

RWTH AACHEN  
UNIVERSITY

# Computer Vision – Lecture 7

## Sliding-Window based Object Detection

13.05.2019

Bastian Leibe  
Visual Computing Institute  
RWTH Aachen University  
<http://www.vision.rwth-aachen.de/>  
leibe@vision.rwth-aachen.de

Computer Vision Summer'19

RWTH AACHEN  
UNIVERSITY

## Course Outline

- Image Processing Basics
- Segmentation
  - Segmentation and Grouping
  - Segmentation as Energy Minimization
- Recognition & Categorization
  - Sliding-Window Object Detection
- Local Features & Matching
- Deep Learning
- 3D Reconstruction

2

Computer Vision Summer'19

RWTH AACHEN  
UNIVERSITY

## Recap: MRFs/CRFs for Image Segmentation

- MRF/CRF formulation

Unary potentials  
 $\phi(x_i, y_i)$

Pairwise potentials  
 $\psi(x_i, x_j)$

⇒ Minimize the energy

$$E(\mathbf{x}, \mathbf{y}) = \sum_i \phi(x_i, y_i) + \sum_{i,j} \psi(x_i, x_j)$$

Data (D)

Unary likelihood

Pair-wise Terms

MAP Solution

3

Computer Vision Summer'19

RWTH AACHEN  
UNIVERSITY

## Recap: Energy Formulation

- Energy function

$$E(\mathbf{x}, \mathbf{y}) = \underbrace{\sum_i \phi(x_i, y_i)}_{\text{Unary potentials}} + \underbrace{\sum_{i,j} \psi(x_i, x_j)}_{\text{Pairwise potentials}}$$

- Unary potentials  $\phi$ 
  - Encode local information about the given pixel/patch
  - How likely is a pixel/patch to belong to a certain class (e.g. foreground/background)?
- Pairwise potentials  $\psi$ 
  - Encode neighborhood information
  - How different is a pixel/patch's label from that of its neighbor? (e.g. based on intensity/color/texture difference, edges)

4

Computer Vision Summer'19

RWTH AACHEN  
UNIVERSITY

## Recap: How to Set the Potentials?

- Unary potentials
  - E.g. color model, modeled with a Mixture of Gaussians
$$\phi(x_i, y_i; \theta_\phi) = \log \sum_k \theta_\phi(x_i, k) p(k|x_i) \mathcal{N}(y_i; \bar{y}_k, \Sigma_k)$$

⇒ Learn color distributions for each label

5

Computer Vision Summer'19

RWTH AACHEN  
UNIVERSITY

## Recap: How to Set the Potentials?

- Pairwise potentials
  - Potts Model
 
$$\psi(x_i, x_j; \theta_\psi) = \theta_\psi \delta(x_i \neq x_j)$$
    - Simplest discontinuity preserving model.
    - Discontinuities between any pair of labels are penalized equally.
    - Useful when labels are unordered or number of labels is small.
  - Extension: "Contrast sensitive Potts model"
 
$$\psi(x_i, x_j, g_{ij}(\mathbf{y}); \theta_\psi) = -\theta_\psi g_{ij}(\mathbf{y}) \delta(x_i \neq x_j)$$

where

$$g_{ij}(\mathbf{y}) = e^{-\beta \|y_i - y_j\|^2} \quad \beta = \frac{1}{2} (\text{avg}(\|y_i - y_j\|^2))^{-1}$$
    - ⇒ Discourages label changes except in places where there is also a large change in the observations.

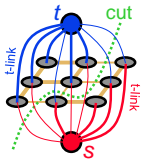
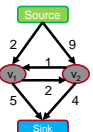
6

Computer Vision Summer'19

Computer Vision Summer'19

## Recap: Graph-Cuts Energy Minimization

- Solve an equivalent graph cut problem
  - Introduce extra nodes: source and sink
  - Weight connections to source/sink (t-links) by  $\phi(x_i = s)$  and  $\phi(x_i = t)$ , respectively.
  - Weight connections between nodes (n-links) by  $\psi(x_i, x_j)$ .
  - Find the minimum cost cut that separates source from sink.
    - ⇒ Solution is equivalent to minimum of the energy.
- s-t Mincut can be solved efficiently
  - Dual to the well-known max flow problem
  - Very efficient algorithms available for regular grid graphs (1-2 MPixels/s)
  - Globally optimal result for 2-class problems

RWTH AACHEN UNIVERSITY

Computer Vision Summer'19

## Recap: When Can s-t Graph Cuts Be Applied?

$$E(L) = \sum_p E_p(L_p) + \sum_{p,q \in N} E(L_p, L_q)$$

Unary potentials (t-links)      Pairwise potentials (n-links)       $L_p \in \{s, t\}$

- s-t graph cuts can only globally minimize **binary energies** that are **submodular**. [Boros & Hummer, 2002, Kolmogorov & Zabih, 2004]

$E(L) \text{ can be minimized by s-t graph cuts} \iff E(s,s) + E(t,t) \leq E(s,t) + E(t,s)$ 

Submodularity ("convexity")

- Submodularity is the discrete equivalent to convexity.
  - Implies that every local energy minimum is a global minimum.
  - ⇒ Solution will be globally optimal.

B. Leibe

Computer Vision Summer'19

## Dealing with Non-Binary Cases

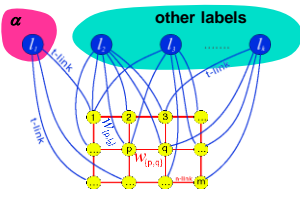
- Limitation to binary energies is often a nuisance.
  - ⇒ E.g. binary segmentation only...
- We would like to solve also multi-label problems.
  - The bad news: Problem is NP-hard with 3 or more labels!
- There exist some approximation algorithms which extend graph cuts to the multi-label case:
  - α-Expansion
  - αβ-Swap
- They are no longer guaranteed to return the globally optimal result.
  - But α-Expansion has a guaranteed approximation quality (2-approx) and converges in a few iterations.

B. Leibe

Computer Vision Summer'19

## α-Expansion Move

- Basic idea:
  - Break multi-way cut computation into a sequence of binary s-t cuts.



Slide credit: Yuri Boykov

B. Leibe

Computer Vision Summer'19

## Topics of This Lecture

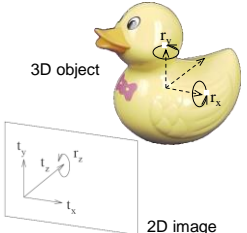
- Object Recognition and Categorization**
  - Problem Definitions
  - Challenges
- Sliding-Window based Object Detection**
  - Detection via Classification
  - Global Representations
  - Classifier Construction
- Classification with SVMs**
  - Support Vector Machines
  - HOG Detector
- Classification with Boosting**
  - AdaBoost
  - Viola-Jones Face Detection

B. Leibe

Computer Vision Summer'19

## Object Recognition: Challenges

- Viewpoint changes
  - Translation
  - Image-plane rotation
  - Scale changes
  - Out-of-plane rotation
- Illumination
- Noise
- Clutter
- Occlusion



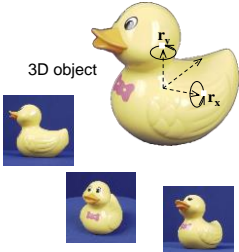
B. Leibe

Computer Vision Summer'19

## Appearance-Based Recognition

RWTH AACHEN UNIVERSITY

- Basic assumption
  - Objects can be represented by a set of images ("appearances").
  - For recognition, it is sufficient to just compare the 2D appearances.
  - No 3D model is needed.



3D object


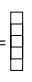

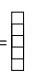

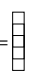
⇒ Fundamental paradigm shift in the 90's

B. Leibe 13

Computer Vision Summer'19

## Global Representation

RWTH AACHEN UNIVERSITY

- Idea
  - Represent each object (view) by a global descriptor.
    -  = 
    -  = 
    -  = 
  - For recognizing objects, just match the descriptors.
  - Some modes of variation are built into the descriptor, the others have to be incorporated in the training data.
    - E.g., a descriptor can be made invariant to image-plane rotations.
    - Other variations:
 

Viewpoint changes	Illumination
– Translation	Noise
– Scale changes	Clutter
– Out-of-plane rotation	Occlusion

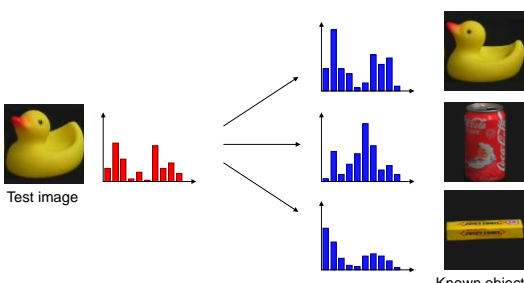
B. Leibe 14

Computer Vision Summer'19

## Appearance based Recognition

RWTH AACHEN UNIVERSITY

- Recognition as feature vector matching



Test image


Known objects

B. Leibe 15

Computer Vision Summer'19

## Appearance based Recognition

RWTH AACHEN UNIVERSITY

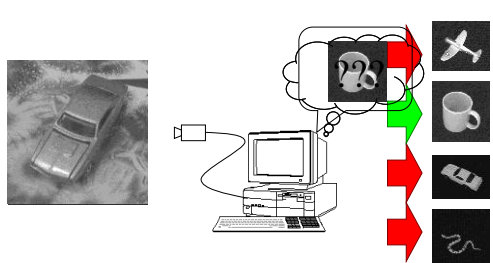
- With multiple training views
  -  with a red feature vector histogram is compared against a vertical stack of blue feature vector histograms for multiple views of the duck and other objects." data-bbox="560 430 880 600"/>

B. Leibe 16

Computer Vision Summer'19

## Identification vs. Categorization

RWTH AACHEN UNIVERSITY



B. Leibe 17

Computer Vision Summer'19

## Identification vs. Categorization

RWTH AACHEN UNIVERSITY


- Find *this particular* object
  - 
- Recognize ANY car
  - 
- Recognize ANY cow
  - 

B. Leibe 18


**RWTH AACHEN UNIVERSITY**

## Object Categorization – Potential Applications


There is a wide range of applications, including...




Autonomous robots




Navigation, driver safety

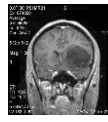


Consumer electronics



Content-based retrieval and analysis for images and videos





Medical image analysis

Slide adapted from Kristen Grauman

**RWTH AACHEN UNIVERSITY**

## Topics of This Lecture

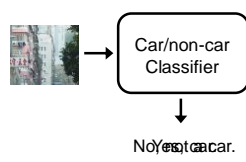
- Object Categorization
  - Problem Definition
  - Challenges
- Sliding-Window based Object Detection
  - Detection via Classification
  - Global Representations
  - Classifier Construction
- Classification with SVMs
  - Support Vector Machines
  - HOG Detector
- Classification with Boosting
  - AdaBoost
  - Viola-Jones Face Detection

B. Leibe

**RWTH AACHEN UNIVERSITY**

## Detection via Classification: Main Idea

- Basic component: a binary classifier




Slide credit: Kristen Grauman

**RWTH AACHEN UNIVERSITY**

## Detection via Classification: Main Idea

- If the object may be in a cluttered scene, slide a window around looking for it.



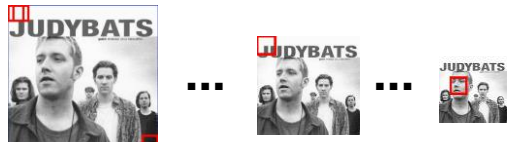
- Essentially, this is a brute-force approach with many local decisions.

Slide credit: Kristen Grauman

**RWTH AACHEN UNIVERSITY**

## What is a Sliding Window Approach?

- Search over space and scale



- Detection as subwindow classification problem
- *"In the absence of a more intelligent strategy, any global image classification approach can be converted into a localization approach by using a sliding-window search."*

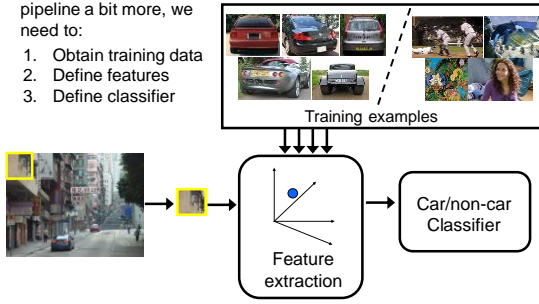
B. Leibe

**RWTH AACHEN UNIVERSITY**

## Detection via Classification: Main Idea

Fleshing out this pipeline a bit more, we need to:

1. Obtain training data
2. Define features
3. Define classifier

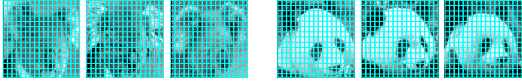


Slide credit: Kristen Grauman


RWTH AACHEN UNIVERSITY

## Feature Extraction: Global Appearance

- Pixel-based representations are sensitive to small shifts



- Color or grayscale-based appearance description can be sensitive to illumination and intra-class appearance variation



Cartoon example: an albino koala

25

Computer Vision Summer'19

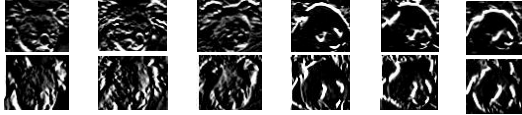
Slide credit: Kristen Grauman

B. Leibe

RWTH AACHEN UNIVERSITY

## Gradient-based Representations

- Idea
  - Consider edges, contours, and (oriented) intensity gradients



26

Computer Vision Summer'19

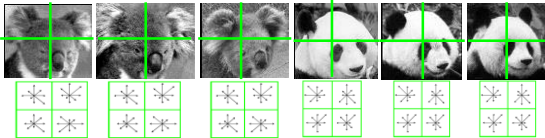
Slide credit: Kristen Grauman

B. Leibe

RWTH AACHEN UNIVERSITY

## Gradient-based Representations

- Idea
  - Consider edges, contours, and (oriented) intensity gradients



- Summarize local distribution of gradients with histograms
  - Locally orderless: offers invariance to small shifts and rotations
  - Localized histograms offer more spatial information than a single global histogram (tradeoff invariant vs. discriminative)
  - Contrast-normalization: try to correct for variable illumination

27

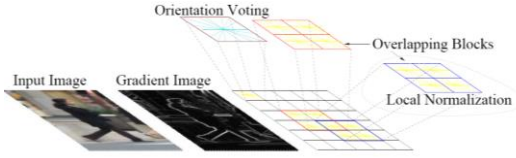
Computer Vision Summer'19

Slide credit: Kristen Grauman

B. Leibe

RWTH AACHEN UNIVERSITY

## Gradient-based Representations: Histograms of Oriented Gradients (HoG)



- Map each grid cell in the input window to a histogram counting the gradients per orientation.
- Code available: <http://pascal.inrialpes.fr/soft/olt/>

[Dalal & Triggs, CVPR 2005]

28

Computer Vision Summer'19

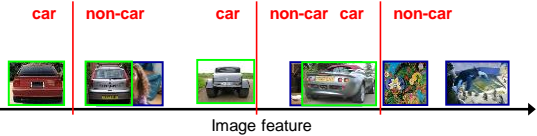
Slide credit: Kristen Grauman

B. Leibe

RWTH AACHEN UNIVERSITY

## Classifier Construction

- How to compute a decision for each subwindow?



29

Computer Vision Summer'19

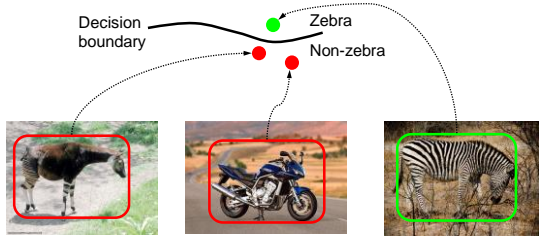
Slide credit: Kristen Grauman

B. Leibe

RWTH AACHEN UNIVERSITY

## Discriminative Methods

- Learn a decision rule (classifier) assigning image features to different classes



30

Computer Vision Summer'19

Slide adapted from Svetlana Lazebnik

B. Leibe

RWTH AACHEN UNIVERSITY

## Classifier Construction: Many Choices...

### Nearest Neighbor

Berg, Berg, Malik 2005, Chum, Zisserman 2007, Boiman, Shechtman, Irani 2008, ...

### Neural networks

LeCun, Bottou, Bengio, Haffner 1998, Rowley, Baluja, Kanade 1998, ...

### Boosting

Viola, Jones 2001, Torralba et al. 2004, Opelt et al. 2006, Benenson 2012, ...

### Support Vector Machines

Vapnik, Schölkopf 1995, Papageorgiou, Poggio '01, Dalal, Triggs 2005, Vedaldi, Zisserman 2012

### Randomized Forests

Amit, Geman 1997, Breiman 2001, Lepetit, Fua 2006, Gall, Lempitsky 2009, ...

Computer Vision Summer'19 | Slide adapted from Kristen Grauman | B. Leibe | 31

RWTH AACHEN UNIVERSITY

## Linear Classifiers

Let  $w = \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$   $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$$w_1 x_1 + w_2 x_2 + b = 0$$

$$\iff w^T x + b = 0$$

Computer Vision Summer'19 | Slide adapted from: Kristen Grauman | B. Leibe | 32

RWTH AACHEN UNIVERSITY

## Linear Classifiers

- Find linear function to separate positive and negative examples

$x_n$  positive:  $w^T x_n + b \geq 0$   
 $x_n$  negative:  $w^T x_n + b < 0$

Which line is best?

Computer Vision Summer'19 | Slide credit: Kristen Grauman | B. Leibe | 33

RWTH AACHEN UNIVERSITY

## Support Vector Machines (SVMs)

- Discriminative classifier based on *optimal separating hyperplane* (i.e. line for 2D case)
- Maximize the *margin* between the positive and negative training examples

Computer Vision Summer'19 | Slide credit: Kristen Grauman | B. Leibe | 34

RWTH AACHEN UNIVERSITY

## Support Vector Machines

see lecture Machine Learning!

- Want line that maximizes the margin.

$x_n$  positive ( $t_n = 1$ ):  $w^T x_n + b \geq 1$   
 $x_n$  negative ( $t_n = -1$ ):  $w^T x_n + b < -1$

For support vectors,  $w^T x_n + b = \pm 1$

Quadratic optimization problem

Minimize  $\frac{1}{2} w^T w$   
 Subject to  $t_n (w^T x_n + b) \geq 1$

Support vectors | Margin | Packages available for that...

Computer Vision Summer'19 | C. Burges, A Tutorial on Support Vector Machines for Pattern Recognition, Data Mining and Knowledge Discovery, 1998 | 35

RWTH AACHEN UNIVERSITY

## Finding the Maximum Margin Line

- Solution:  $w = \sum_{n=1}^N a_n t_n x_n$

Learned weight | Support vector

Computer Vision Summer'19 | C. Burges, A Tutorial on Support Vector Machines for Pattern Recognition, Data Mining and Knowledge Discovery, 1998 | 36

RWTH AACHEN UNIVERSITY

## Finding the Maximum Margin Line

- Solution:  $\mathbf{w} = \sum_{n=1}^N a_n t_n \mathbf{x}_n$
- Classification function:
 
$$f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

$$= \text{sign}\left(\sum_{n=1}^N a_n t_n \mathbf{x}_n^T \mathbf{x} + b\right)$$

If  $f(x) < 0$ , classify as neg.,  
if  $f(x) > 0$ , classify as pos.

  - Notice that this relies on an *inner product* between the test point  $\mathbf{x}$  and the support vectors  $\mathbf{x}_n$
  - (Solving the optimization problem also involves computing the inner products  $\mathbf{x}_n^T \mathbf{x}_m$  between all pairs of training points)

Computer Vision Summer'19

C. Burges, *A Tutorial on Support Vector Machines for Pattern Recognition*, Data Mining and Knowledge Discovery, 1998

37

RWTH AACHEN UNIVERSITY

## Extension: Non-Linear SVMs

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

More on that in the Machine Learning lecture...

Computer Vision Summer'19

Slide from Andrew Moore's tutorial: <http://www.autonlab.org/tutorials/svm.html>

38

RWTH AACHEN UNIVERSITY

## Nonlinear SVMs

- *The kernel trick*: instead of explicitly computing the lifting transformation  $\phi(\mathbf{x})$ , define a kernel function  $K$  such that
 
$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$
- This gives a nonlinear decision boundary in the original feature space:
 
$$\sum_n a_n t_n K(\mathbf{x}_n, \mathbf{x}) + b$$
- Since the optimization formulation uses the data points only in the form of inner products  $\phi(\mathbf{x}_n)^T \phi(\mathbf{x}_m)$ , we never need to actually compute the lifting transformation  $\phi(\mathbf{x})$ .

Computer Vision Summer'19

C. Burges, *A Tutorial on Support Vector Machines for Pattern Recognition*, Data Mining and Knowledge Discovery, 1998

39

RWTH AACHEN UNIVERSITY

## Some Often-Used Kernel Functions

- Linear:
 
$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$$
- Polynomial of power  $p$ :
 
$$K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^p$$
- Gaussian (Radial-Basis Function):
 
$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

Computer Vision Summer'19

Slide from Andrew Moore's tutorial: <http://www.autonlab.org/tutorials/svm.html>

40

RWTH AACHEN UNIVERSITY

## Summary: SVMs for Recognition

1. Define your representation for each example.
2. Select a kernel function.
3. Compute pairwise kernel values between labeled examples
4. Pass this "kernel matrix" to SVM optimization software to identify support vectors & weights.
5. To classify a new example: compute kernel values between new input and support vectors, apply weights, check sign of output.

Computer Vision Summer'19

Slide credit: Kristen Grauman

B. Leibe

41

RWTH AACHEN UNIVERSITY

## Topics of This Lecture

- Object Categorization
  - Problem Definition
  - Challenges
- Sliding-Window based Object Detection
  - Detection via Classification
  - Global Representations
  - Classifier Construction
- Classification with SVMs
  - Support Vector Machines
  - HOG Detector
- Classification with Boosting
  - AdaBoost
  - Viola-Jones Face Detection

Computer Vision Summer'19

B. Leibe

42

RWTH AACHEN UNIVERSITY

## HOG Descriptor Processing Chain

Computer Vision Summer'19



Image Window

43

Slide adapted from Navneet Dalal

RWTH AACHEN UNIVERSITY

## HOG Descriptor Processing Chain

- Optional: Gamma compression
  - Goal: Reduce effect of overly strong gradients
  - Replace each pixel color/intensity by its square-root
 
$$x \mapsto \sqrt{x}$$

⇒ Small performance improvement




Image Window

44

Slide adapted from Navneet Dalal

RWTH AACHEN UNIVERSITY

## HOG Descriptor Processing Chain

- Gradient computation
  - Compute gradients on all color channels and take strongest one
  - Simple finite difference filters work best (no Gaussian smoothing)


$$\begin{bmatrix} -1 & 0 & 1 \\ 0 & 1 & -1 \end{bmatrix}$$


Image Window

45

Slide adapted from Navneet Dalal

RWTH AACHEN UNIVERSITY

## HOG Descriptor Processing Chain

- Spatial/Orientation binning
  - Compute localized histograms of oriented gradients
  - Typical subdivision: 8x8 cells with 8 or 9 orientation bins

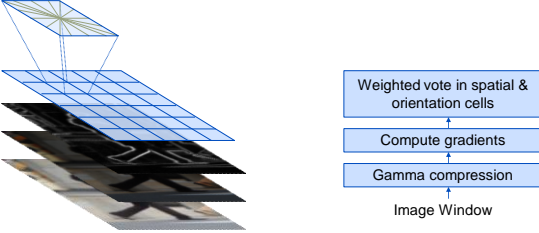


Image Window

46

Slide adapted from Navneet Dalal

RWTH AACHEN UNIVERSITY

## HOG Descriptor Processing Chain

- 2-Stage contrast normalization
  - L2 normalization, clipping, L2 normalization

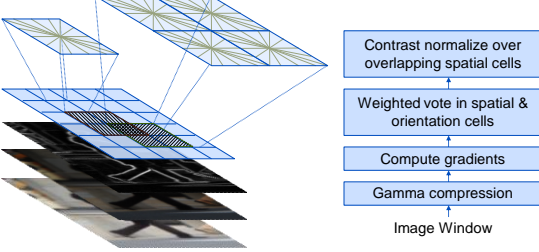


Image Window

49

Slide adapted from Navneet Dalal

RWTH AACHEN UNIVERSITY

## HOG Descriptor Processing Chain

- Feature vector construction
  - Collect HOG blocks into vector

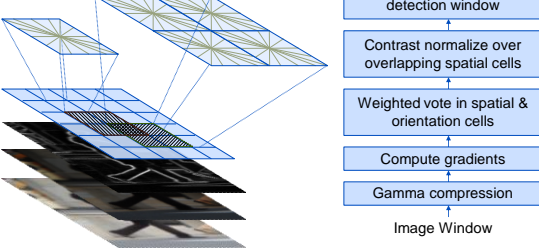
$$[ \dots, \dots, \dots, \dots ]$$


Image Window

50

Slide adapted from Navneet Dalal



RWTH AACHEN UNIVERSITY

## HOG Descriptor Processing Chain

- SVM Classification
  - Typically using a linear SVM
  - [ ..., ..., ..., ... ]

Object/Non-object

Linear SVM

Collect HOGs over detection window

Contrast normalize over overlapping spatial cells

Weighted vote in spatial & orientation cells

Compute gradients

Gamma compression

Image Window

51

Computer Vision Summer'19

Slide adapted from Navneet Dalal

RWTH AACHEN UNIVERSITY

## Pedestrian Detection with HOG

- Intuition
  - Train a pedestrian template using a linear SVM
  - At test time, convolve feature map with learned template  $w$

HOG feature map      Template      Detector response map

N. Dalal and B. Triggs, [Histograms of Oriented Gradients for Human Detection](#), CVPR 2005

Slide credit: Svetlana Lazebnik

Computer Vision Summer'19

RWTH AACHEN UNIVERSITY

## Non-Maximum Suppression

Clip detection score

Map each detection to 3D  $[x, y, scale]$  space

Apply robust mode detection, e.g. mean shift

Non-maximum suppression

Fusion of multiple detections

Goal

53

Computer Vision Summer'19

B. Leibe      Image source: Navneet Dalal, PhD Thesis

RWTH AACHEN UNIVERSITY

## Pedestrian detection with HoGs & SVMs

- N. Dalal, B. Triggs, [Histograms of Oriented Gradients for Human Detection](#), CVPR 2005.

54

Computer Vision Summer'19

Slide credit: Kristen Grauman      B. Leibe

RWTH AACHEN UNIVERSITY

## Classifier Construction: Many Choices...

**Nearest Neighbor**

Shakhnarovich, Viola, Darrell 2003  
Berg, Berg, Malik 2005,  
Boiman, Shechtman, Irani 2008, ...

**Neural networks**

LeCun, Bottou, Bengio, Haffner 1998  
Rowley, Baluja, Kanade 1998  
...

**Boosting**

Viola, Jones 2001,  
Torralba et al. 2004,  
Opelt et al. 2006,  
Benenson 2012, ...

**Support Vector Machines**

Vapnik, Schölkopf 1995,  
Papageorgiou, Poggio '01,  
Dalal, Triggs 2005,  
Vedaldi, Zisserman 2012

**Randomized Forests**

Amit, Geman 1997,  
Breiman 2001,  
Lepetit, Fua 2006,  
Gall, Lempitsky 2009, ...

55

Computer Vision Summer'19

Slide adapted from Kristen Grauman      B. Leibe

RWTH AACHEN UNIVERSITY

## Boosting

- Idea
  - Build a strong classifier  $H$  by combining a number of "weak classifiers"  $h_1, \dots, h_M$ , which need only be better than chance.
  - Sequential learning process: at each iteration, add a weak classifier
- Flexible to choice of weak learner
  - including fast simple classifiers that alone may be inaccurate
- We'll look at Freund & Schapire's AdaBoost algorithm
  - Easy to implement
  - Base learning algorithm for Viola-Jones face detector

Y. Freund and R. Schapire, [A short introduction to boosting](#), *Journal of Japanese Society for Artificial Intelligence*, 14(5):771-780, 1999.

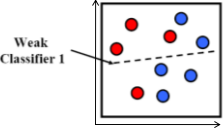
56

Computer Vision Summer'19

Slide credit: Kristen Grauman      B. Leibe

Computer Vision Summer'19

## AdaBoost: Intuition



Weak Classifier 1

Consider a 2D feature space with **positive** and **negative** examples.

Each weak classifier splits the training examples with at least 50% accuracy.

Examples misclassified by a previous weak learner are given more emphasis at future rounds.

Figure adapted from Freund and Schapire

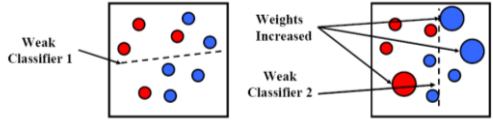
Slide credit: Kristen Grauman

B. Leibe

57

Computer Vision Summer'19

## AdaBoost: Intuition



Weak Classifier 1

Weights Increased

Weak Classifier 2

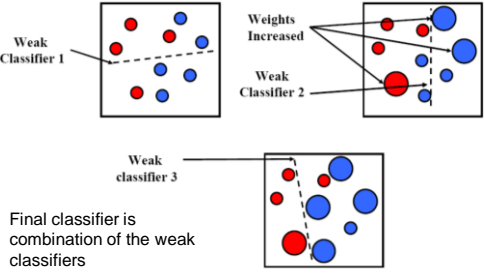
Slide credit: Kristen Grauman

B. Leibe

58

Computer Vision Summer'19

## AdaBoost: Intuition



Weak Classifier 1

Weights Increased

Weak Classifier 2

Weak classifier 3

Final classifier is combination of the weak classifiers

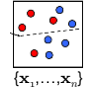
Slide credit: Kristen Grauman

B. Leibe

59

Computer Vision Summer'19

## AdaBoost – Formalization



- 2-class classification problem
  - Given: training set  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  with target values  $\mathbf{T} = \{t_1, \dots, t_N\}$ ,  $t_n \in \{-1, 1\}$ .
  - Associated weights  $\mathbf{W} = \{w_1, \dots, w_N\}$  for each training point.
- Basic steps
  - In each iteration, AdaBoost trains a new weak classifier  $h_m(\mathbf{x})$  based on the current weighting coefficients  $\mathbf{W}^{(m)}$ .
  - We then adapt the weighting coefficients for each point
    - Increase  $w_n$  if  $\mathbf{x}_n$  was misclassified by  $h_m(\mathbf{x})$ .
    - Decrease  $w_n$  if  $\mathbf{x}_n$  was classified correctly by  $h_m(\mathbf{x})$ .
  - Make predictions using the final combined model

$$H(\mathbf{x}) = \text{sign} \left( \sum_{m=1}^M \alpha_m h_m(\mathbf{x}) \right)$$

Slide credit: Kristen Grauman

B. Leibe

60

Computer Vision Summer'19

## AdaBoost: Detailed Training Algorithm

see lecture Machine Learning!

- Initialization: Set  $w_n^{(1)} = \frac{1}{N}$  for  $n = 1, \dots, N$ .
- For  $m = 1, \dots, M$  iterations
  - Train a new weak classifier  $h_m(\mathbf{x})$  using the current weighting coefficients  $\mathbf{W}^{(m)}$  by minimizing the weighted error function
 
$$J_m = \sum_{n=1}^N w_n^{(m)} I(h_m(\mathbf{x}_n) \neq t_n) \quad I(A) = \begin{cases} 1, & \text{if } A \text{ is true} \\ 0, & \text{else} \end{cases}$$
  - Estimate the weighted error of this classifier on  $\mathbf{X}$ :
 
$$\epsilon_m = \frac{\sum_{n=1}^N w_n^{(m)} I(h_m(\mathbf{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}}$$
  - Calculate a weighting coefficient for  $h_m(\mathbf{x})$ :
 
$$\alpha_m = \ln \left\{ \frac{1 - \epsilon_m}{\epsilon_m} \right\}$$
  - Update the weighting coefficients:
 
$$w_n^{(m+1)} = w_n^{(m)} \exp \{ \alpha_m I(h_m(\mathbf{x}_n) \neq t_n) \}$$

Slide credit: Kristen Grauman

B. Leibe

61

Computer Vision Summer'19

## AdaBoost: Recognition

- Evaluate all selected weak classifiers on test data.
 
$$h_1(\mathbf{x}), \dots, h_m(\mathbf{x})$$
- Final classifier is weighted combination of selected weak classifiers:
 
$$H(\mathbf{x}) = \text{sign} \left( \sum_{m=1}^M \alpha_m h_m(\mathbf{x}) \right)$$
- Very simple procedure!
  - Less than 10 lines in Matlab!
  - But works extremely well in practice...

Slide credit: Kristen Grauman


B. Leibe

62

Computer Vision Summer'19

## Example: Face Detection

- Frontal faces are a good example of a class where global appearance models + a sliding window detection approach fit well:
  - Regular 2D structure
  - Center of face almost shaped like a "patch"/window




- Now we'll take AdaBoost and see how the Viola-Jones face detector works

Slide credit: Kristen Grauman B. Leibe 63

Computer Vision Summer'19

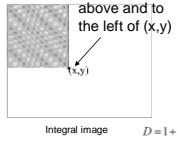
## Feature extraction

"Rectangular" filters



Feature output is difference between adjacent regions

Value at (x,y) is sum of pixels above and to the left of (x,y)



Efficiently computable with integral image: any sum can be computed in constant time

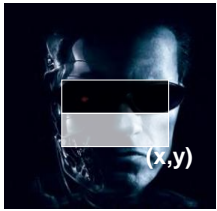
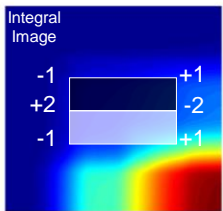
Avoid scaling images → scale features directly for same cost

$$D = 1 + 4 - (2 + 3) = A + (A + B + C + D) - (A + C + A + B) = D$$

Slide credit: Kristen Grauman B. Leibe 64 [Viola & Jones, CVPR 2001]

Computer Vision Summer'19

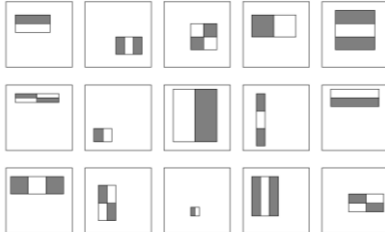
## Example

Slide credit: Svetlana Lazebnik B. Leibe 65

Computer Vision Summer'19

## Large Library of Features



Considering all possible filter parameters: position, scale, and type: 180,000+ possible features associated with each 24 x 24 window

Use AdaBoost both to select the informative features and to form the classifier

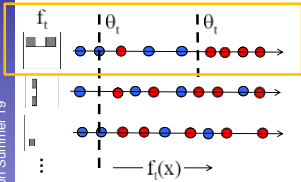
Weak classifier: feature output >  $\theta$ ?

Slide credit: Kristen Grauman B. Leibe 66 [Viola & Jones, CVPR 2001]

Computer Vision Summer'19

## AdaBoost for Feature+Classifier Selection

Want to select the single rectangle feature and threshold that best separates positive (faces) and negative (non-faces) training examples, in terms of weighted error.



Resulting weak classifier:

$$h_i(x) = \begin{cases} +1 & \text{if } f_i(x) > \theta_i \\ -1 & \text{otherwise} \end{cases}$$

For next round, reweight the examples according to errors, choose another filter/threshold combo.

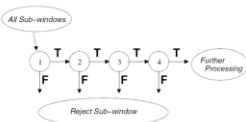
Outputs of a possible rectangle feature on faces and non-faces.

Slide credit: Kristen Grauman B. Leibe 67 [Viola & Jones, CVPR 2001]

Computer Vision Summer'19

## Cascading Classifiers for Detection

- Even if the filters are fast to compute, each new image has a lot of possible windows to search.
- For efficiency, apply less accurate but faster classifiers first to immediately discard windows that clearly appear to be negative; e.g.,
  - Filter for promising regions with an initial inexpensive classifier
  - Build a chain of classifiers, choosing cheap ones with low false negative rates early in the chain



[Fleuret & Geman, IJCV 2001]  
[Rowley et al., PAMI 1998]  
[Viola & Jones, CVPR 2001]

Slide credit: Kristen Grauman B. Leibe 69 Figure from Viola & Jones CVPR

RWTH AACHEN UNIVERSITY

## Viola-Jones Face Detector: Summary

- Train with 5K positives, 350M negatives
- Real-time detector using 38 layer cascade
- 6061 features in final layer
- [Implementation available in OpenCV: <http://sourceforge.net/projects/opencvlibrary/>]

71

Slide credit: Kristen Grauman B. Leibe

RWTH AACHEN UNIVERSITY

## Viola-Jones Face Detector: Results

73

Slide credit: Kristen Grauman B. Leibe

RWTH AACHEN UNIVERSITY

## Viola-Jones Face Detector: Results

74

Slide credit: Kristen Grauman B. Leibe

RWTH AACHEN UNIVERSITY

## You Can Try It At Home...

- The Viola & Jones detector was a huge success
  - First real-time face detector available
  - Many derivative works and improvements
- C++ implementation available in OpenCV [Lienhart, 2002]
  - <http://sourceforge.net/projects/opencvlibrary/>
- Matlab wrappers for OpenCV code available, e.g. here
  - <http://www.mathworks.com/matlabcentral/fileexchange/19912>

P. Viola, M. Jones, [Robust Real-Time Face Detection](#), IJCV, Vol. 57(2), 2004

76

Slide credit: Kristen Grauman B. Leibe

RWTH AACHEN UNIVERSITY

## Example Application

Frontal faces detected and then tracked, character names inferred with alignment of script and subtitles.

Everingham, M., Sivic, J. and Zisserman, A. "Hello! My name is... Buffy" - Automatic naming of characters in TV video, BMVC 2006. <http://www.robots.ox.ac.uk/~vgg/research/nface/index.html>

77

Slide credit: Kristen Grauman B. Leibe

RWTH AACHEN UNIVERSITY

## Summary: Sliding-Windows

- Pros
  - Simple detection protocol to implement
  - Good feature choices critical
  - Past successes for certain classes
  - Good detectors available (Viola & Jones, HOG, etc.)
- Cons/Limitations
  - High computational complexity
    - For example: 250,000 locations x 30 orientations x 4 scales = 30,000,000 evaluations!
    - This puts tight constraints on the classifiers we can use.
    - If training binary detectors independently, this means cost increases linearly with number of classes.
  - With so many windows, false positive rate better be low


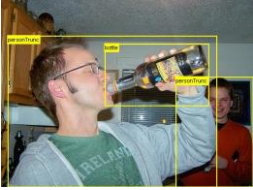
78

Slide credit: Kristen Grauman B. Leibe

RWTH AACHEN  
UNIVERSITY

## Limitations (continued)

- Not all objects are "box" shaped


Computer Vision Summer'19 79

Slide credit: Kristen Grauman B. Leibe

RWTH AACHEN  
UNIVERSITY

## Limitations (continued)

- Non-rigid, deformable objects not captured well with representations assuming a fixed 2D structure; or must assume fixed viewpoint
- Objects with less-regular textures not captured well with holistic appearance-based descriptions



Computer Vision Summer'19 80

Slide credit: Kristen Grauman B. Leibe

RWTH AACHEN  
UNIVERSITY

## Limitations (continued)

- If considering windows in isolation, context is lost




Computer Vision Summer'19 81

Figure credit: Derek Hoiem B. Leibe

RWTH AACHEN  
UNIVERSITY

## Limitations (continued)

- In practice, often entails large, cropped training set (expensive)
- Requiring good match to a global appearance description can lead to sensitivity to partial occlusions




Computer Vision Summer'19 82

Image credit: Adam, Rivin, & Shimshoni K. Grauman, B. Leibe

RWTH AACHEN  
UNIVERSITY

## References and Further Reading

- Read the HOG paper
  - N. Dalal, B. Triggs, [Histograms of Oriented Gradients for Human Detection](#), CVPR, 2005.
- HOG Detector
  - Code available: <http://pascal.inrialpes.fr/soft/olt/>

Computer Vision Summer'19