

#### THE UNIVERSITY OF BRITISH COLUMBIA

# **CPSC 425: Computer Vision**



Lecture 28: Image Classification (final)

## Menu for Today (November 18, 2020)

#### **Topics:**

- Boosting
- Sliding Window Object Detection

#### **Redings:**

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

#### **Reminders:**

- Friendly Assignment 5 competition
- Quiz 5 will be due end-of-day on Friday



#### Face detection with Viola & Jones - Object Proposals

# Forsyth & Ponce (2nd ed.) 17.1–17.2

— Assignment 5: Scene Recognition with Bag of Words due November 20



### Today's "fun" Example:

# Wait for it ... :)

## Lecture 27: Re-cap Decision Trees

#### A random forest is an ensemble of decision trees.

Randomness is incorporated via training set sampling and/or generation of the candidate binary tests

The prediction of the random forest is obtained by averaging over all decision trees.



Forsyth & Ponce (2nd ed.) Figure 14.19. Original credit: J. Shotton et al., 2011 4









Figure credit: J. Shotton et al., 2011



$$f_{\theta}(I, \mathbf{x}) = d_I \left( \mathbf{x} + \frac{1}{d} \right)$$



 $\left( \frac{\mathbf{u}}{d_I(\mathbf{x})} \right) - d_I \left( \mathbf{x} + \frac{\mathbf{v}}{d_I(\mathbf{x})} \right)$ 

Figure credit: J. Shotton et al., 2011

 $f_{\theta}(I, \mathbf{x}) > \Theta_j$ 



 $f_{\theta}(I, \mathbf{x}) = d_I \left( \mathbf{x} + \frac{\mathbf{u}}{d_I(\mathbf{x})} \right) - d_I \left( \mathbf{x} + \frac{\mathbf{v}}{d_I(\mathbf{x})} \right)$ 

7





.... ...

 $f_{\theta}$ 

 $f_{\theta}(I, \mathbf{x}) > \Theta_j$ 



$$q(I, \mathbf{x}) = d_I \left( \mathbf{x} + \frac{\mathbf{u}}{d_I(\mathbf{x})} \right) - d_I \left( \mathbf{x} + \frac{\mathbf{v}}{d_I(\mathbf{x})} \right)$$





 $f_{\theta}$ 

••••

. . .

 $f_{\theta}(I, \mathbf{x}) > \Theta_j$ 



$$q(I, \mathbf{x}) = d_I \left( \mathbf{x} + \frac{\mathbf{u}}{d_I(\mathbf{x})} \right) - d_I \left( \mathbf{x} + \frac{\mathbf{v}}{d_I(\mathbf{x})} \right)$$





•••• ••• •••

 $f_{\theta}$ 

 $f_{\theta}(I, \mathbf{x}) > \Theta_j$ 



$$q(I, \mathbf{x}) = d_I \left( \mathbf{x} + \frac{\mathbf{u}}{d_I(\mathbf{x})} \right) - d_I \left( \mathbf{x} + \frac{\mathbf{v}}{d_I(\mathbf{x})} \right)$$





$$q(I, \mathbf{x}) = d_I \left( \mathbf{x} + \frac{\mathbf{u}}{d_I(\mathbf{x})} \right) - d_I \left( \mathbf{x} + \frac{\mathbf{v}}{d_I(\mathbf{x})} \right)$$





# Combining **Classifiers**

One common strategy to obtain a better classifier is to combine multiple classifiers.

A simple approach is to train an ensemble of independent classifiers, and average their predictions.

**Boosting** is another approach.

— Train an ensemble of classifiers sequentially.

 Bias subsequent classifiers to correctly predict training examples that previous classifiers got wrong.

- The final boosted classifier is a weighted combination of the individual classifiers.











Final classifier is a combination of weak classifiers





#### THE UNIVERSITY OF BRITISH COLUMBIA

# **CPSC 425: Computer Vision**



Lecture 28: Object Detection

19

# **Object Detection**: Introduction

We have been discussing image classification, where we pass a whole image into a classifier and obtain a class label as output

We assumed the image contained a single, central object

object class in an image

- The task of **object detection** is to detect and localize all instances of a target
- Localization typically means putting a tight bounding box around the object

Train an image classifier as described previously. 'Slide' a fixed-sized detection window across the image and evaluate the classifier on each window.



Train an image classifier as described previously. 'Slide' a fixed-sized detection window across the image and evaluate the classifier on each window. Is there a car?

# 



Train an image classifier as described previously. 'Slide' a fixed-sized detection window across the image and evaluate the classifier on each window.

#### Is there a car?



Train an image classifier as described previously. 'Slide' a fixed-sized detection window across the image and evaluate the classifier on each window.



#### Is there a car?

Train an image classifier as described previously. 'Slide' a fixed-sized detection window across the image and evaluate the classifier on each window. Is there a car?



Train an image classifier as described previously. 'Slide' a fixed-sized detection window across the image and evaluate the classifier on each window. Is there a car?

# 



Train an image classifier as described previously. 'Slide' a fixed-sized detection window across the image and evaluate the classifier on each window.



#### Is there a car?

Train an image classifier as described previously. 'Slide' a fixed-sized detection window across the image and evaluate the classifier on each window.



Is there a car?

Train an image classifier as described previously. 'Slide' a fixed-sized detection window across the image and evaluate the classifier on each window.



Is there a car?

Train an image classifier as described previously. 'Slide' a fixed-sized detection window across the image and evaluate the classifier on each window.



This is a search over location — We have to search over scale as well — We may also have to search over aspect ratios

# What data we **train** a classifier on? Image Classifiers





# **Image** classifiers can be applied to regions/windows, but do not work so well in practice ...

# What data we **train** a classifier on? Image Classifiers





# What data we **train** a classifier on? Image Classifiers





#### **Object** Classifiers



















## What data we train a classifier on?

**Object** classifiers work a lot better ... but require expensive bounding box annotations ...

#### **Object** Classifiers



















# Today's "fun" Example: Detection with No Data

#### **Expensive** — do this for few classes (e.g., 20)



#### **Cheap** – do this for many classes (e.g., 80 or 20,000)

LSDA: Large Scale Detection through Adaptation (Hoffman et al, 2014)



# Today's "fun" Example: Detection with No Data

Our **UniT**: Unified Knowledge Transfer

Method

LSDA [22] LSDA+Semantic [52] LSDA+MIL [23] DOCK [26]

UniT (Ours)

Full Supervision [26]

$AP_{50}$	$AP_S$	$\operatorname{AP}_M$	$AP_L$
4.6	1.2	5.1	7.8
4.7	1.1	5.1	8.0
5.9	1.5	8.3	10.7
14.4	2.0	12.8	24.9
16.7	3.2	16.6	27.3
25.2	5.8	26.0	41.6
# Today's "fun" Example: Detection with No Data

### Our **UniT**: Unified Knowledge Transfer



(Note, above result is obtained without human ever annotating bounding boxes or pixel-level segments for any of these objects)



### Let's assume we have **object** labeled data ...

**Object** classifiers work a lot better ... but require expensive bounding box annotations ...

### **Object** Classifiers



















### The **Viola-Jones** face detector is a classic sliding window detector that learns both efficient features and a classifier

- A key strategy is to use features that are fast to evaluate to reject most windows early

The Viola-Jones detector computes 'rectangular' features within each window

A 'rectangular' feature is computed by summing up pixel values within rectangular regions and then differencing those region sums



Figure credit: P. Viola and M. Jones, 2001

A 'rectangular' feature is computed by summing up pixel values within rectangular regions and then differencing those region sums



Figure credit: P. Viola and M. Jones, 2001

original image

 $A(x,y) = \sum_{x' \le x, y' \le y} I(x',y')$ 

A(x,y)			
1	6	8	
3	12	15	
5	15	19	

### integral image

What is the sum of the bottom right 2x2 square?



original image

A(x,y)		
1	6	8
3	12	15
5	15	19

### integral image

What is the sum of the bottom right 2x2 square?



original image

A(x,y)		
1	6	8
3	12	15
5	15	19

### integral image

What is the sum of the bottom right 2x2 square?

$$I(x, y)$$
  
1 5 2  
2 4 1  
2 1 1

original image

A(1, 1, 3, 3) = A(3, 3) -4 0



integral image

$$egin{array}{rl} A(1,3) - A(3,1) + A(1,1) \ & 8 & - & 5 & + & 1 \end{array}$$

Image Credit: Ioannis (Yannis) Gkioulekas (CMU)

45



original image

Can find the **sum** of any block using **3** operations

$$A(x_1, y_1, x_2, y_2) = A(x_2, y_2) - A(x_1, y_2) - A(x_2, y_1) + A(x_1, y_1)$$

 $A(x,y) = \sum I(x',y')$  $x' \leq x, y' \leq y$ 



integral image

Given an integral image, the sum within a rectangular region in I can be computed with just 3 additions/subtractions



### Sum = A - B - C + D



**Constant time:** does not depend on the size of the region. We can avoid scaling images - just scale features directly (remember template matching!)

Figure credit: P. Viola

### **Integral Image Layer for Deep Neural Networks**

In a classical paper [1] from 2001, Viola and Jones popularized the use of large rectangular image filters in order to obtain features for image recognition. The use of very large filters allowed Viola and Jones to compute features over very large receptive fields without blowing up the computation cost. For the next 10+ years, such features remained the staple of fast computer vision (e.g. [2]). The advent of deep learning made the use of integral-image features far less popular. Currently, state-of-the-art architectures invariably relying on very deep architectures. In these architectures sufficiently large receptive fields are obtained via the use of downsampling with subsequent upsampling [3] or via dilated convolutions [4]. All such tricks however have their downsides and usually necessitate the use of very deep networks.

The goal of this project is to implement an integral image-based filtering as a layer for deep architectures in Torch deep learning package, and to evaluate it for the task of learning very fast object detectors (as an alternative to e.g. [5]) and semantic segmentation systems (as an alternative to e.g. [3,4]). The hope is to obtain much shallower architectures, which at least for simple classes (e.g. roadsigns or upright pedestrians) will approach the performance of much deeper ones.

The project is supervised by Victor Lempitsky at Skoltech, Moscow, Russia.

https://github.com/shrubb/integral-layer

### **Deep Neural Networks for Object Detection**

Christian Szegedy Alexander Toshev Dumitru Erhan Google, Inc. {szegedy, toshev, dumitru}@google.com

Deep Neural Networks (DNNs) have recently shown outstanding performance on image classification tasks [14]. In this paper we go one step further and address the problem of object detection using DNNs, that is not only classifying but also precisely localizing objects of various classes. We present a simple and yet powerful formulation of object detection as a regression problem to object bounding box masks. We define a multi-scale inference procedure which is able to produce high-resolution object detections at a low cost by a few network applications. State-of-the-art performance of the approach is shown on Pascal VOC.

### Abstract

### Training Dataset:

 $(x_1, 1)$  $(x_2, 1)$ 



Faces



**Not-faces** 

Evaluate a Harr Wavelet filter on each training example



 $(x_1, 1)$  $(x_2, 1)$ 



Faces

### $(x_3, 0)$ $(x_4, 0)$ $(x_5, 0)$ $(x_6, 0)$



**Not-faces** 

Evaluate a Harr Wavelet filter on each training example

	$(x_1, 1)$	$(x_2, 1)$
$\alpha_1$		
	0.8	0.7

We can build a simple classifier by just selecting a threshold on the filter response (e.g. Harr filter response > 0.6 = face; Harr filter response <= 0.6 = not face)



### Note: it is easy to find an optimal threshold. Just requires linear search over training example responses.



Evaluate a Harr Wavelet filter on each training example

	$(x_1, 1)$	$(x_2, 1)$
$\alpha_1$		
	0.8	0.7

Weak classifier 
$$h_j(x) = \begin{cases} 1 & \text{if } \\ 0 & 0 \end{cases}$$

 $(x_3, 0)$  $(x_4, 0)$   $(x_5, 0)$  $(x_6, 0)$ 

 $\cdots (x_n, y_n)$ 

0.2

0.3

0.8

0.1

 $f_j(x) > \theta_j$ 

threshold

Evaluate a Harr Wavelet filter on each training example

 $(x_1, 1)$ 



Faces

### $(x_2,1)$ $(x_3,0)$ $(x_4,0)$ $(x_5,0)$ $(x_6, 0)$



**Not-faces** 

Evaluate a Harr Wavelet filter on each training example

 $(x_1, 1)$ 



Faces

### Note: we can easily compare different Harr Wavelet features under their individual best thresholds to see which is the most informative (has highest classification)

### $(x_2,1)$ $(x_3,0)$ $(x_4,0)$ $(x_5,0)$ $(x_6, 0)$



**Not-faces** 







**Figure credit**: B. Freeman Many possible rectangular features (180,000+ were used in the original paper)

Evaluate a Harr Wavelet filter on each training example

$$(x_1, 1)$$
  $(x_2, 1)$ 



Faces

### Note: we can easily compare different Harr Wavelet features under their individual best thresholds to see which is the most informative (has highest classification)

However, no one feature is likely to be good enough

# $(x_3,0)$ $(x_4,0)$ $(x_5,0)$ $(x_6,0)$



**Not-faces** 





feature against a threshold







### Use **boosting** to both select the informative features and form the classifier. Each round chooses a weak classifier that simply compares a single rectangular

Figure credit: P. Viola and M. Jones, 2001

1. Select best filter/threshold combination

a. Normalize the weights

b. For each feature, j

c. Choose the classifier,  $h_t$  with the lowest error  $\varepsilon$ 

2. Re-weight examples

$$W_{t+1,i} = W_{t,i}\beta_t^{1-|h_t(x_i)-y_i|} \qquad \beta_t = \frac{\varepsilon_t}{1-\varepsilon_t}$$

$$w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^{n} w_{t,j}}$$

$$h_j(x) = \begin{cases} 1 & \text{if } f_j(x) > \theta_j \\ 0 & \text{otherwise} \end{cases}$$

$$\varepsilon_j = \sum_i w_i \left| h_j(x_i) - y_i \right|$$

1. Select best filter/threshold combination

a. Normalize the weights

$$w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^{n} w_{t,j}}$$

### We start with all sample weights = 1

1. Select best filter/threshold combination

a. Normalize the weights

b. For each feature, j

**Note:** the second term is 0/1

$$w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^{n} w_{t,j}}$$

$$h_j(x) = \begin{cases} 1 & \text{if } f_j(x) > \theta_j \\ 0 & \text{otherwise} \end{cases}$$

$$\varepsilon_j = \sum_i w_i \left| h_j(x_i) - y_i \right|$$
Weighed sum of miss-classified training examples

 O predicted label and true label are same 1 predicted label and true label are different (error)

1. Select best filter/threshold combination

a. Normalize the weights

b. For each feature, j

c. Choose the classifier,  $h_t$  with the lowest error  $\mathcal{E}$ 

$$w_{t,i} \leftarrow \frac{w_{t,i}}{\sum_{j=1}^{n} w_{t,j}}$$

$$h_j(x) = \begin{cases} 1 & \text{if } f_j(x) > \theta_j \\ 0 & \text{otherwise} \end{cases}$$

$$\varepsilon_j = \sum_i w_i \left| h_j(x_i) - y_i \right|$$

1. Select best filter/threshold combination

a. Normalize the weights

b. For each feature, j

c. Choose the classifier,  $h_t$  with the lowest error  $\varepsilon$ 

2. Re-weight examples

$$W_{t+1,i} = W_{t,i} \beta_t^{1-|h_t(x_i)-y_i|} \qquad \beta_t = \frac{\varepsilon_t}{1-\varepsilon_t}$$

**Case 1**: Classification for the sample i is **correct** 

 $\mathbf{w}_{t+1,i} = \mathbf{w}_{t,i} \ \beta_t$ 

**Case 2**: Classification for the sample i is **incorrect**  $\mathbf{w}_{t+1,i} = \mathbf{w}_{t,i}$ 

2. Re-weight examples

$$W_{t+1,i} = W_{t,i} \beta_t^{1-|h_t(x_i)-y_i|} \qquad \beta_t = \frac{\varepsilon_t}{1-\varepsilon_t}$$

**Case 1**: Classification for the sample i is **correct** 

 $\mathbf{w}_{t+1,i} = \mathbf{w}_{t,i} \ \beta_t$ 

**Case 2**: Classification for the sample i is **incorrect**  $\mathbf{w}_{t+1,i} = \mathbf{w}_{t,i}$ 

2. Re-weight examples

$$w_{t+1,i} = w_{t,i} \beta_t^{1-|h_t(x_i)-y_i|}$$

### **Note**: the Beta is < 1

$$\beta_t = \frac{\varepsilon_t}{1 - \varepsilon_t}$$

1. Select best filter/threshold combination

a. Normalize the weights

b. For each feature, j

c. Choose the classifier,  $h_t$  with the lowest error  $\varepsilon$ 

2. Re-weight examples

$$W_{t+1,i} = W_{t,i}\beta_t^{1-|h_t(x_i)-y_i|} \qquad \beta_t = \frac{\varepsilon_t}{1-\varepsilon_t}$$

# Viola & Jones algorithm

3. The final strong classifier is

$$h(x) = \begin{cases} 1 & \sum_{t=1}^{T} \alpha_t h_t(x) \ge \frac{1}{2} \sum_{t=1}^{T} \alpha_t \\ 0 & \text{otherwise} \end{cases} \quad \alpha_t = \log \frac{1}{\beta_t}$$

The final strong classifier is a weighted linear combination of the T weak classifiers where the weights are inversely proportional to the training errors



# **Example**: Face Detection Summary



Figure credit: K. Grauman

# **Example**: Face Detection Summary



Figure credit: K. Grauman

### **Observations**:

- On average only 0.01% of all sub-windows are positive (faces)
- Equal computation time is spent on all sub-window
- Shouldn't we spend most time only on **potentially positive** sub-windows?

### windows are positive (faces) on all sub-window ly on **notentially positive** sub-windo

### **Observations**:

- On average only 0.01% of all sub-windows are positive (faces)
- Equal computation time is spent on all sub-window
- Shouldn't we spend most time only on **potentially positive** sub-windows?

A simple 2-feature classifier can achieve almost 100% detection rate (0% false negatives) with 50% false positive rate

### **Observations**:

- On average only 0.01% of all sub-windows are positive (faces)
- Equal computation time is spent on all sub-window
- Shouldn't we spend most time only on **potentially positive** sub-windows?

### **Solution:**

most negative (clearly non-face) windows

- 2nd layer with 10 features can tackle "harder" negative-windows which survived the 1st layer, and so on...

A simple 2-feature classifier can achieve almost 100% detection rate (0% false negatives) with 50% false positive rate

- A simple 2-feature classifier can act as a 1st layer of a series to filter out
# **Cascading** Classifiers



To make detection **faster**, features can be reordered by increasing complexity of evaluation and the thresholds adjusted so that the early (simpler) tests have few or no false negatives

Any window that is rejected by early tests can be discarded quickly without computing the other features

This is referred to as a **cascade** architecture

## **Cascading** Classifiers



### A classifier in the cascade is not necessarily restricted to a single feature

### **Example**: Face Detection

# Viola & Jones algorithm

3. The final strong classifier is

$$h(x) = \begin{cases} 1 & \sum_{t=1}^{T} \alpha_t h_t(x) \ge \frac{1}{2} \sum_{t=1}^{T} \alpha_t \\ 0 & \text{otherwise} \end{cases} \quad \alpha_t = \log \frac{1}{\beta_t}$$

The final strong classifier is a weighted linear combination of the T weak classifiers where the weights are inversely proportional to the training errors

Image Credit: Ioannis (Yannis) Gkioulekas (CMU)



## **Example**: Face Detection Summary



Figure credit: K. Grauman

# Hard Negative Mining





### **Image From**: Jamie Kang

# **Example**: Face Detection

### Just for fun:



recognition technology"

### "CV Dazzle, a project focused on finding fashionable ways to thwart facial-

Figure source: Wired, 2015





## **Recall:** Sliding Window

Train an image classifier as described previously. 'Slide' a fixed-sized detection window across the image and evaluate the classifier on each window.



Image credit: KITTI Vision Benchmark

# **Recall:** Sliding Window

Train an image classifier as described previously. 'Slide' a fixed-sized window.



looking for.

# detection window across the image and evaluate the classifier on each

Image credit: KITTI Vision Benchmark

### This is a lot of possible windows! And most will not contain the object we are

- **Object proposal** algorithms generat object-like properties
- These regions are likely to contain background texture
- The object detector then considers the exhaustive sliding window search

Object proposal algorithms generate a short list of regions that have generic

- These regions are likely to contain some kind of foreground object instead of

The object detector then considers these candidate regions only, instead of

.

### First introduced by Alexe et al., who asked 'what is an object?' and defined an 'objectness' score based on several visual cues





### First introduced by Alexe et al., who asked 'what is an object?' and defined an 'objectness' score based on several visual cues



This work argued that objects typically - are unique within the image and stand out as salient have a contrasting appearance from surroundings and/or - have a well-defined closed boundary in space



### Multiscale Saliency

### - Favors regions with a unique appearance within the image





### High scale

### Low scale

### Successful Case

### Failure Case



### **Colour Contrast**

# Favors regions with a contrasting colour appearance from immediate surroundings



### Successful Cases

Failure Case

Figure credit: Alexe et al., 2012



2

### Superpixels Straddling

- Favors regions with a well-defined closed boundary
- contain pixels both inside and outside of the window



(b)



# — Measures the extent to which superpixels (obtained by image segmentation)

(c)



### Superpixels Straddling

- Favors regions with a well-defined closed boundary
- contain pixels both inside and outside of the window







Successful Cases Failure Case

# — Measures the extent to which superpixels (obtained by image segmentation)

(b)





TABLE 2: For each detector [11, 18, 33] we report its performance (left column) and that of our algorithm 1 using the same window scoring function (right column). We show the average number of windows evaluated per image #win and the detection performance as the mean average precision (mAP) over all 20 classes.

	[11] O	BJ- [11]	[18] C	BJ- [18]	ESS-BOW[33]	OBJ-BOW
mAP	0.186	0.162	0.268	0.225	0.127	0.125
#win	79945	1349	18562 -	1358	183501	

Speeding up [11] HOG pedestrian detector [18] Deformable part model detector [33] Bag of words detector

Table credit: Alexe et al., 2012

# Summary

Detection scores in the deformable part model are based on both appearance and location

The deformable part model is trained iteratively by alternating the steps 1. Assume components and part locations given; compute appearance and

- 1. Assume components and part lo offset models
- 2. Assume appearance and offset part locations

An object **proposal** algorithm generates a short list of regions with generic object-like properties that can be evaluated by an object detector in place of an exhaustive sliding window search

2. Assume appearance and offset models given; compute components and