

THE UNIVERSITY OF BRITISH COLUMBIA

CPSC 425: Computer Vision



Lecture 20: Classification (part 3)

Menu for Today

Topics:

- Boosting
- Introduction to neural networks

Redings:

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

Reminders:

- Assignment 5 is due on Wednsday
- Assignment 3 grades are posted

Forsyth & Ponce (2nd ed.) 17.1–17.2





Lecture 19: Re-cap (Decision Tree)

A decision tree is a simple non-linear parametric classifier

A data point starts at the root and recursively proceeds to the child node determined by the feature test, until it reaches a leaf node

- Consists of a tree in which each internal node is associated with a feature test
- The leaf node stores a class label or a probability distribution over class labels

Lecture 19: Re-cap (Decision Tree)

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



Kinect allows users of Microsoft's Xbox 360 console to interact with games using natural body motions instead of a traditional handheld controller. The pose (joint positions) of the user is predicted using a random forest trained on depth features.



Figure credit: J. Shotton et al., 2011

1



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

Simple test: threshold on the difference of two depth values at an offset from a target pixel ...



Figure credit: J. Shotton et al., 2011



.

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



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





....

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



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





....

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



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





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













Figure credit: J. Shotton et al., 2011

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



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



















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

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

Object Classifiers



















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

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

Object Classifiers

(for convenience we will normalize patches to 64x64 ... or 128x128)









Example: Face Detection

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

Example: Face Detection

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

Training Dataset:

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



Faces



Not-faces

Evaluate a Harr Wavelet filter on each training example



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



Faces

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



Not-faces

Evaluate a Harr Wavelet filter on each training example

	$(x_1, 1)$	$(x_2, 1)$
α_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



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



 $f_j(x) > \theta_j$ 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 ε

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

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

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

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

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 ε

2. Re-weight examples

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

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

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

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

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 ε

2. Re-weight examples

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

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

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

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

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 **potentially positive** sub-windows?

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

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



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 Summary



Figure credit: K. Grauman



Hard Negative Mining





Image From: Jamie Kang



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



looking for.

Image credit: KITTI Vision Benchmark

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

Object Proposals

- object-like properties
- background texture
- 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

Object Proposals

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



Figure credit: Alexe et al., 2012



 \mathbf{O}

-
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



Superpixels Straddling

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



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





Superpixels Straddling

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



(a)



Successful Cases

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

(b)



Failure Case



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

.