8 ± 0.6	77.2 ± 0.4	80.7 ± 0.3

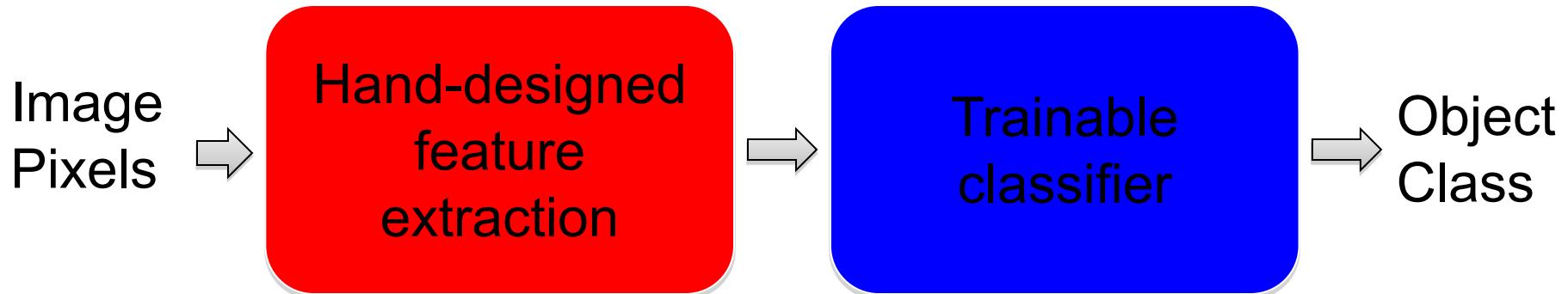
Spatial pyramids

- Caltech101 classification results

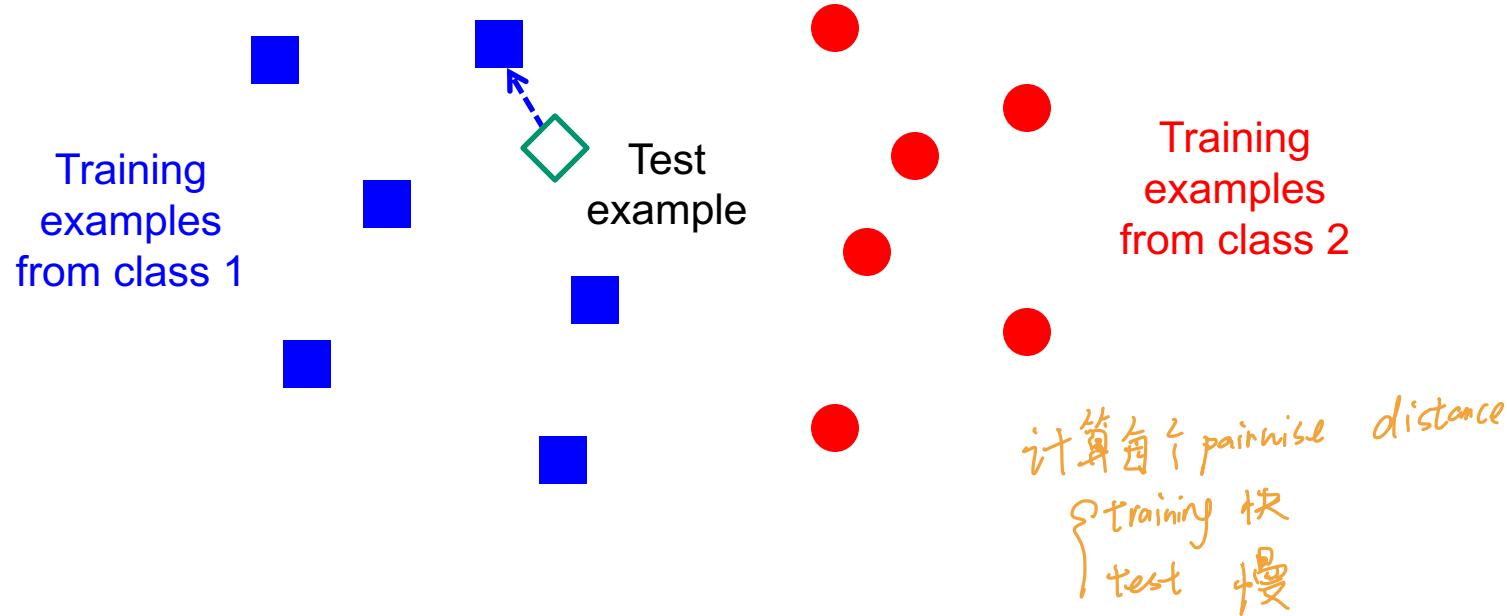


	Weak features (16)		Strong features (200)	
Level	Single-level	Pyramid	Single-level	Pyramid
0	15.5 ± 0.9		41.2 ± 1.2	
1	31.4 ± 1.2	32.8 ± 1.3	55.9 ± 0.9	57.0 ± 0.8
2	47.2 ± 1.1	49.3 ± 1.4	63.6 ± 0.9	$\mathbf{64.6} \pm 0.8$
3	52.2 ± 0.8	$\mathbf{54.0} \pm 1.1$	60.3 ± 0.9	64.6 ± 0.7

Traditional recognition pipeline



Classifiers: Nearest neighbor



$f(x)$ = label of the training example nearest to x

All we need is a distance function for our inputs

No training required!

First classifier: Nearest Neighbor



Training data with labels



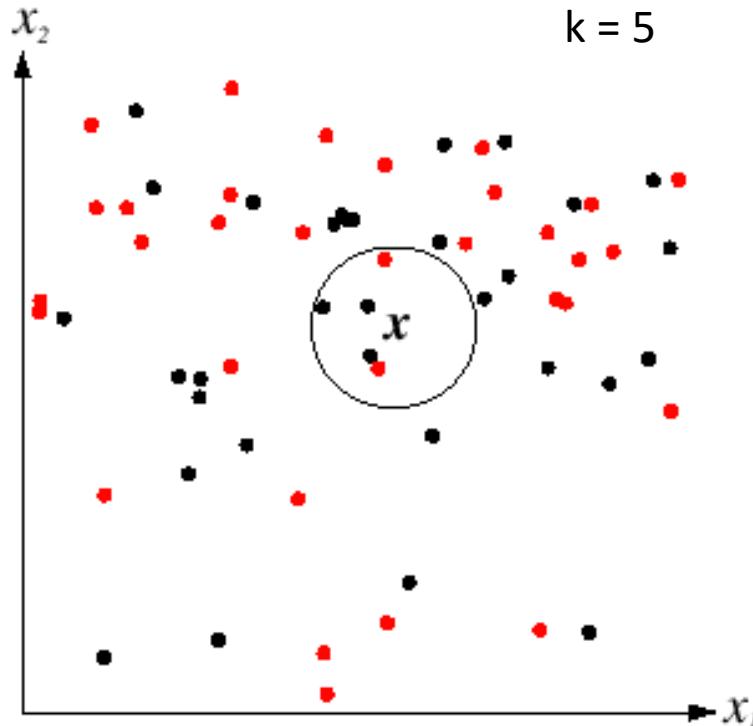
query data

Distance Metric | , → ℝ



K-nearest neighbor classifier

- For a new point, find the k closest points from training data
- Vote for class label with labels of the k points



Distance Metric to compare images

L1 distance:

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

test image

56	32	10	18
90	23	128	133
24	26	178	200
2	0	255	220

training image

10	20	24	17
8	10	89	100
12	16	178	170
4	32	233	112

-

pixel-wise absolute value differences

46	12	14	1
82	13	39	33
12	10	0	30
2	32	22	108

=

add → 456

```

import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred

```

Nearest Neighbor classifier

无参数 \Rightarrow 无学习过程 Memoriza
时间都用在 prediction

```
import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

        return Ypred
```

Nearest Neighbor classifier

Memorize training data

Fei-Fei Li, Jiajun Wu, Ruohan Gao

```
import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred
```

Nearest Neighbor classifier

For each test image:
Find closest train image
Predict label of nearest image

```

import numpy as np

class NearestNeighbor:
    def __init__(self):
        pass

    def train(self, X, y):
        """ X is N x D where each row is an example. Y is 1-dimension of size N """
        # the nearest neighbor classifier simply remembers all the training data
        self.Xtr = X
        self.ytr = y

    def predict(self, X):
        """ X is N x D where each row is an example we wish to predict label for """
        num_test = X.shape[0]
        # lets make sure that the output type matches the input type
        Ypred = np.zeros(num_test, dtype = self.ytr.dtype)

        # loop over all test rows
        for i in xrange(num_test):
            # find the nearest training image to the i'th test image
            # using the L1 distance (sum of absolute value differences)
            distances = np.sum(np.abs(self.Xtr - X[i,:]), axis = 1)
            min_index = np.argmin(distances) # get the index with smallest distance
            Ypred[i] = self.ytr[min_index] # predict the label of the nearest example

    return Ypred

```

Nearest Neighbor classifier

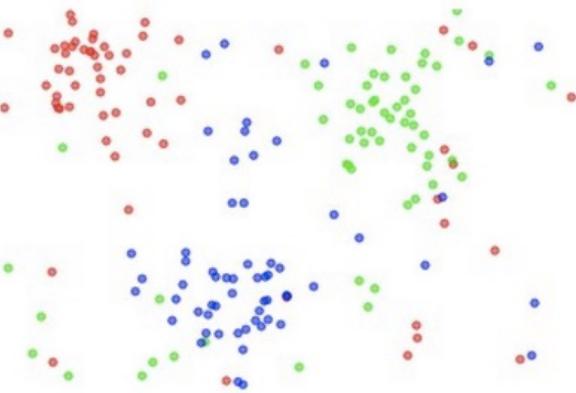
Q: With N examples, how fast are training and prediction?

Ans: Train $O(1)$, predict $O(N)$

This is bad: we want classifiers that are **fast** at prediction; **slow** for training is ok

K-nearest neighbor classifier

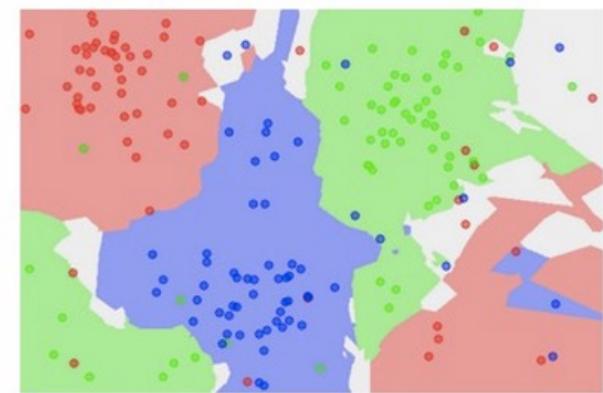
the data



NN classifier



5-NN classifier



Which classifier is more robust to *outliers*?

Hyperparameters

What is the best value of k to use?

What is the best **distance** to use?

These are **hyperparameters**: choices about the algorithms themselves.

Very problem/dataset-dependent.

Must try them all out and see what works best.

Setting Hyperparameters

cross-validation

Idea #1: Choose hyperparameters
that work best on the **training data**

train

Setting Hyperparameters

Idea #1: Choose hyperparameters
that work best on the **training data**

BAD: K = 1 always works
perfectly on training data

train

Setting Hyperparameters

Idea #1: Choose hyperparameters
that work best on the **training data**



train

BAD: K = 1 always works
perfectly on training data

Idea #2: choose hyperparameters
that work best on **test** data



train

test

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the **training data**

BAD: K = 1 always works perfectly on training data



train

Idea #2: choose hyperparameters that work best on **test data**

BAD: No idea how algorithm will perform on new data



train

test

Never do this!

Setting Hyperparameters

Idea #1: Choose hyperparameters that work best on the **training data**

BAD: K = 1 always works perfectly on training data



train

Idea #2: choose hyperparameters that work best on **test** data

BAD: No idea how algorithm will perform on new data



train

test

Idea #3: Split data into **train**, **val**; choose hyperparameters on val and evaluate on test

Better!



train

validation

test

Setting Hyperparameters

train

Idea #4: Cross-Validation: Split data into **folds**,
try each fold as validation and average the results

fold 1	fold 2	fold 3	fold 4	fold 5	test
--------	--------	--------	--------	--------	------

fold 1	fold 2	fold 3	fold 4	fold 5	test
--------	--------	--------	--------	--------	------

fold 1	fold 2	fold 3	fold 4	fold 5	test
--------	--------	--------	--------	--------	------

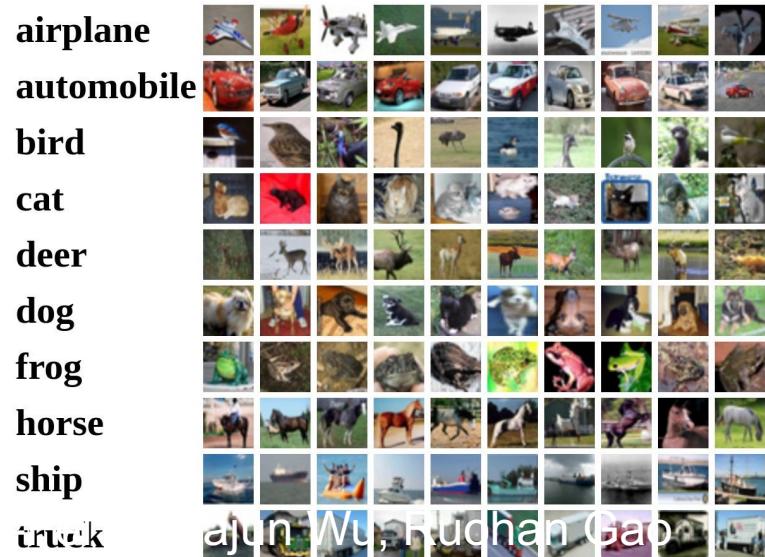
Useful for small datasets, but not used too frequently in deep learning

Example Dataset: CIFAR10

10 classes

50,000 training images

10,000 testing images



Alex Krizhevsky, "Learning Multiple Layers of Features from Tiny Images", Technical Report, 2009.

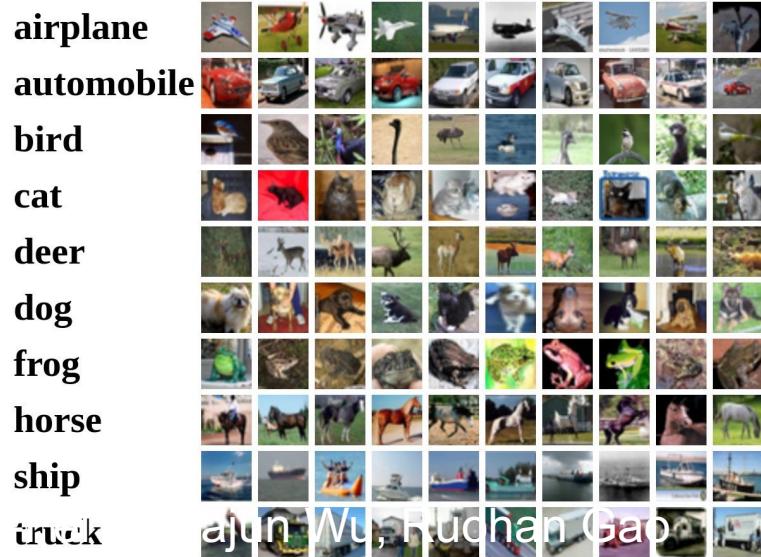
ajun Wu, Ruchan Guo

Example Dataset: CIFAR10

10 classes

50,000 training images

10,000 testing images

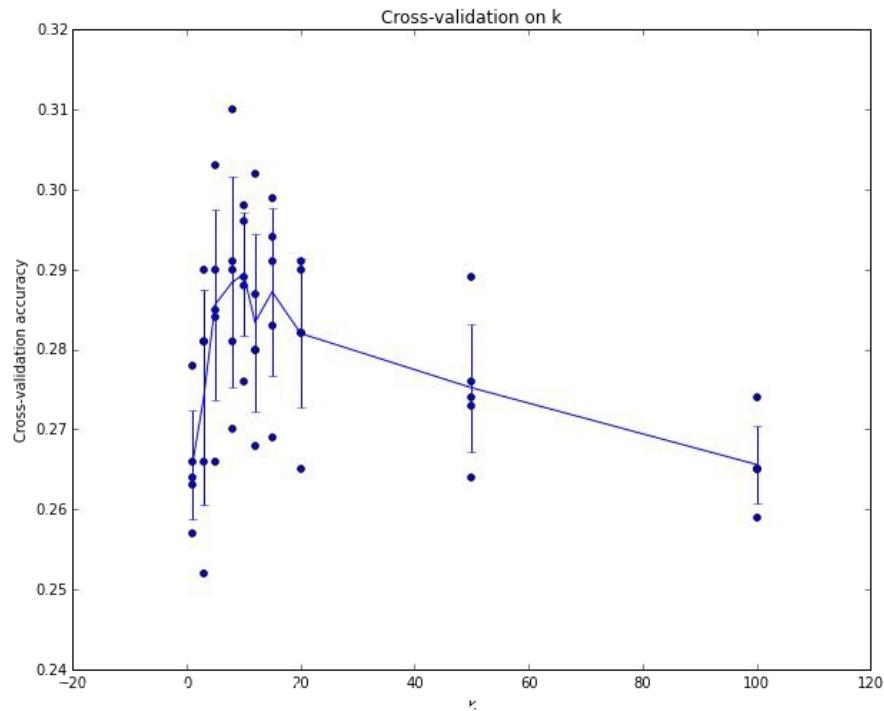


Test images and nearest neighbors



Alex Krizhevsky, "Learning Multiple Layers of Features from Tiny Images", Technical Report, 2009.

Setting Hyperparameters



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 $k \sim 7$ works best
for this data)

What does this look like?



uph ↗ Gao

77

What does this look like?



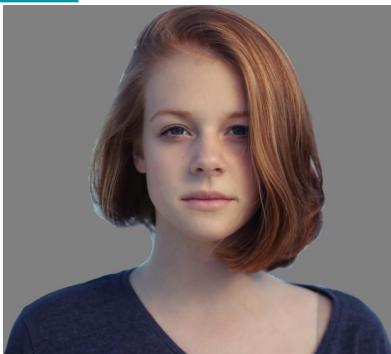
78

k-Nearest Neighbor with pixel distance **never** used.

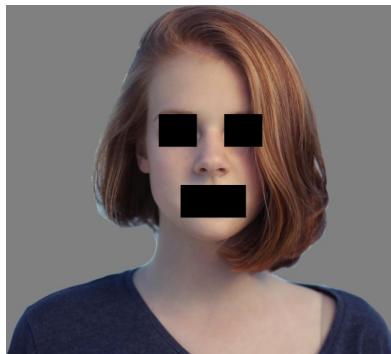
- Distance metrics on pixels are not informative

[OriginalImage.js](#)
CC0 public domain

Original



Occluded



Shifted (1 pixel)



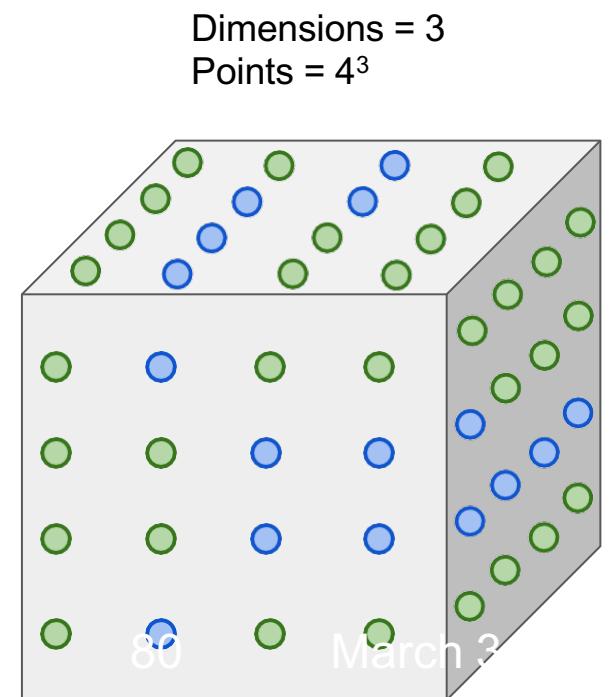
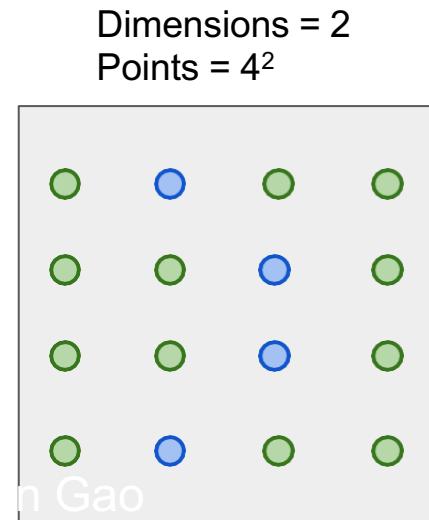
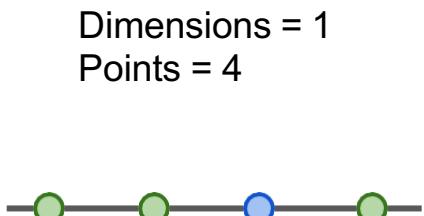
Tinted



(All three images on the right have the same pixel distances to the one on the left)

k-Nearest Neighbor with pixel distance **never** used.

- Curse of dimensionality



K-Nearest Neighbors: Summary

In **image classification** we start with a **training set** of images and labels, and must predict labels on the **test set**

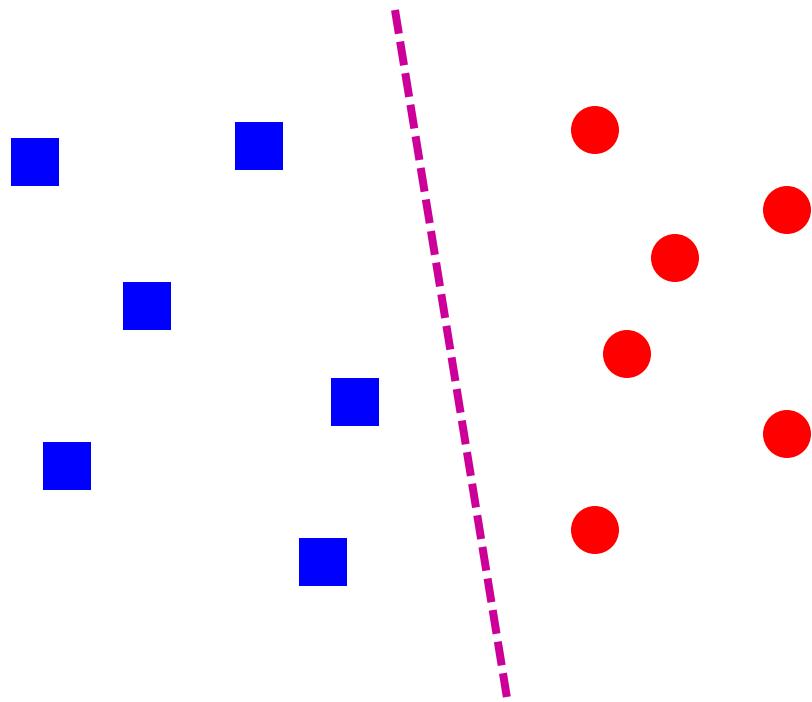
The **K-Nearest Neighbors** classifier predicts labels based on the K nearest training examples

Distance metric and K are **hyperparameters**

Choose hyperparameters using the **validation set**;

Only run on the test set once at the very end!

Linear classifiers

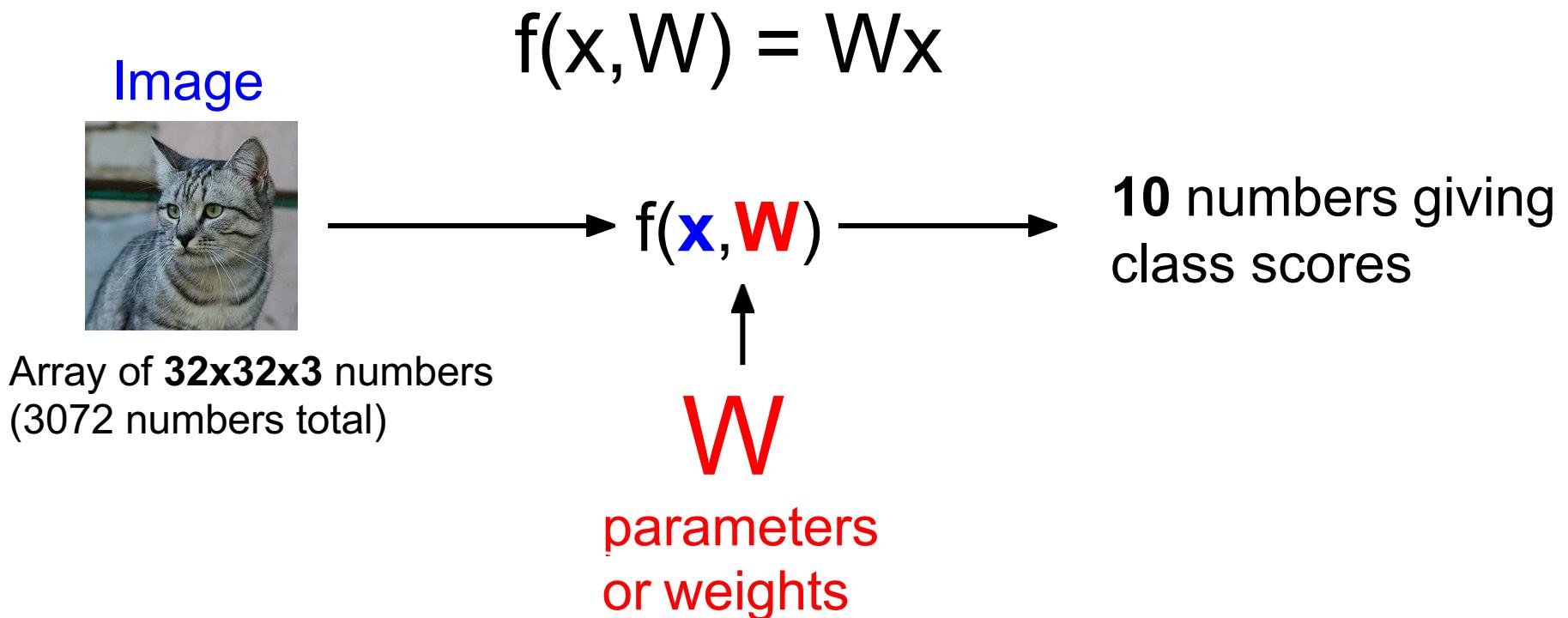


Find a *linear function* to separate the classes:

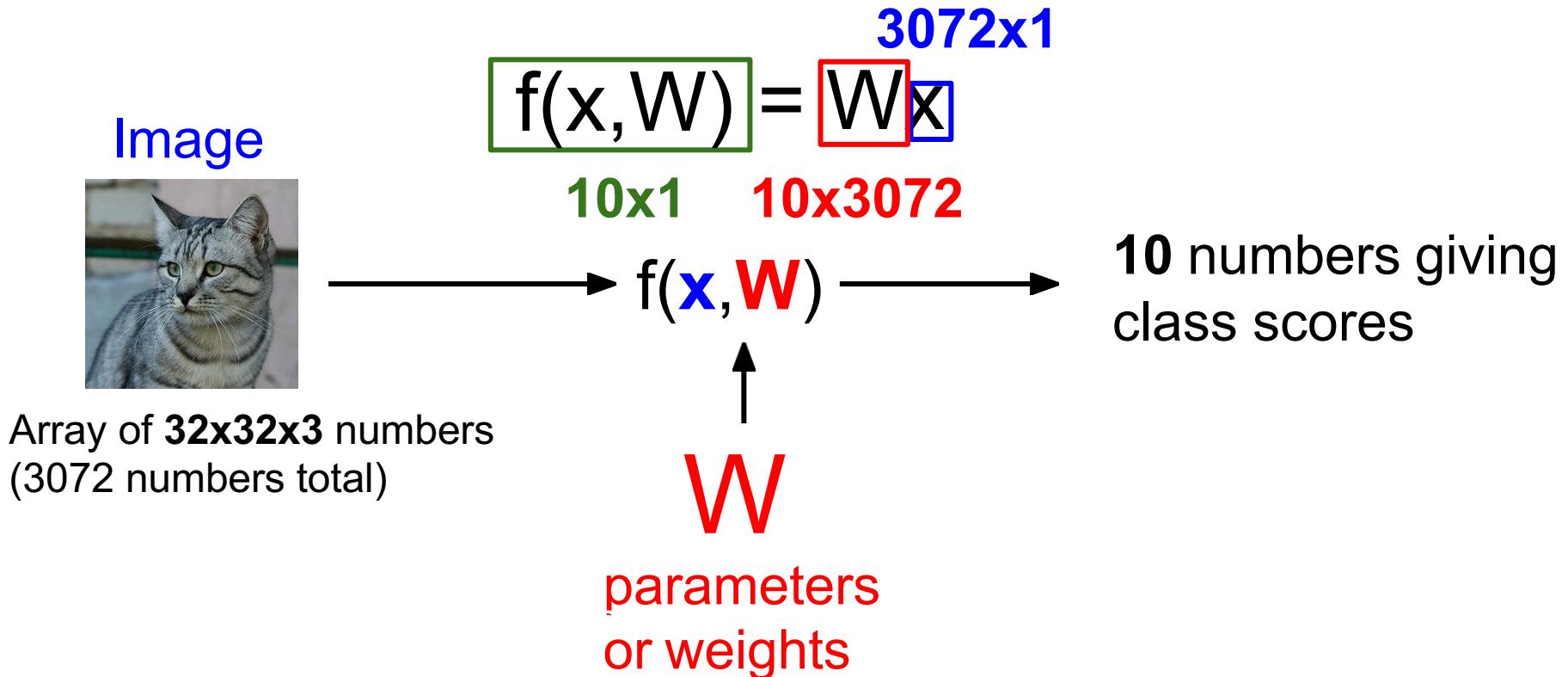
$$f(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x} + b)$$

$$\begin{cases} > 0 & 1 \\ \leq 0 & -1 \end{cases}$$

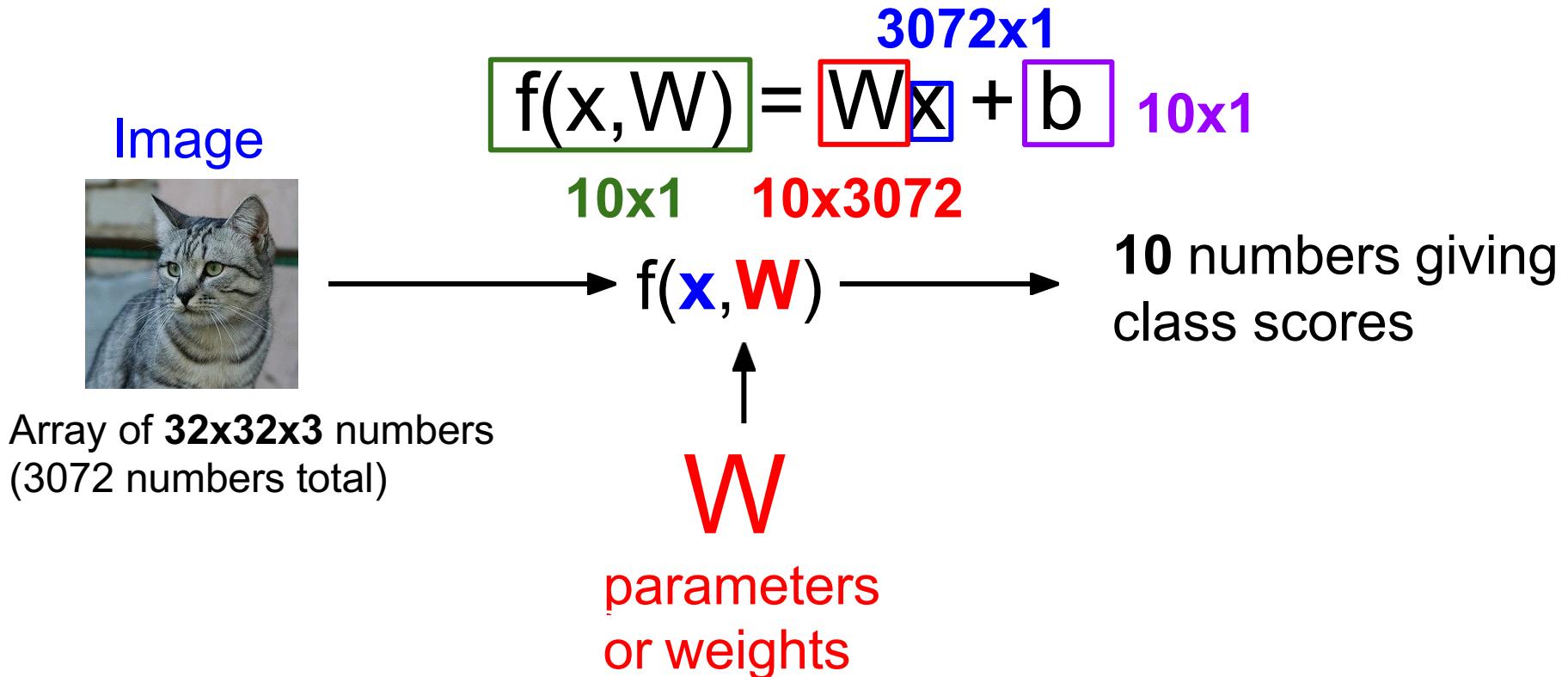
Parametric Approach: Linear Classifier



Parametric Approach: Linear Classifier

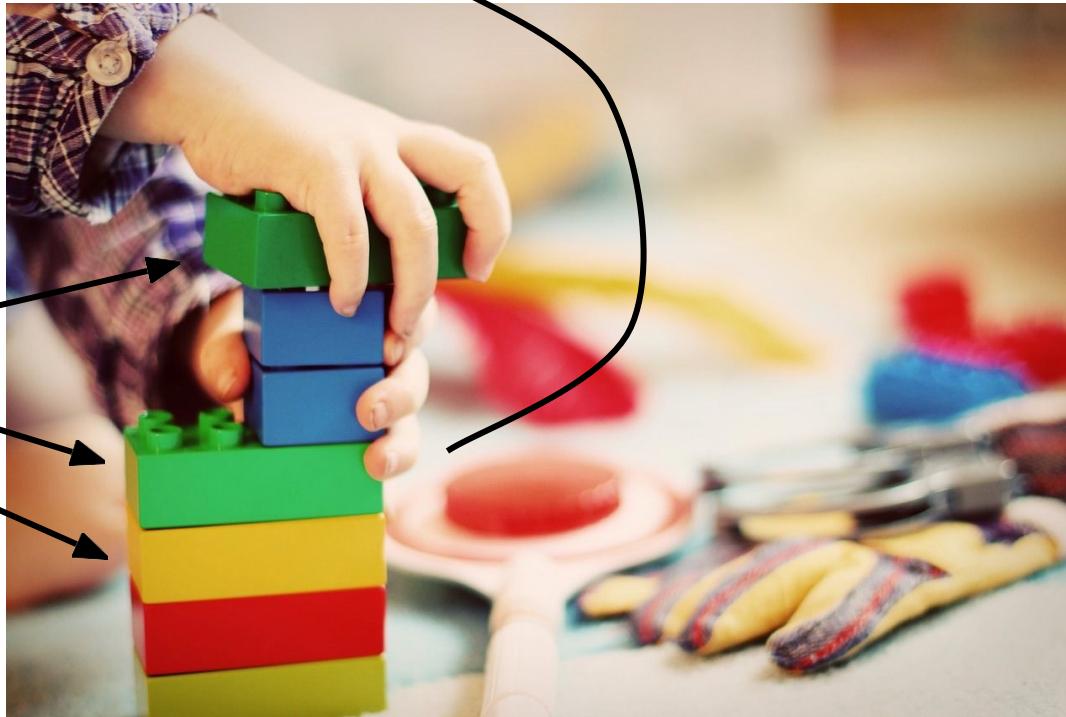


Parametric Approach: Linear Classifier

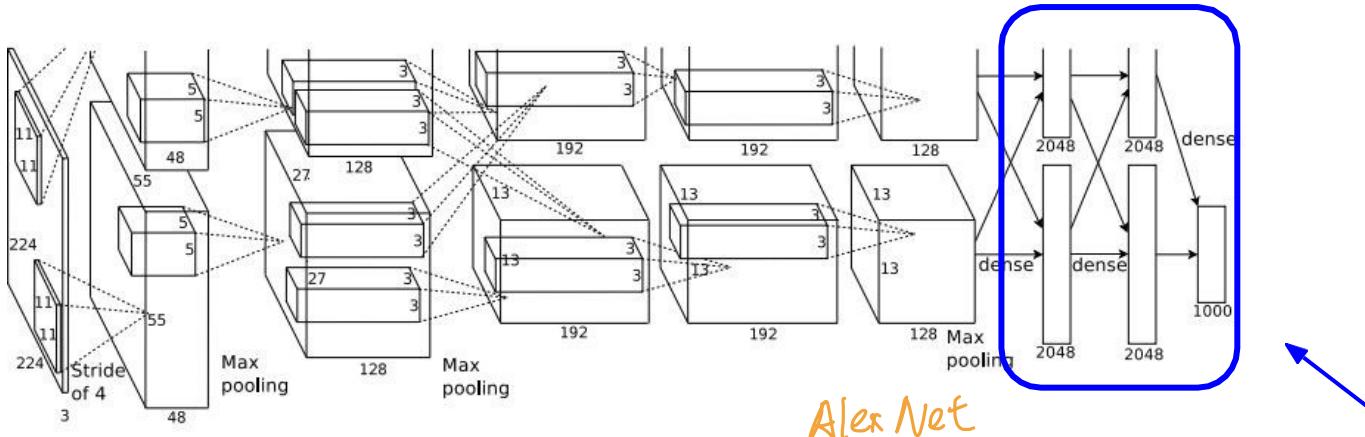


Neural Network

Linear
classifiers



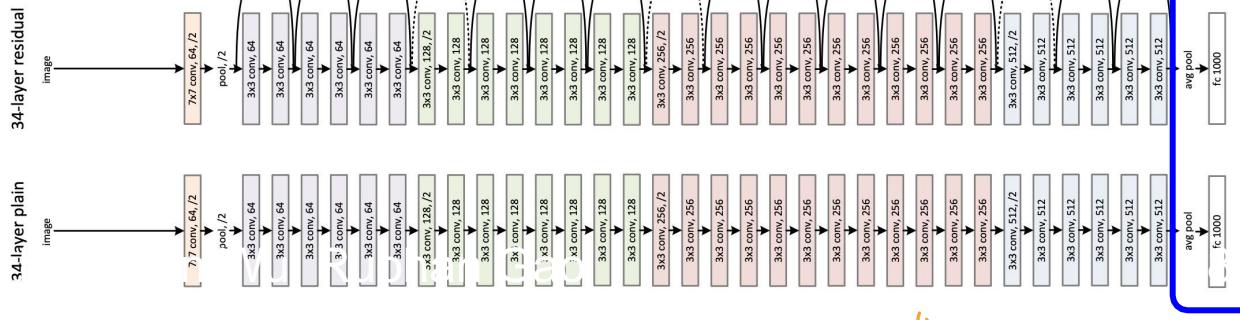
[This image](#) is CC0 1.0 public domain



[Krizhevsky et al. 2012]

Alex Net

Linear layers



[He et al. 2015]

ResNet

Recall CIFAR10

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



truck

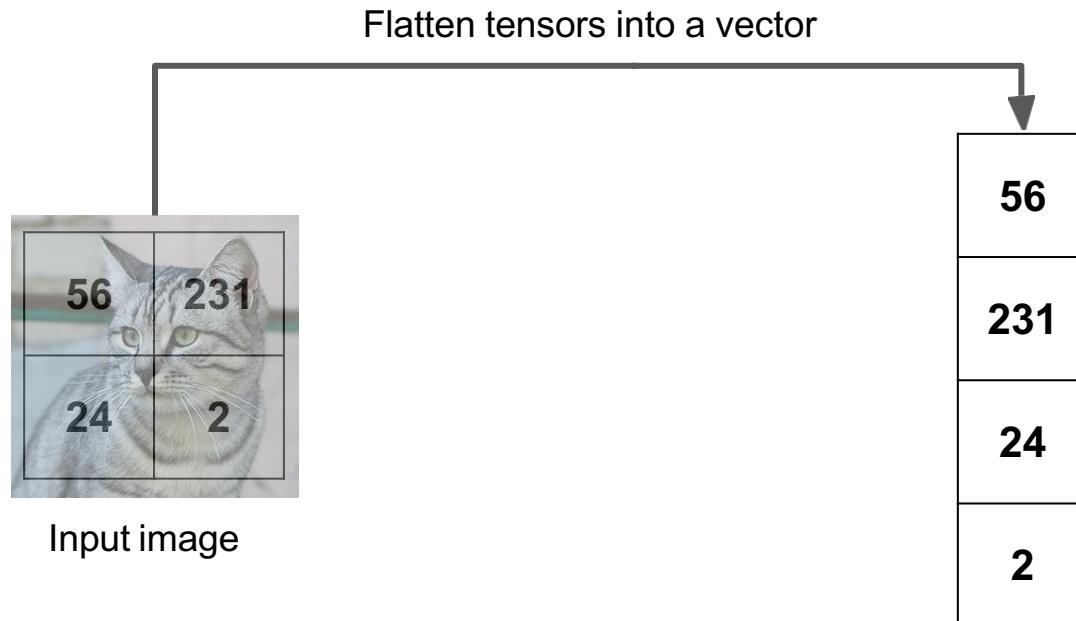


ajuh Wu, Ruohan Gab

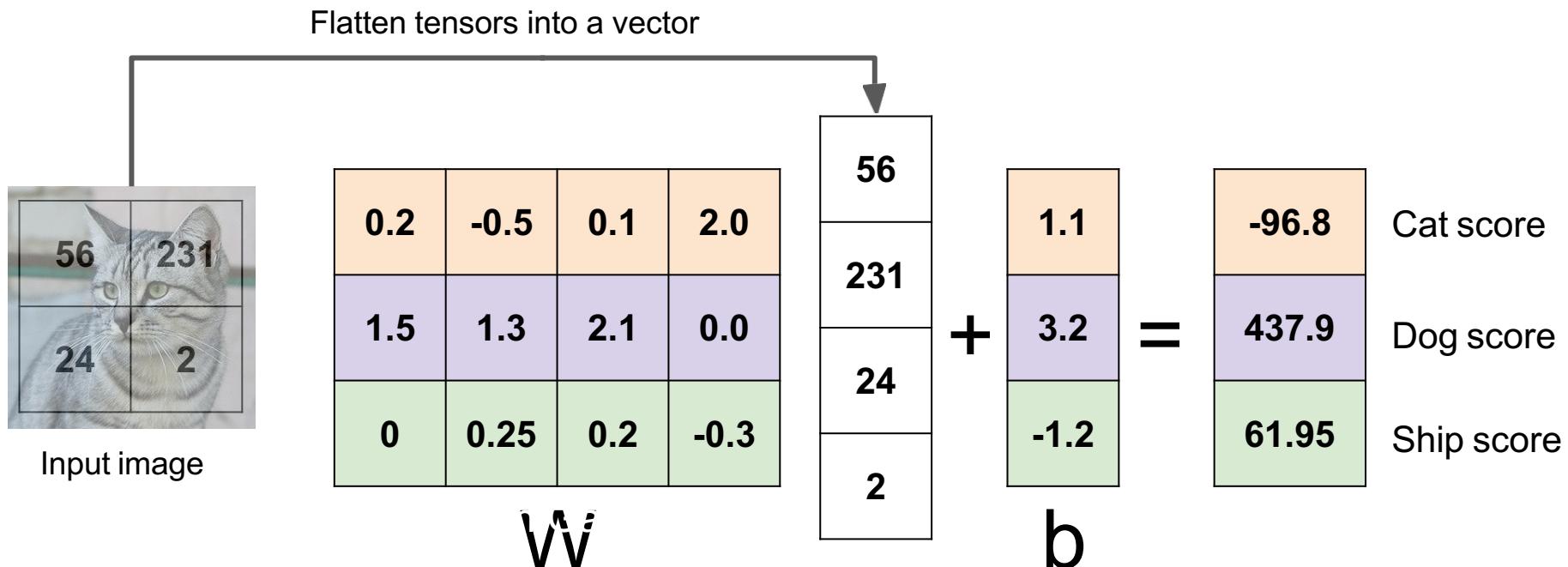
50,000 training images
each image is **32x32x3**

10,000 test images.

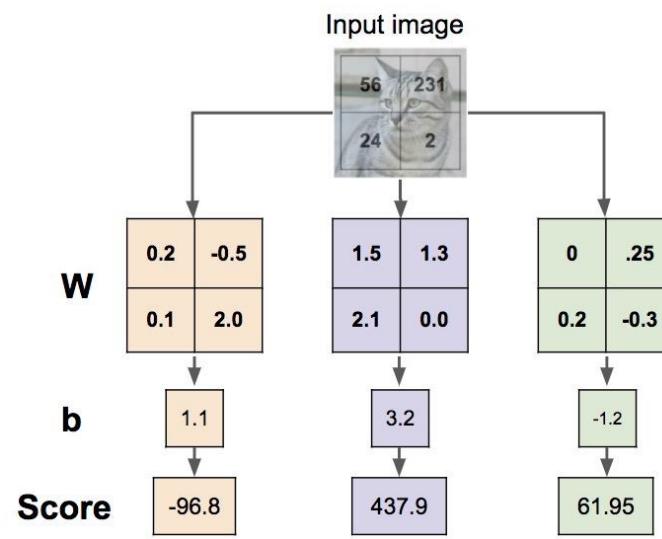
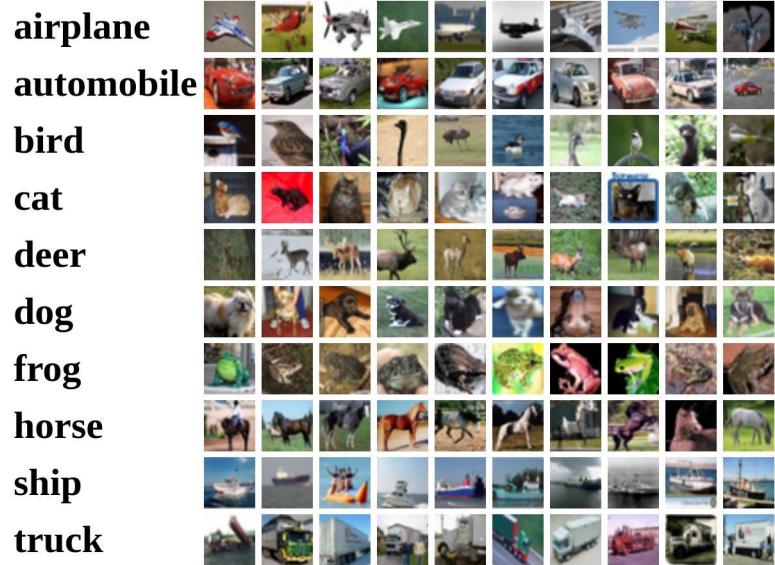
Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**)



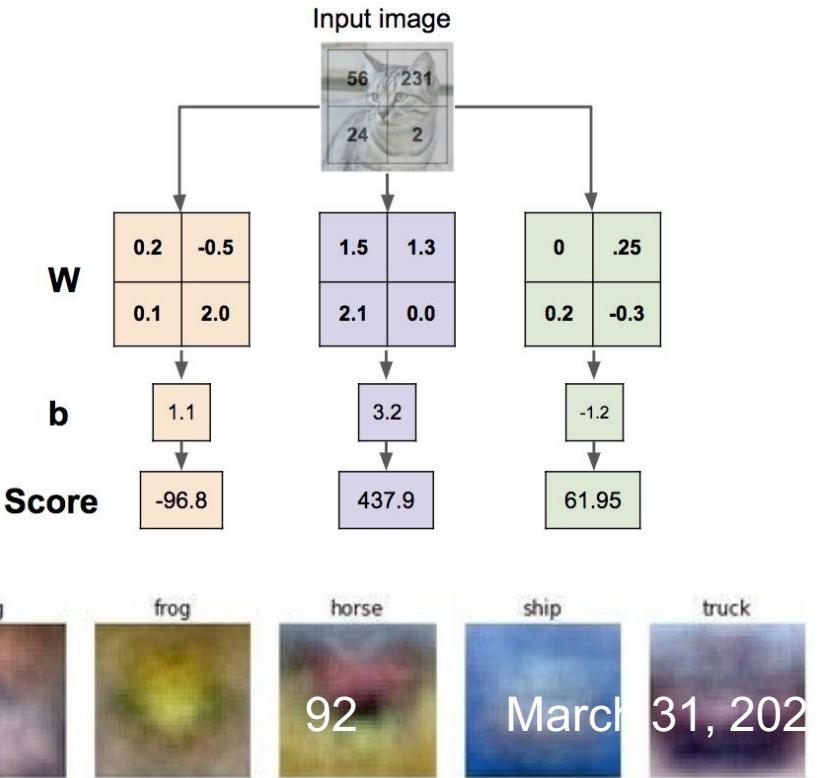
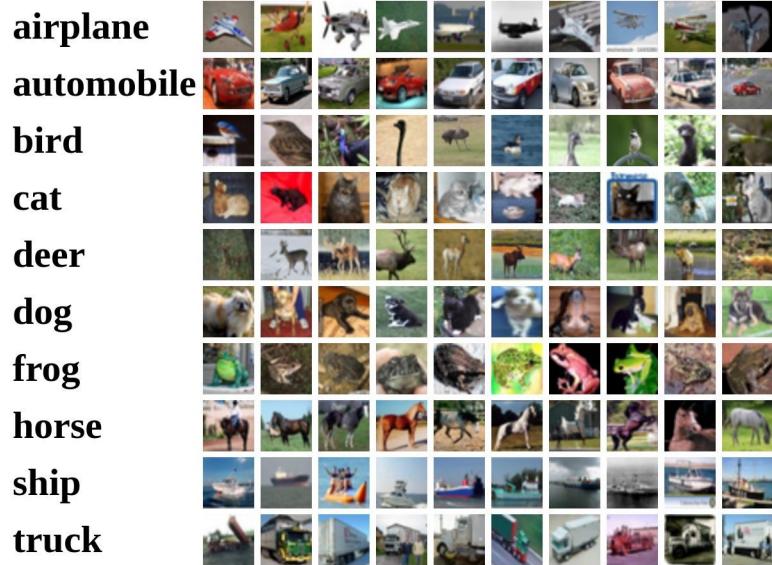
Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**) Algebraic Viewpoint



Interpreting a Linear Classifier



Interpreting a Linear Classifier: Visual Viewpoint



plane
ei-Fei L

car
Jiajun Wu

bird
Ruohan Gao

cat

deer

dog

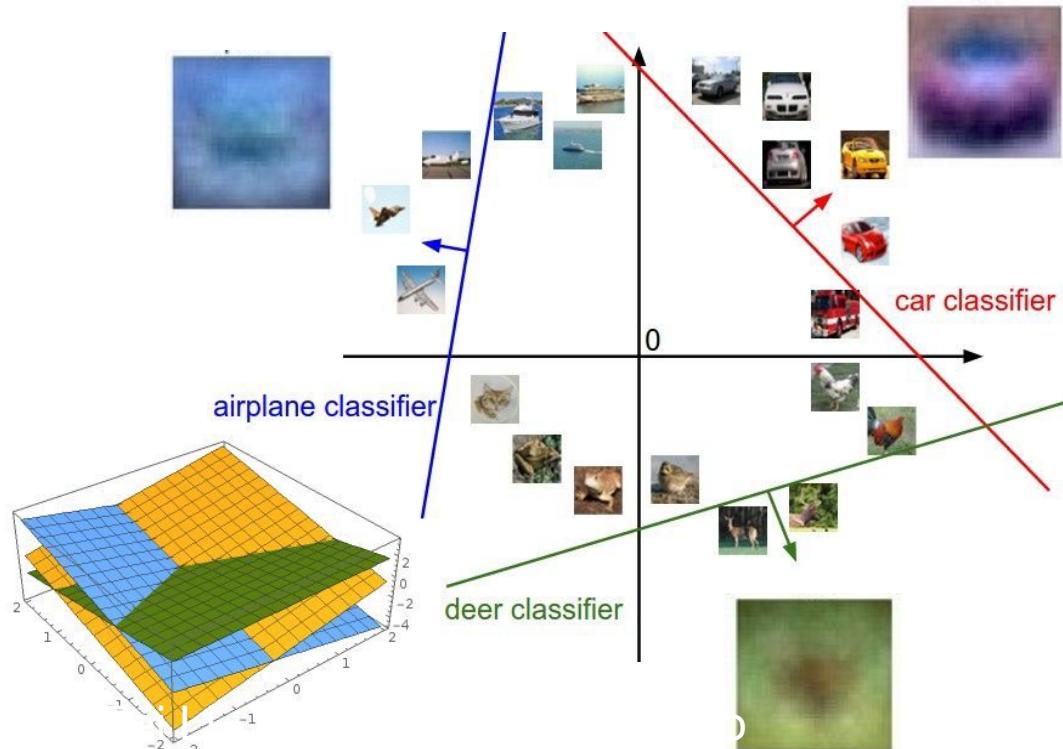
frog

horse
92

ship
March

truck
31, 202

Interpreting a Linear Classifier: Geometric Viewpoint



$$f(x, W) = Wx + b$$



Array of **32x32x3** numbers
(3072 numbers total)

Plot created using [Wolfram Cloud](#)

Cat image by [Nikita](#) is licensed under [CC-BY 2.0](#)

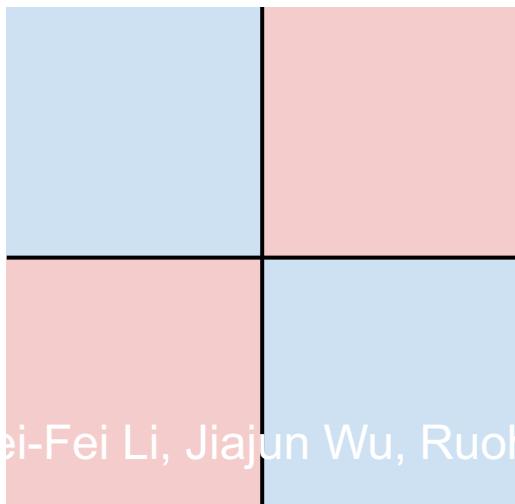
Hard cases for a linear classifier

Class 1:

First and third quadrants

Class 2:

Second and fourth quadrants



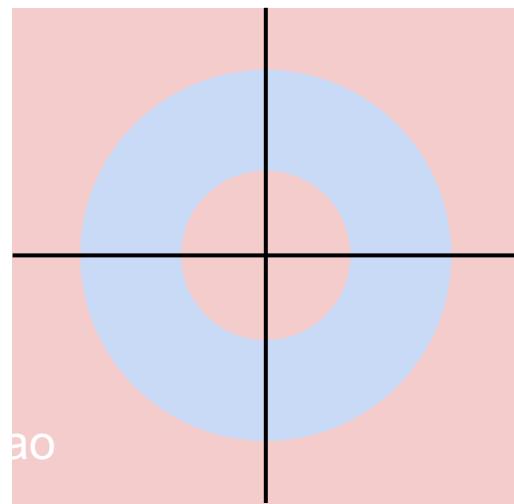
Fei-Fei Li, Jiajun Wu, Ruohao

Class 1:

$1 \leq L_2 \text{ norm} \leq 2$

Class 2:

Everything else



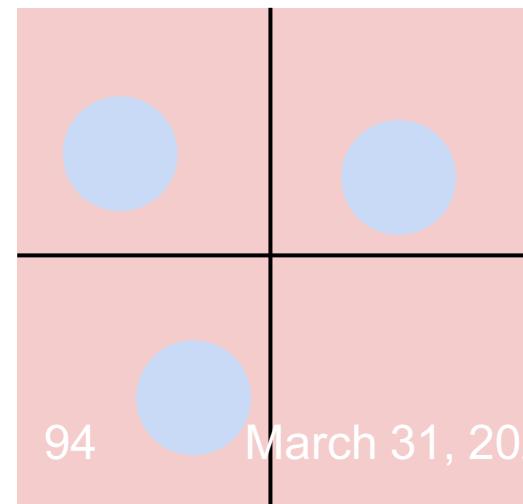
ao

Class 1:

Three modes

Class 2:

Everything else



94

March 31, 202

Linear Classifier – Choose a good W



airplane	-3.45	-0.51	3.42
automobile	-8.87	6.04	4.64
bird	0.09	5.31	2.65
cat	2.9	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	-4.34
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

TODO:

1. Define a **loss function** that quantifies our unhappiness with the scores across the training data.
2. Come up with a way of efficiently finding the parameters that minimize the loss function. (**optimization**)

Suppose: 3 training examples, 3 classes.

With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
-----	------------	-----	-----

car	5.1	4.9	2.5
-----	-----	------------	-----

frog	-1.7	2.0	-3.1
------	------	-----	-------------

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



A **loss function** tells how good our current classifier is

cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

A **loss function** tells how good our current classifier is

Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where x_i is image and
 y_i is (integer) label

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

A **loss function** tells how good our current classifier is

Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where x_i is image and
 y_i is (integer) label

Loss over the dataset is a average of loss over examples:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

Nearest neighbor vs. linear classifiers

- NN pros:

- Simple to implement
- Decision boundaries not necessarily linear
- Works for any number of classes
- *Nonparametric* method

- NN cons:

- Need good distance function
- Slow at test time

- Linear pros:

- Low-dimensional *parametric* representation
- Very fast at test time

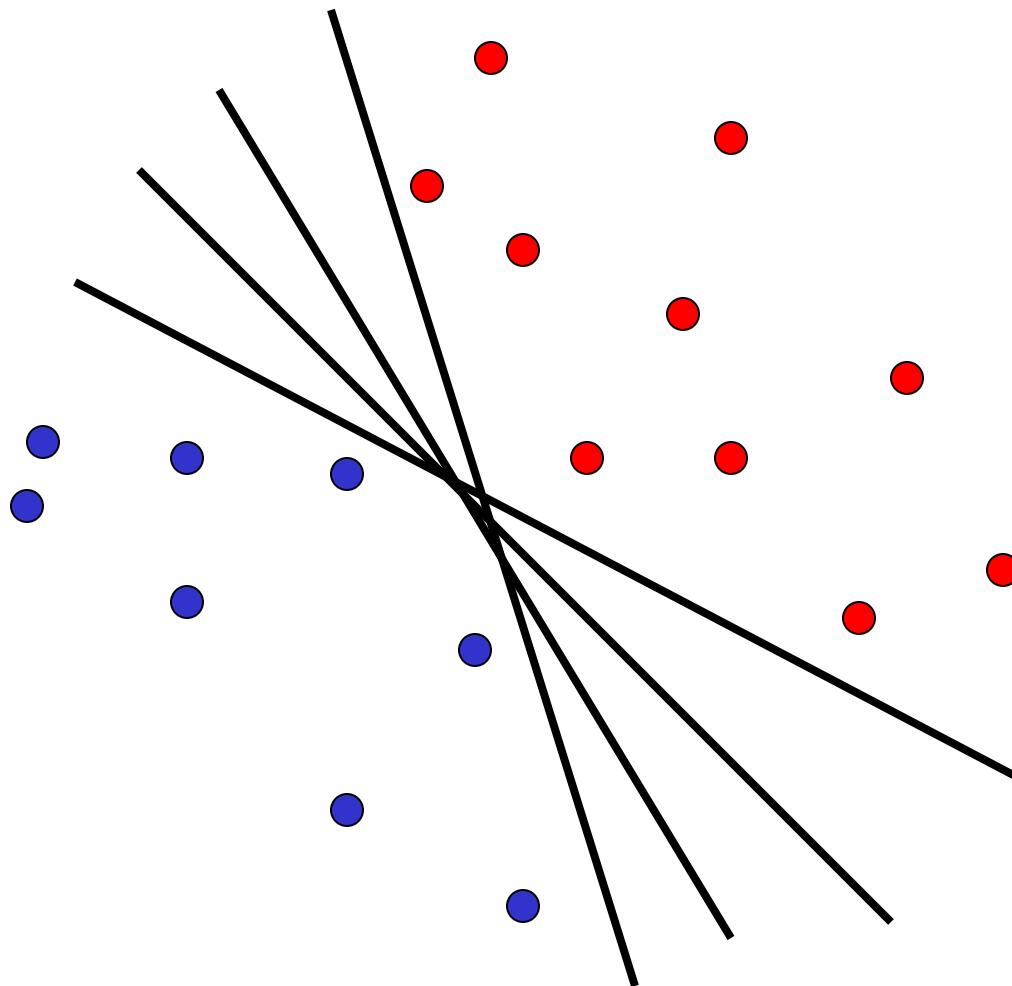
- Linear cons:

- Works for two classes
- How to train the linear function?  激活函数
- What if data is not linearly separable?

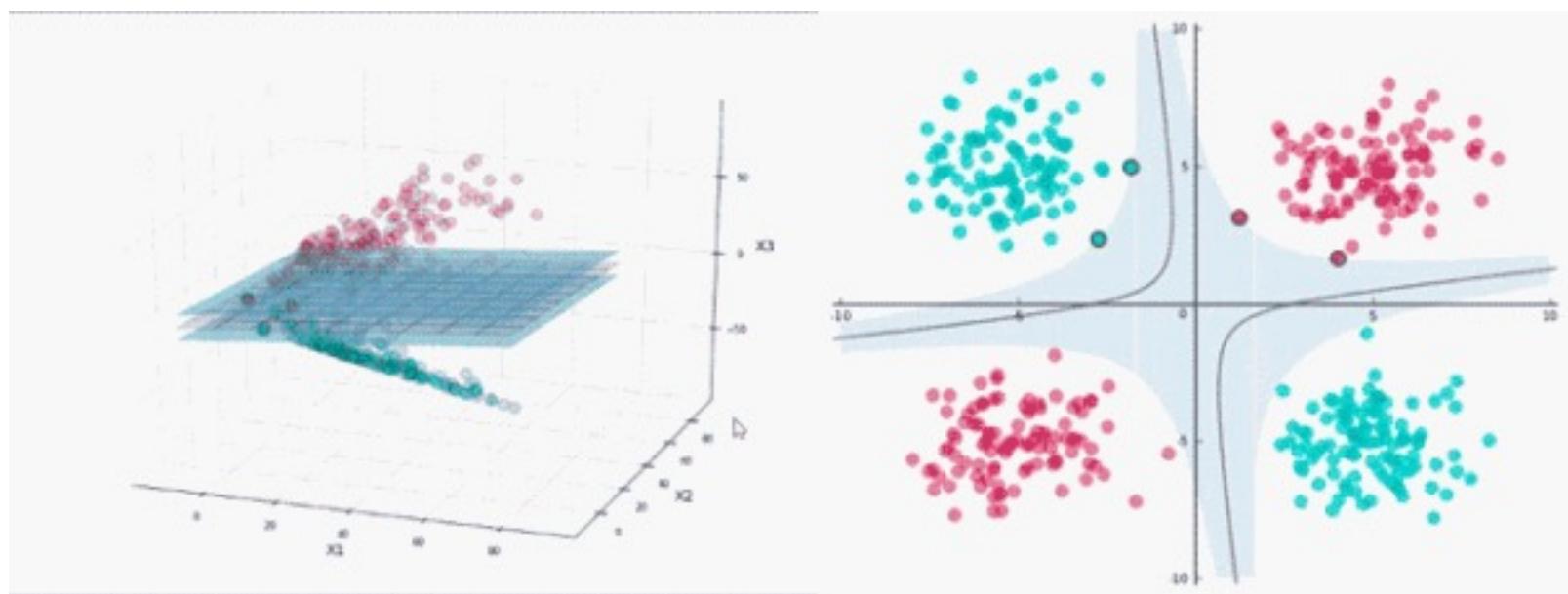
Support vector machines

① 找超平面切兩類
② 讓其與最近點最遠

- When the data is linearly separable, there may be more than one separator (hyperplane)

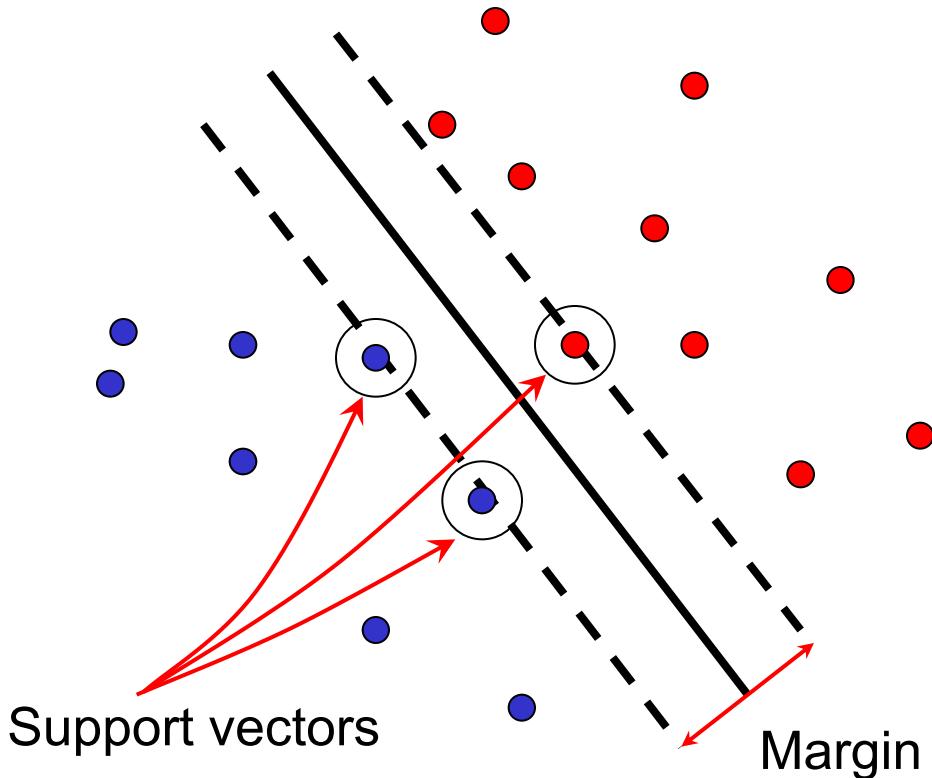


Which separator
is best?



Support vector machines

- Find hyperplane that maximizes the *margin* between the positive and negative examples



$$\mathbf{x}_i \text{ positive } (y_i = 1) : \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1) : \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

For support vectors, $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$

Distance between point and hyperplane:
$$\frac{|\mathbf{x}_i \cdot \mathbf{w} + b|}{\|\mathbf{w}\|}$$

Therefore, the margin is $2 / \|\mathbf{w}\|$

Finding the maximum margin hyperplane

1. Maximize margin $2 / \|\mathbf{w}\|$
2. Correctly classify all training data:

$$\left\{ \begin{array}{ll} \mathbf{x}_i \text{ positive } (y_i = 1) : & \mathbf{x}_i \cdot \mathbf{w} + b \geq 1 \\ \mathbf{x}_i \text{ negative } (y_i = -1) : & \mathbf{x}_i \cdot \mathbf{w} + b \leq -1 \end{array} \right.$$

Quadratic optimization problem:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$

SVM parameter learning

- Separable data: $\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$ subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$

A red curly brace positioned below the term $\frac{1}{2} \|\mathbf{w}\|^2$.

Maximize
margin

A red curly brace positioned below the term $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1$.

Classify training data correctly

- Non-separable data:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b))$$

A red curly brace positioned below the term $\frac{1}{2} \|\mathbf{w}\|^2$.

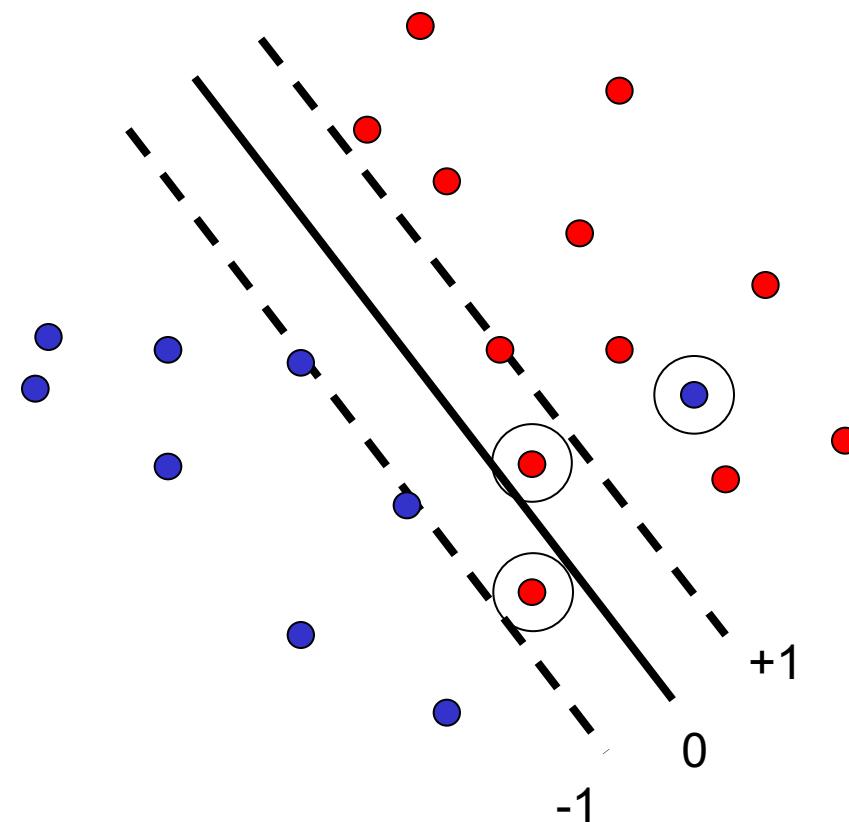
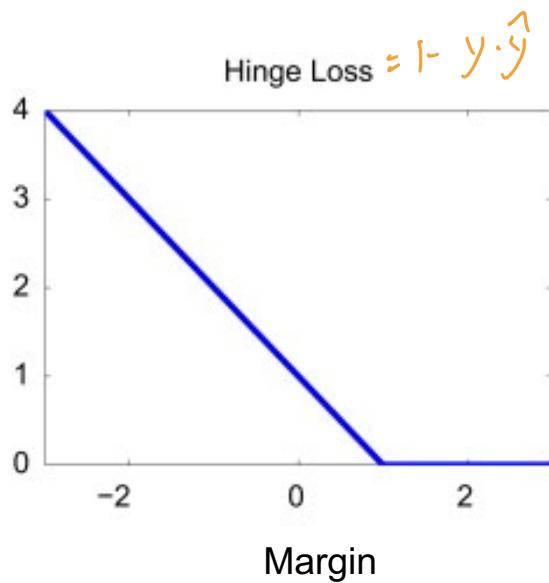
Maximize
margin

A red curly brace positioned below the term $\max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b))$.

Minimize classification mistakes

SVM parameter learning

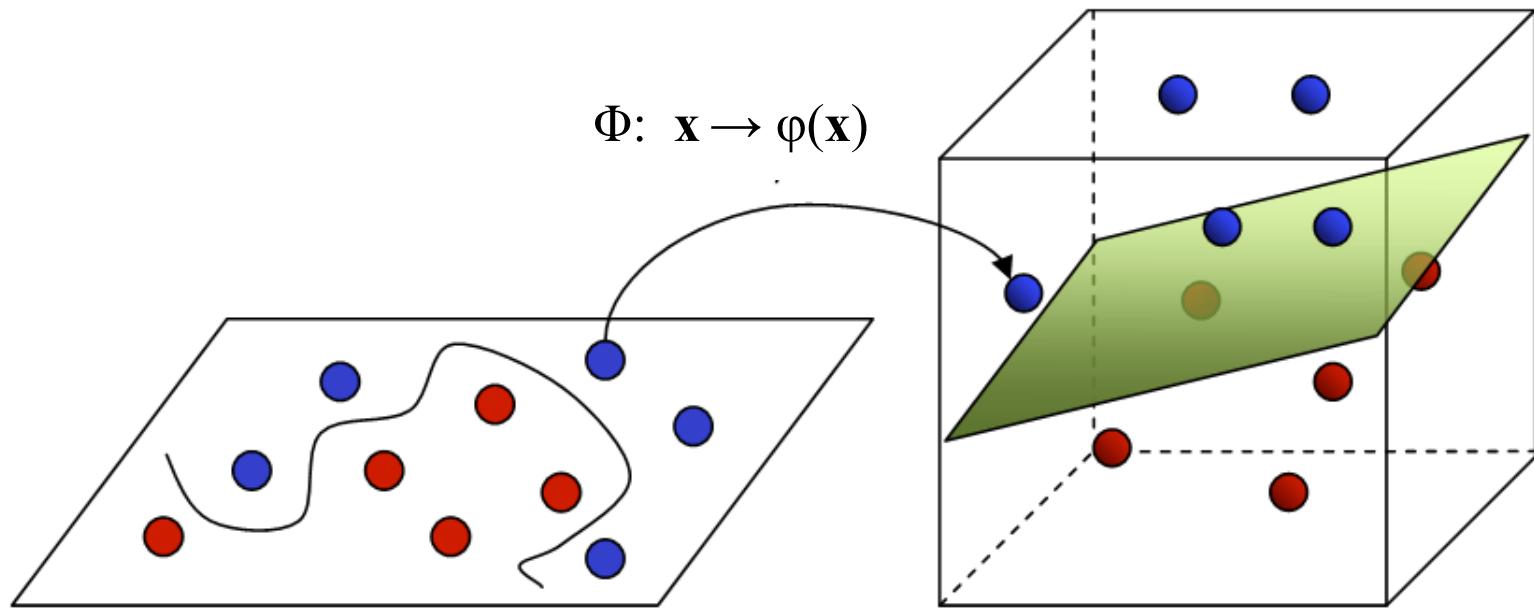
$$\min_{\mathbf{w}, b} \quad \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b))$$



Demo: <http://cs.stanford.edu/people/karpathy/svmjs/demo>

Nonlinear SVMs

- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable



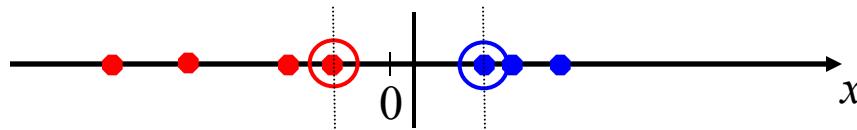
Input Space

Feature Space

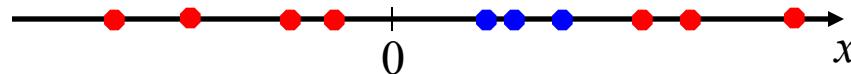
[Image source](#)

Nonlinear SVMs

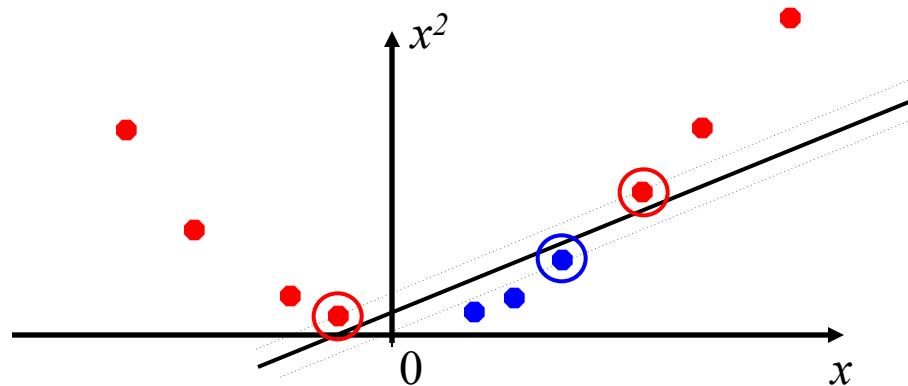
- Linearly separable dataset in 1D:



- Non-separable dataset in 1D:



- We can map the data to a *higher-dimensional space*:



The kernel trick

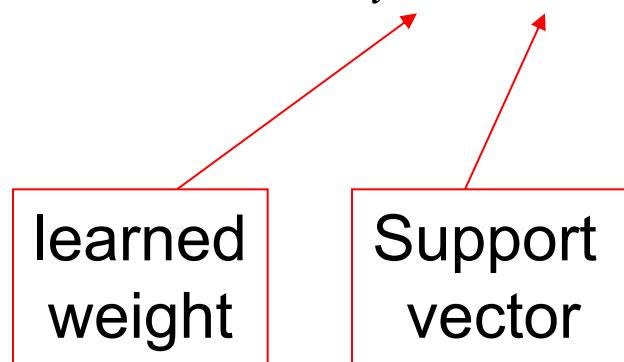
- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable
- **The kernel trick:** instead of explicitly computing the lifting transformation $\varphi(\mathbf{x})$, define a kernel function K such that

$$K(\mathbf{x}, \mathbf{y}) = \varphi(\mathbf{x}) \cdot \varphi(\mathbf{y})$$

The kernel trick

- Linear SVM decision function:

$$\mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$$



The kernel trick

- Linear SVM decision function:

$$\mathbf{w} \cdot \mathbf{x} + b = \sum_i \underline{\alpha}_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$$

lalrange 乘子
learnt para

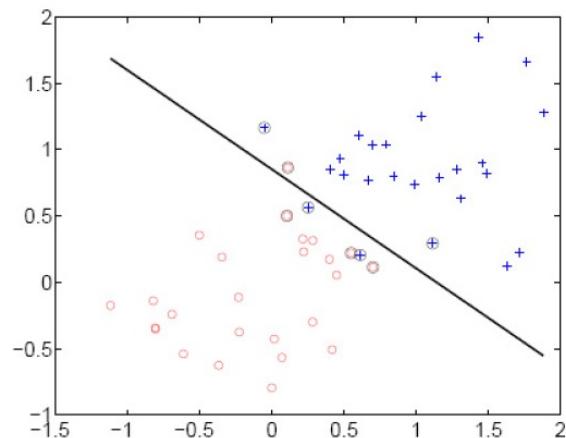
- Kernel SVM decision function:

$$\sum_i \alpha_i y_i \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}) + b = \sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

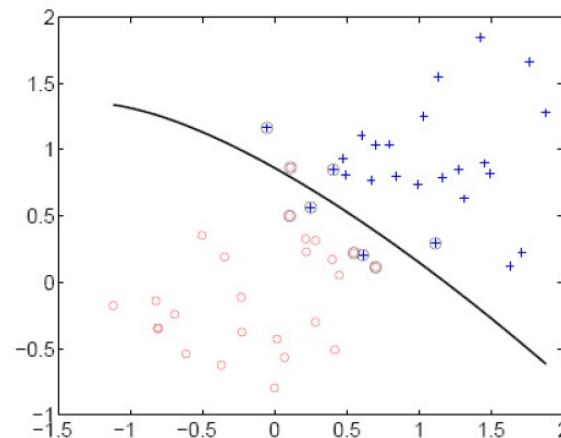
核函数代替内积

- This gives a nonlinear decision boundary in the original feature space

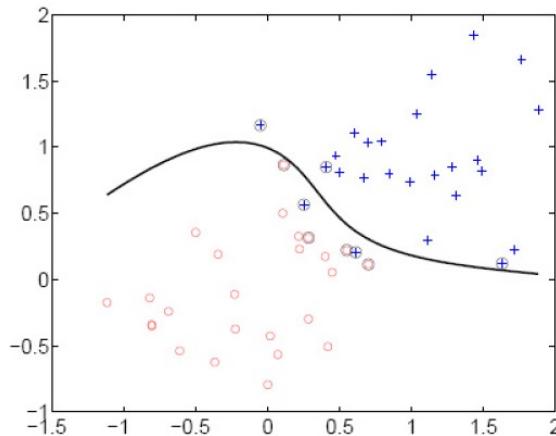
Polynomial kernel: $K(\mathbf{x}, \mathbf{y}) = (c + \mathbf{x} \cdot \mathbf{y})^d$



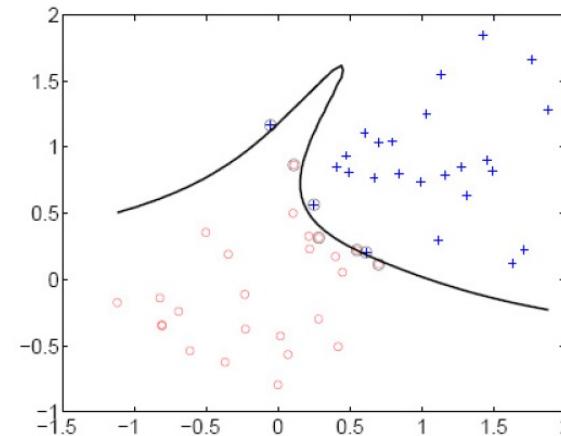
linear



2nd order polynomial



4th order polynomial

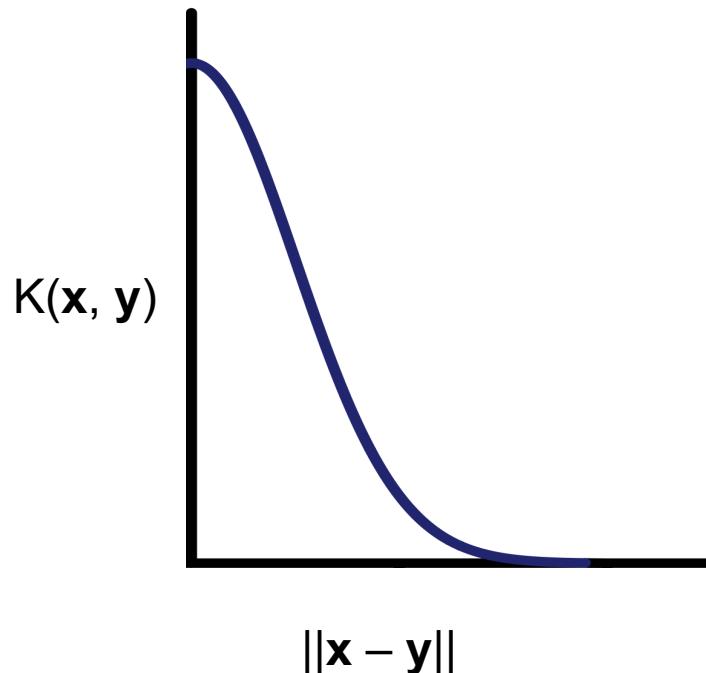


8th order polynomial

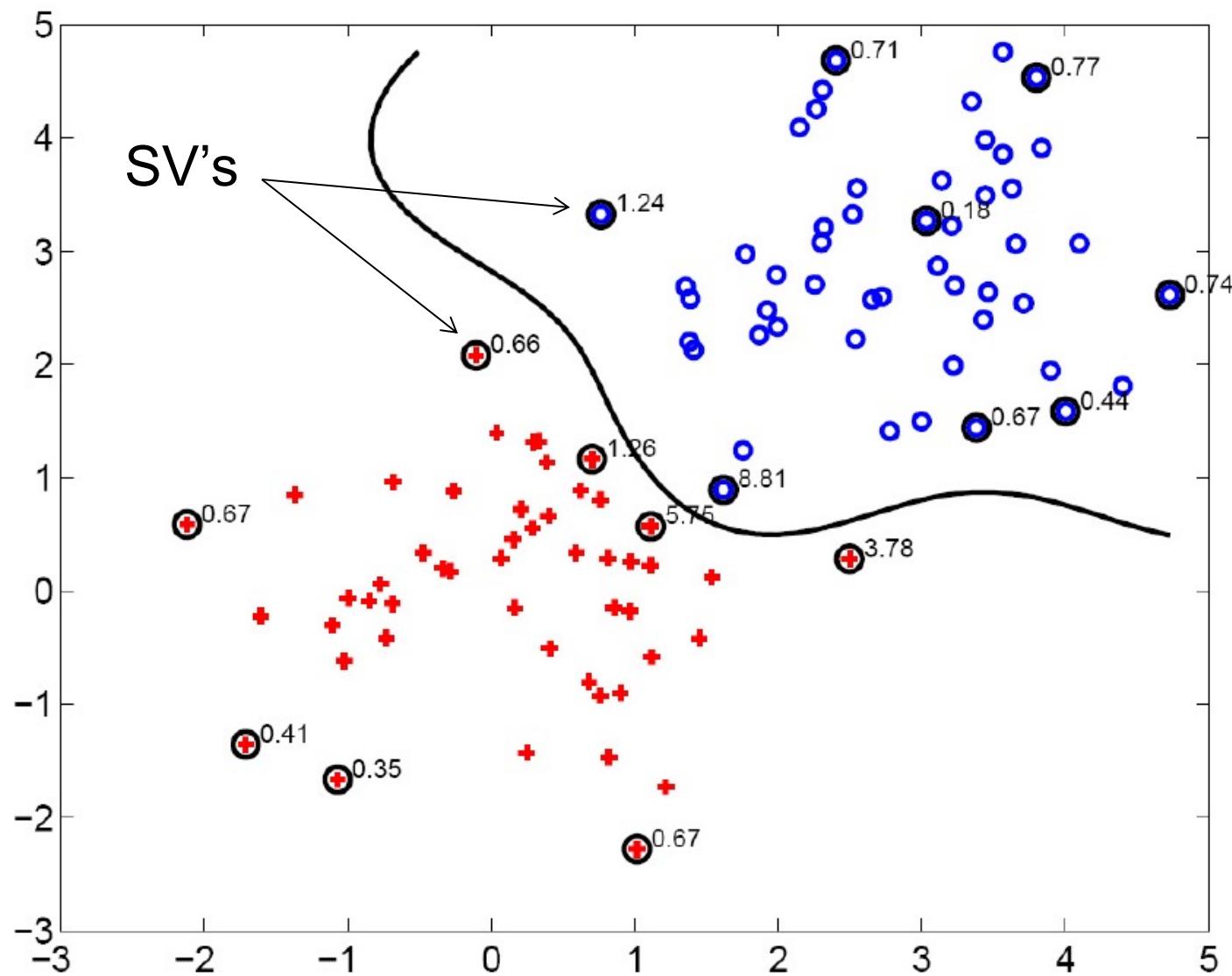
Gaussian kernel

- Also known as the radial basis function (RBF) kernel:

$$K(\mathbf{x}, \mathbf{y}) = \exp\left(-\frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{y}\|^2\right)$$



Gaussian kernel



SVMs: Pros and cons

- **Pros**

- Kernel-based framework is very powerful, flexible
- Training is convex optimization, globally optimal solution can be found
- Amenable to theoretical analysis
- SVMs work very well in practice, even with very small training sample sizes

- **Cons**

一个SVM 只能二分类

- No “direct” multi-class SVM, must combine two-class SVMs (e.g., with one-vs-others)
- Computation, memory (esp. for nonlinear SVMs)

Generalization

- Generalization refers to the ability to correctly classify never before seen examples
- Can be controlled by turning “knobs” that affect the complexity of the model



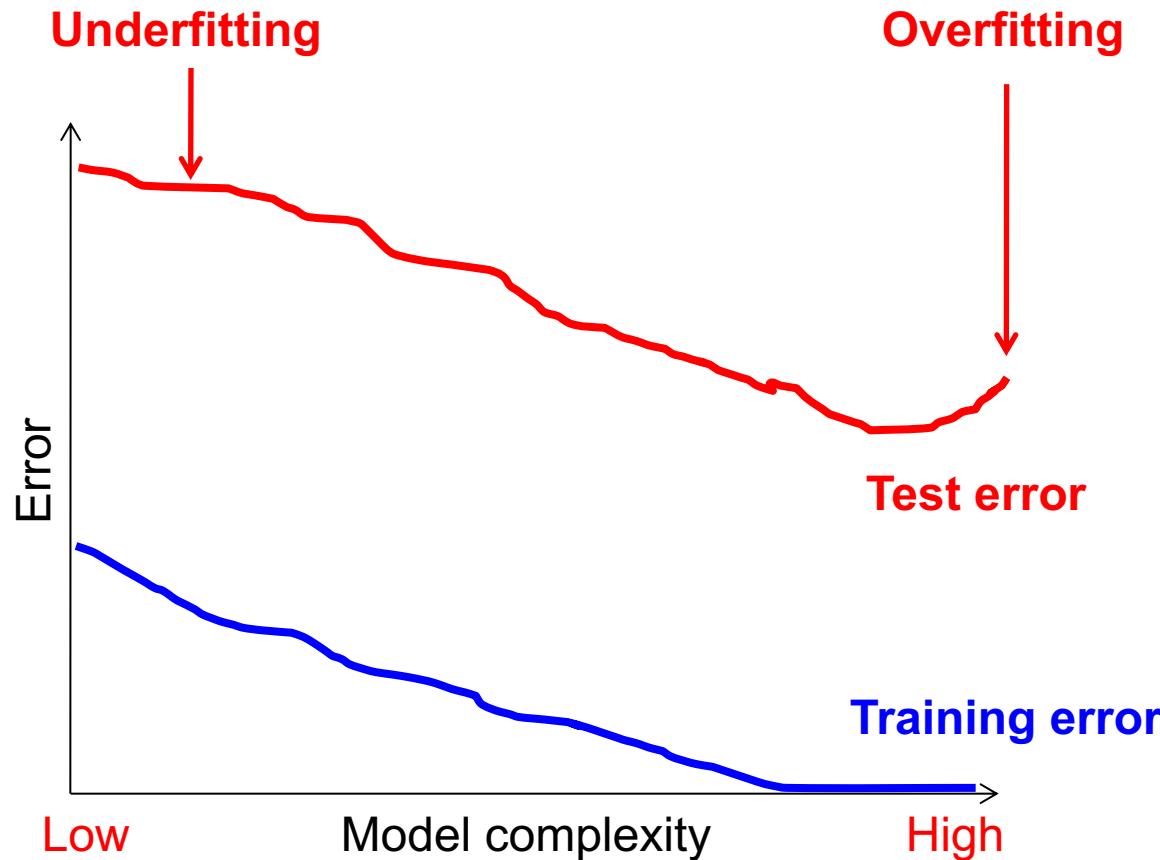
Training set (labels known)



Test set (labels unknown)

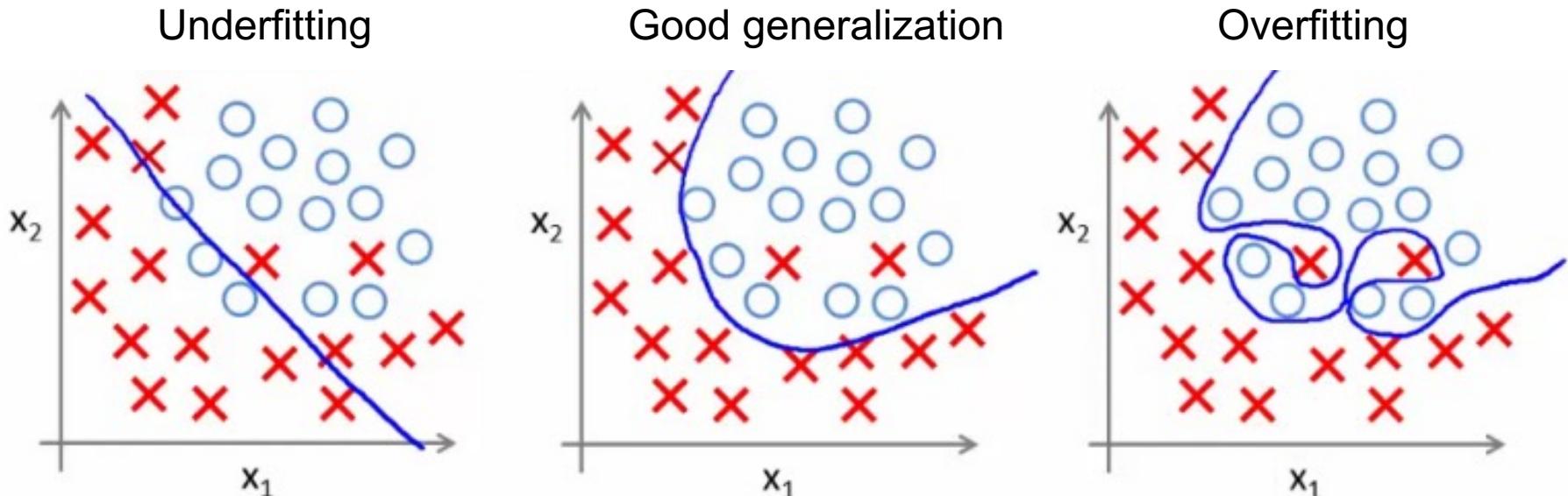
Diagnosing generalization ability

- **Training error:** how does the model perform on the data on which it was trained?
- **Test error:** how does it perform on never before seen data?

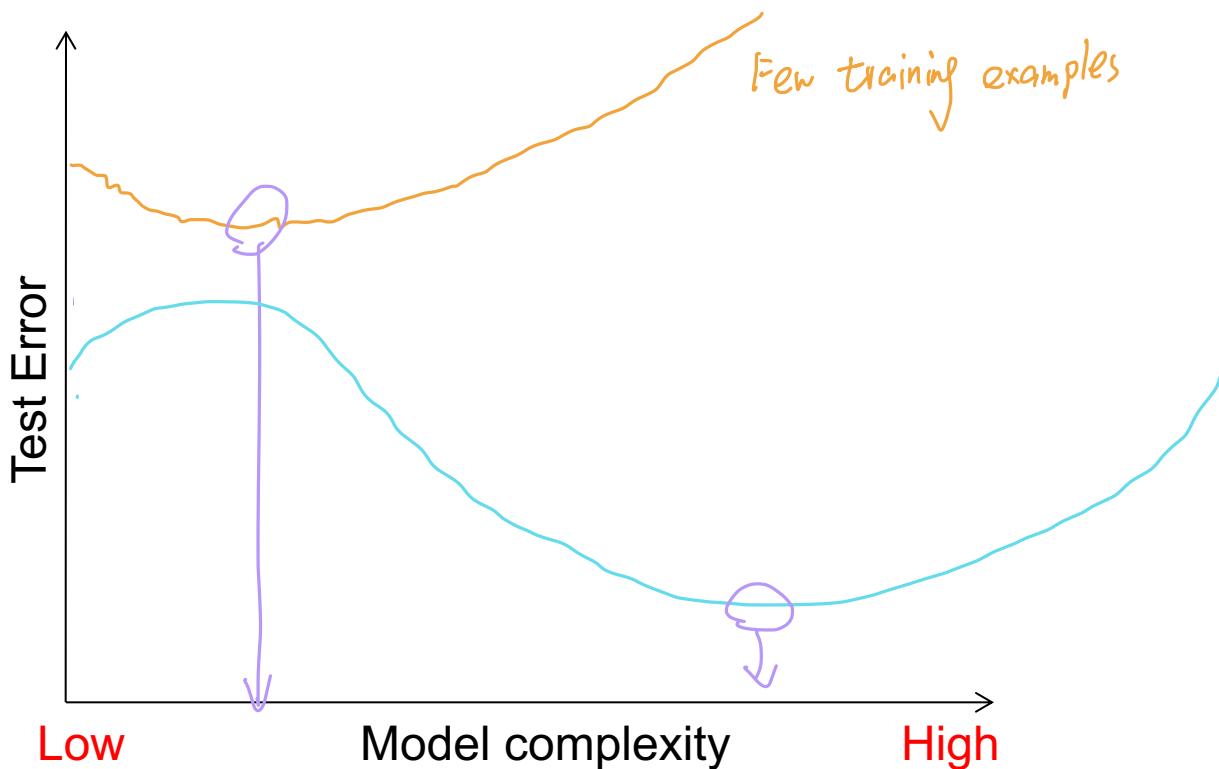


Underfitting and overfitting

- **Underfitting:** training and test error are both *high*
 - Model does an equally poor job on the training and the test set
 - Either the training procedure is ineffective or the model is too “simple” to represent the data
- **Overfitting:** Training error is *low* but test error is *high*
 - Model fits irrelevant characteristics (noise) in the training data
 - Model is too complex or amount of training data is insufficient



Effect of training set size



Validation

- Split the data into **training**, **validation**, and **test** subsets
- Use training set to **optimize model parameters**
- Use validation test to **choose the best model**
- Use test set only to **evaluate performance**

