



Discriminative classifiers  
for image recognition  
May 23<sup>rd</sup>, 2019

---

---

---

---

---

---

---

---

### Announcements

- PS3 out tonight; due 6/4, 11:59 pm

2

---

---

---

---

---

---

---

---

### Outline

- **Last time:** window-based generic object detection
  - basic pipeline
  - face detection with boosting as case study

3

---

---

---

---

---


---

---

---

**Window-based models**  
**Building an object model**

Given the representation, train a binary classifier

 → Car/non-car Classifier  
 ↓  
 Yes it is a car.

4  
Kristen Grauman

---

---

---

---



---

---

---

---

**Window-based models**  
**Generating and scoring candidates**

 →  → Car/non-car Classifier

5  
Kristen Grauman

---

---

---

---

---

---

---

---


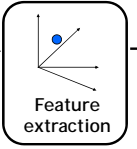
**Window-based object detection: recap**



**Training:**

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

**Given new image:**

1. Slide window
2. Score by classifier

 Training examples  
 ↓ ↓ ↓  
 Feature extraction → Car/non-car Classifier

 →  → Feature extraction → Car/non-car Classifier

6  
Kristen Grauman

---

---

---

---

---

---

---

---

### Viola-Jones detector: summary

- A seminal approach to real-time object detection
- Training is slow, but detection is very fast
- Key ideas
  - > *Integral images* for fast feature evaluation



P. Viola and M. Jones. [Rapid object detection using a boosted cascade of simple features](#). CVPR 2001.

P. Viola and M. Jones. [Robust real-time face detection](#). IJCV 57(2), 2004.

7

---

---

---

---

---

---

---

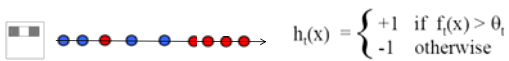
---

---

---

### Viola-Jones detector: summary

- A seminal approach to real-time object detection
- Training is slow, but detection is very fast
- Key ideas
  - > *Integral images* for fast feature evaluation
  - > *Boosting* for feature selection



P. Viola and M. Jones. [Rapid object detection using a boosted cascade of simple features](#). CVPR 2001.

P. Viola and M. Jones. [Robust real-time face detection](#). IJCV 57(2), 2004.

8

---

---

---

---

---

---

---

---

---

---

### Viola-Jones detector: summary

- A seminal approach to real-time object detection
- Training is slow, but detection is very fast
- Key ideas
  - > *Integral images* for fast feature evaluation
  - > *Boosting* for feature selection
  - > *Attentional cascade* of classifiers for fast rejection of non-face windows

P. Viola and M. Jones. [Rapid object detection using a boosted cascade of simple features](#). CVPR 2001.

P. Viola and M. Jones. [Robust real-time face detection](#). IJCV 57(2), 2004.

9

---

---

---

---

---

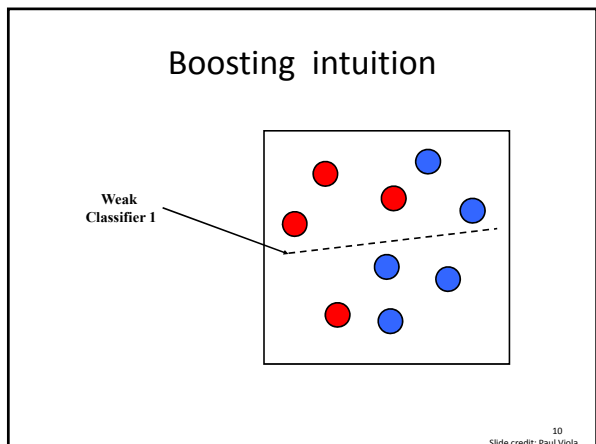
---

---

---

---

---



---

---

---

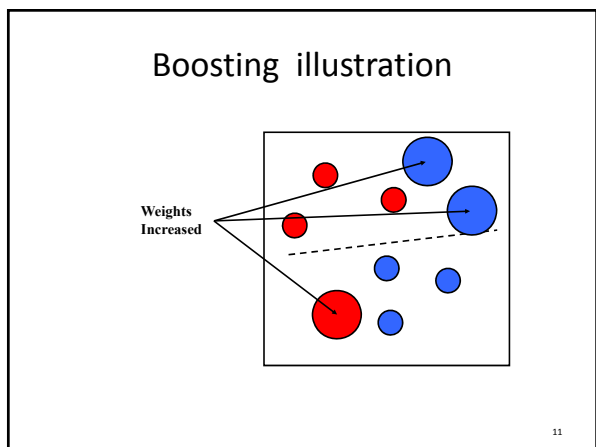
---

---

---

---

---



---

---

---

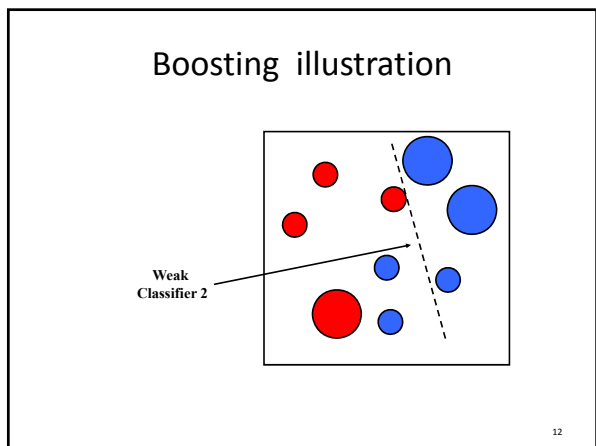
---

---

---

---

---



---

---

---

---

---

---

---

---

### Boosting illustration

Weights Increased

13

This diagram illustrates a weak classifier's decision boundary (dashed line) separating red and blue circles. The circles vary in size, representing their weights. The text 'Weights Increased' has arrows pointing to the larger circles that have been misclassified by the current decision boundary.

---

---

---

---

---

---

---

---

### Boosting illustration

Weak Classifier 3

14

This diagram shows the same set of red and blue circles as the previous slide. A vertical dashed line represents the decision boundary of 'Weak Classifier 3'. An arrow points from the text 'Weak Classifier 3' to this dashed line.

---

---

---

---

---

---

---

---

### Boosting illustration

Final classifier is a combination of weak classifiers

15

This diagram shows the same set of red and blue circles. Multiple dashed lines represent the decision boundaries of several weak classifiers. The text 'Final classifier is a combination of weak classifiers' indicates that the final decision boundary is a complex combination of these individual weak classifiers.

---

---

---

---


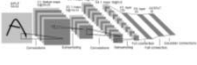
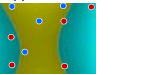


---

---

---

---

### Discriminative classifier construction

<p><b>Nearest neighbor</b></p>  <p>10<sup>6</sup> examples</p> <p>Shakhnarovich, Viola, Darrell 2003 Berg, Berg, Malik 2005...</p>	<p><b>Neural networks</b></p>  <p>LeCun, Bottou, Bengio, Haffner 1998 Rowley, Baluja, Kanade 1998 ...</p>	
<p><b>Support Vector Machines</b></p>  <p>Guyon, Vapnik Heisele, Serre, Poggio, 2001,...</p>	<p><b>Boosting</b></p>  <p>Viola, Jones 2001, Torralba et al. 2004, Opelt et al. 2006,...</p>	<p><b>Conditional Random Fields</b></p>  <p>McCallum, Freitag, Pereira 2000; Kumar, Hebert 2003 ...</p>

16  
Slide adapted from Antonio Torralba

---

---

---

---

---

---


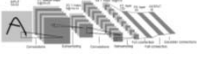
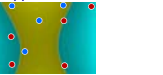


---

---

---

---

### Discriminative classifier construction

<p><b>Nearest neighbor</b></p>  <p>10<sup>6</sup> examples</p> <p>Shakhnarovich, Viola, Darrell 2003 Berg, Berg, Malik 2005...</p>	<p><b>Neural networks</b></p>  <p>LeCun, Bottou, Bengio, Haffner 1998 Rowley, Baluja, Kanade 1998 ...</p>	
<p><b>Support Vector Machines</b></p>  <p>Guyon, Vapnik Heisele, Serre, Poggio, 2001,...</p>	<p><b>Boosting</b></p>  <p>Viola, Jones 2001, Torralba et al. 2004, Opelt et al. 2006,...</p>	<p><b>Conditional Random Fields</b></p>  <p>McCallum, Freitag, Pereira 2000; Kumar, Hebert 2003 ...</p>

17  
Slide adapted from Antonio Torralba

---

---

---

---

---

---

---

---

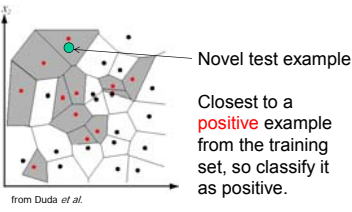
---

---

### Nearest Neighbor classification

- Assign label of nearest training data point to each test data point

Black = negative  
Red = positive



Novel test example

Closest to a positive example from the training set, so classify it as positive.

from Duda et al.

Voronoi partitioning of feature space for 2-category 2D data

18

---

---

---

---

---

---

---

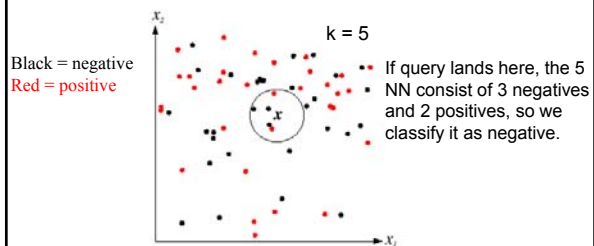
---

---

---

### K-Nearest Neighbors classification

- For a new point, find the k closest points from training data
- Labels of the k points "vote" to classify



19  
Source: D. Lowe

---

---

---

---

---

---

---

---

### A nearest neighbor recognition example

20

---

---

---

---

---

---

---

---

### Where in the World?



[Hays and Efros. **im2gps**: Estimating Geographic Information from a Single Image. CVPR 2008.]  
Slides: James Hays

---

---

---

---

---

---

---

---

Where in the World?



22  
Slides: James Hays

---

---

---

---

---

---

---

---

Where in the World?



23  
Slides: James Hays

---

---

---

---

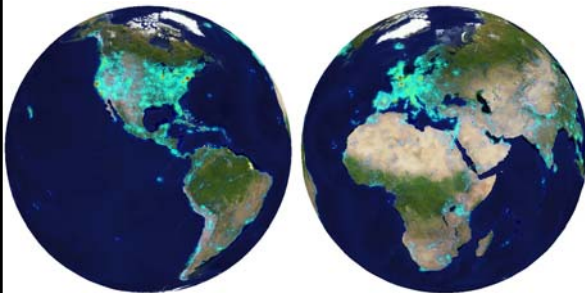
---

---

---

---

6+ million geotagged photos  
by 109,788 photographers



Annotated by Flickr users

24  
Slides: James Hays

---

---

---

---

---

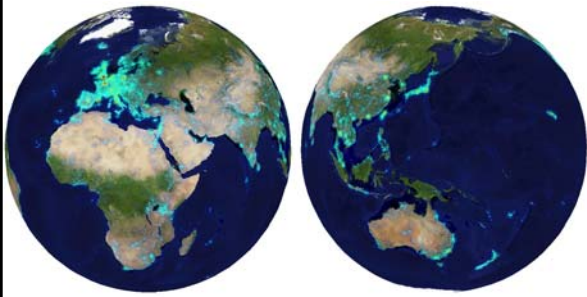
---

---

---



6+ million geotagged photos  
by 109,788 photographers



Annotated by Flickr users

Slides: James Hays

The slide features two globes side-by-side. The left globe shows the Western Hemisphere (Africa, Europe, Americas) with numerous small red dots indicating geotagged photo locations. The right globe shows the Eastern Hemisphere (Asia, Australia, Oceania) with similar red dots. The text above the globes states '6+ million geotagged photos by 109,788 photographers'. Below the globes, it says 'Annotated by Flickr users' and 'Slides: James Hays'.

---

---

---

---

---

---

---

---

Quantitative Evaluation Test Set



Quantitative Evaluation Test Set

The slide displays a grid of 48 small, diverse scene images arranged in 6 rows and 8 columns. The images include various landscapes, buildings, and objects. At the bottom center of the grid, there are three small black dots.

---

---

---

---

---

---

---

---

Which scene properties are relevant?

Which scene properties are relevant?

27

The slide contains the text 'Which scene properties are relevant?' at the top. At the bottom right corner, there is a small number '27'.

---

---

---

---

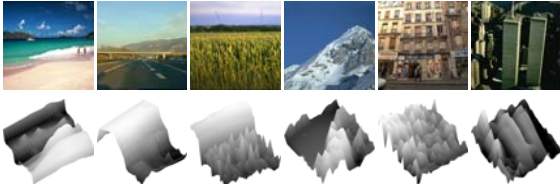
---

---

---

---

**Spatial Envelope Theory of Scene Representation**  
Oliva & Torralba (2001)



↓

A scene is a single surface that can be represented by global (statistical) descriptors

28  
Slide Credit: Aude Oliva

---

---

---

---

---

---

---

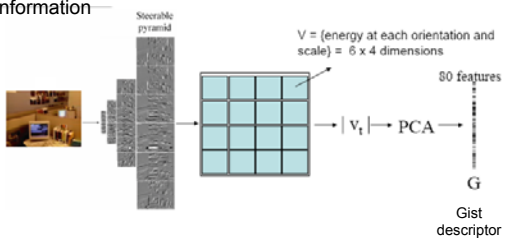
---

---

---

**Global texture:  
capturing the “Gist” of the scene**

Capture global image properties while keeping some spatial information —



29  
Oliva & Torralba IJCV 2001, Torralba et al. CVPR 2003

---

---

---

---

---

---

---

---

---

---

Which scene properties are relevant?

- **Gist scene descriptor**
- **Color Histograms** - L\*A\*B\* 4x14x14 histograms
- **Texton Histograms** – 512 entry, filter bank based
- **Line Features** – Histograms of straight line stats

30

---

---

---

---

---

---

---

---

---

---

### Im2GPS: Scene Matches

[Hays and Efros, *im2gps*: Estimating Geographic Information from a Single Image, CVPR 2008.] Slides: James Hays

---

---

---

---

---

---

---

---

32  
Slides: James Hays

---

---

---

---

---

---

---

---

### Im2GPS: Scene Matches

[Hays and Efros, *im2gps*: Estimating Geographic Information from a Single Image, CVPR 2008.] Slides: James Hays

---

---

---

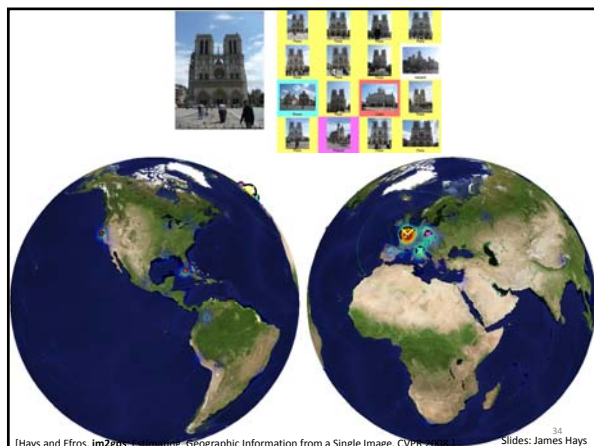
---

---

---

---

---



---

---

---

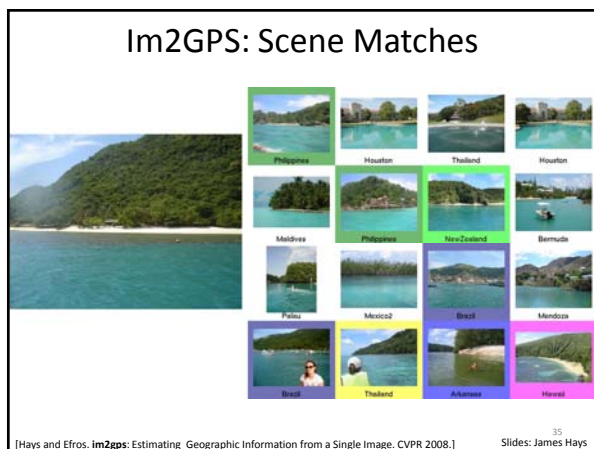
---

---

---

---

---



---

---

---

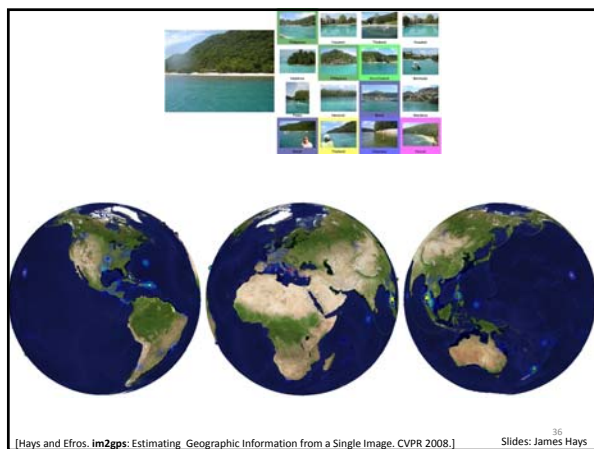
---

---

---

---

---



---

---

---

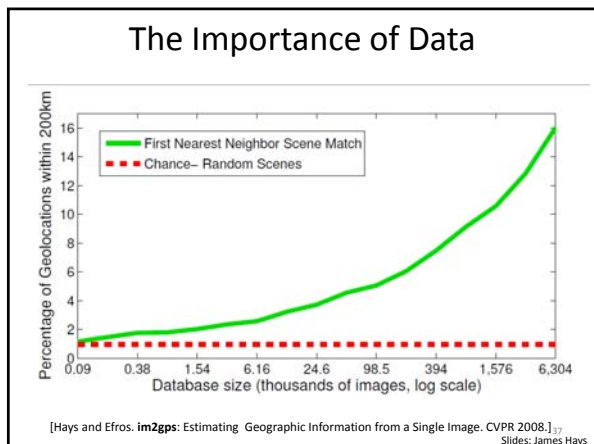
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

- ### Nearest neighbors: pros and cons
- **Pros:**
    - Simple to implement
    - Flexible to feature / distance choices
    - Naturally handles multi-class cases
    - Can do well in practice with enough representative data
  - **Cons:**
    - Large search problem to find nearest neighbors (slow during testing)
    - Storage of data
    - Must have a meaningful distance function
- 38

---

---

---

---

---

---

---

---

---

---

- ### Outline
- Discriminative classifiers
    - Boosting (last time)
    - Nearest neighbors
    - Support vector machines
- 39

---

---

---

---

---

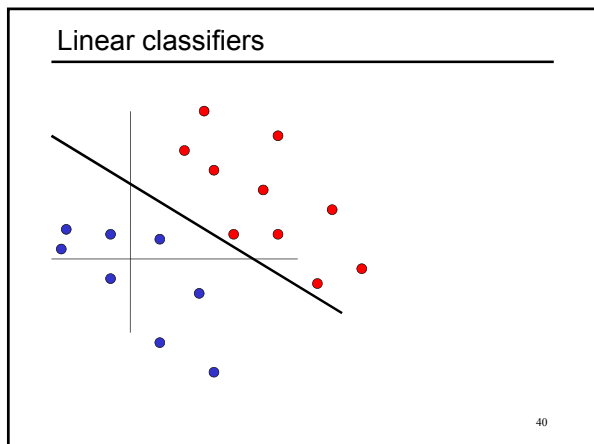
---

---

---

---

---




---

---

---

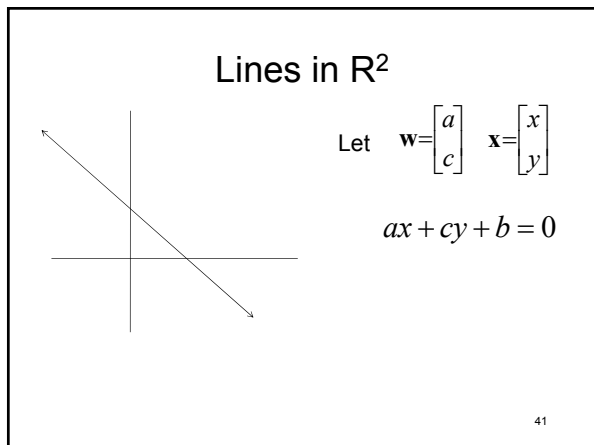
---

---

---

---

---




---

---

---

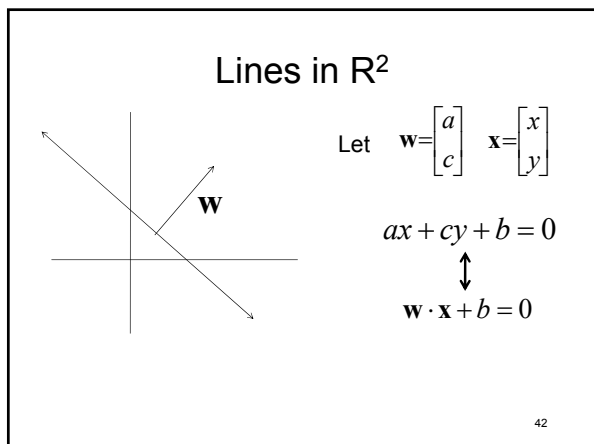
---

---

---

---

---




---

---

---

---

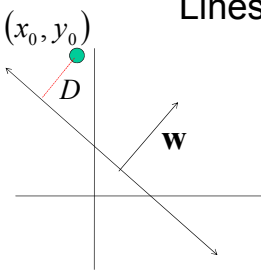
---

---

---

---

**Lines in  $\mathbb{R}^2$**



Let  $\mathbf{w} = \begin{bmatrix} a \\ c \end{bmatrix}$   $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$

$$ax + cy + b = 0$$

$$\updownarrow$$

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

43

---

---

---

---

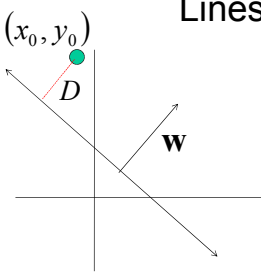
---

---

---

---

**Lines in  $\mathbb{R}^2$**



Let  $\mathbf{w} = \begin{bmatrix} a \\ c \end{bmatrix}$   $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$

$$ax + cy + b = 0$$

$$\updownarrow$$

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

$D = \frac{|ax_0 + cy_0 + b|}{\sqrt{a^2 + c^2}}$

distance from point to line 44

---

---

---

---

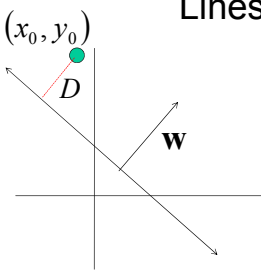
---

---

---

---

**Lines in  $\mathbb{R}^2$**



Let  $\mathbf{w} = \begin{bmatrix} a \\ c \end{bmatrix}$   $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$

$$ax + cy + b = 0$$

$$\updownarrow$$

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

$D = \frac{|ax_0 + cy_0 + b|}{\sqrt{a^2 + c^2}} = \frac{\mathbf{w}^T \mathbf{x}_0 + b}{\|\mathbf{w}\|}$  distance from point to line 45

---

---

---

---

---

---

---

---

### Linear classifiers

- Find linear function to separate positive and negative examples

$x_i$  positive:  $x_i \cdot w + b \geq 0$   
 $x_i$  negative:  $x_i \cdot w + b < 0$

Which line is best?

46

---

---

---

---

---

---

---

---

### Support Vector Machines (SVMs)

- Discriminative classifier based on *optimal separating line* (for 2d case)
- Maximize the *margin* between the positive and negative training examples

47

---

---

---

---

---

---

---

---

### Support vector machines

- Want line that maximizes the margin.

$x_i$  positive ( $y_i = 1$ ):  $x_i \cdot w + b \geq 1$   
 $x_i$  negative ( $y_i = -1$ ):  $x_i \cdot w + b \leq -1$

For support vectors,  $x_i \cdot w + b = \pm 1$

Support vectors Margin

C. Burges, [A Tutorial on Support Vector Machines for Pattern Recognition](#), Data Mining and Knowledge Discovery, 1998

---

---

---

---

---

---

---

---



### Support vector machines

- Want line that maximizes the margin.

The diagram shows a set of data points (red and blue) separated by a solid line. Two dashed lines parallel to the solid line represent the decision boundaries. The points on these dashed lines are labeled as support vectors. The distance between the two dashed lines is labeled as the margin M. The decision boundary is labeled  $w \cdot x + b = 0$ , the upper decision boundary is  $w \cdot x + b = 1$ , and the lower decision boundary is  $w \cdot x + b = -1$ .

$x_i$  positive ( $y_i = 1$ ):  $x_i \cdot w + b \geq 1$   
 $x_i$  negative ( $y_i = -1$ ):  $x_i \cdot w + b \leq -1$   
 For support vectors,  $x_i \cdot w + b = \pm 1$

Distance between point and line:  $\frac{|x_i \cdot w + b|}{\|w\|}$

For support vectors:  $\frac{w^T x + b}{\|w\|} = \frac{\pm 1}{\|w\|}$   $M = \left| \frac{1}{\|w\|} - \frac{-1}{\|w\|} \right| = \frac{2}{\|w\|}$

Support vectors Margin M

49

---

---

---

---

---

---

---

---

---

---

### Support vector machines

- Want line that maximizes the margin.

The diagram shows a set of data points (red and blue) separated by a solid line. Two dashed lines parallel to the solid line represent the decision boundaries. The points on these dashed lines are labeled as support vectors. The distance between the two dashed lines is labeled as the margin M. The decision boundary is labeled  $w \cdot x + b = 0$ , the upper decision boundary is  $w \cdot x + b = 1$ , and the lower decision boundary is  $w \cdot x + b = -1$ .

$x_i$  positive ( $y_i = 1$ ):  $x_i \cdot w + b \geq 1$   
 $x_i$  negative ( $y_i = -1$ ):  $x_i \cdot w + b \leq -1$   
 For support vectors,  $x_i \cdot w + b = \pm 1$

Distance between point and line:  $\frac{|x_i \cdot w + b|}{\|w\|}$

Therefore, the margin is  $2 / \|w\|$

Support vectors Margin M

50

---

---

---

---

---

---

---

---

---

---

### Finding the maximum margin line

- Maximize margin  $2/\|w\|$
- Correctly classify all training data points:
  - $x_i$  positive ( $y_i = 1$ ):  $x_i \cdot w + b \geq 1$
  - $x_i$  negative ( $y_i = -1$ ):  $x_i \cdot w + b \leq -1$

Quadratic optimization problem:

Minimize  $\frac{1}{2} w^T w$

Subject to  $y_i(w \cdot x_i + b) \geq 1$

51

---

---

---

---

---

---

---

---

---

---

### Finding the maximum margin line

---

- Solution:  $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$

learned weight

Support vector

52

---

---

---

---

---

---

---

---

### Finding the maximum margin line

---

- Solution:  $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$   
 $b = y_i - \mathbf{w} \cdot \mathbf{x}_i$  (for any support vector)  
 $\mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$
- Classification function:  
 $f(x) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$  If  $f(x) < 0$ , classify as negative,  
if  $f(x) > 0$ , classify as positive  
 $= \text{sign}\left(\sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b\right)$

53

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

---

---

---

---

---

---

---

---

### Questions

- **What if the features are not 2d?**
- What if the data is not linearly separable?
- What if we have more than just two categories?

54

---

---

---

---

---

---

---

---

## Questions

- **What if the features are not 2d?**
  - Generalizes to d-dimensions – replace line with “hyperplane”
- What if the data is not linearly separable?
- What if we have more than just two categories?

---

---

---

---

---

---

---

---

---

---

---

---

### Histograms of Oriented Gradients for Human Detection

Navneet Dalal and Bill Triggs  
INRIA Rhône-Alpes, 655 avenue de l'Europe, Montbonnot 38334, France  
(Navneet.Dalal,Bill.Triggs)@inrialpes.fr, <http://lear.inrialpes.fr>

**Abstract**

We study the question of feature sets for robust visual object recognition, adapting linear SVM based human detection as a test case. After reviewing existing edge and gradient based descriptors, we show experimentally that grids of Histograms of Oriented Gradient (HOG) descriptors significantly outperform existing feature sets for human detection. We study the influence of each stage of the computation on performance, concluding that fine-scale gradients, fine orientation binning, relatively coarse spatial binning, and high-quality local contrast normalization in overlapping descriptor blocks are all important for good results. The new approach gives near-perfect separation on the original MIT pedestrian database, so we introduce a more challenging dataset containing over 1800 unannotated human images with a large range of pose variations and backgrounds.

**1 Introduction**

• CVPR 2005  
• 26000+ citations

**2 Previous Work**

There is an extensive literature on object detection, but here we mention just a few relevant papers on human detection [18, 17, 22, 16, 20]. See [6] for a survey. Papageorgiou *et al* [18] describe a pedestrian detector based on a polynomial SVM using rectified Haar wavelets as input descriptors, with a parts (subwindow) based variant in [17]. Depoortere *et al* give an optimized version of this [2]. Giavetta & Philonen [8] take a more direct approach, extracting edge images and matching them to a set of learned exemplars using chamfer distance. This has been used in a practical real-time pedestrian detection system [7]. Viola *et al* [22] build an efficient

---

---

---

---

---

---

---

---


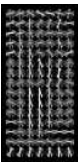
---

---

---

---

## Person detection with HoG's & linear SVM's

- Map each grid cell in the input window to a histogram counting the gradients per orientation.
- Train a linear SVM using training set of pedestrian vs. non-pedestrian windows.

Dalal & Triggs, CVPR 2005

---

---

---

---

---

---

---

---


---

---

---

---

### Person detection with HoG's & linear SVM's



• Histograms of Oriented Gradients for Human Detection, [Navneet Dalal](#), [Bill Triggs](#),  
International Conference on Computer Vision & Pattern Recognition - June 2005 58

---

---

---

---

---

---

---

---

### Questions

- What if the features are not 2d?
- **What if the data is not linearly separable?**
- What if we have more than just two categories?

59

---

---

---

---

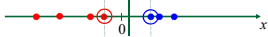
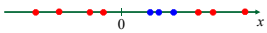
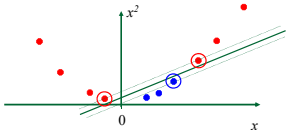
---

---

---

---

### Non-linear SVMs

- Datasets that are linearly separable with some noise work out great: 
- But what are we going to do if the dataset is just too hard? 
- How about... mapping data to a higher-dimensional space: 

60

---

---

---

---

---

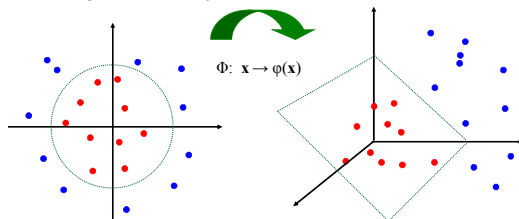
---

---

---

## Non-linear SVMs: feature spaces

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



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

61

## The "Kernel Trick"

- The linear classifier relies on dot product between vectors  $K(x_i, x_j) = x_i^T x_j$

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

62

## Finding the maximum margin line

- Solution:  $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$   
 $b = y_i - \mathbf{w} \cdot \mathbf{x}_i$  (for any support vector)  
 $\mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$

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

63

## The “Kernel Trick”

- The linear classifier relies on dot product between vectors  $K(x_i, x_j) = x_i^T x_j$
- If every data point is mapped into high-dimensional space via some transformation  $\Phi: x \rightarrow \phi(x)$ , the dot product becomes:

$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

- A *kernel function* is similarity function that corresponds to an inner product in some expanded feature space.

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

64

---

---

---

---

---

---

---

---

---

---

## Example

2-dimensional vectors  $x = [x_1 \ x_2]$ ;

let  $K(x_i, x_j) = (1 + x_i^T x_j)^2$

Need to show that  $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$ :

$$\begin{aligned} K(x_i, x_j) &= (1 + x_i^T x_j)^2 \\ &= 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{j1} + 2 x_{i2} x_{j2} \\ &= [1 \ x_{i1}^2 \ \sqrt{2} x_{i1} x_{i2} \ x_{i2}^2 \ \sqrt{2} x_{i1} \ \sqrt{2} x_{i2}]^T \\ &\quad [1 \ x_{j1}^2 \ \sqrt{2} x_{j1} x_{j2} \ x_{j2}^2 \ \sqrt{2} x_{j1} \ \sqrt{2} x_{j2}] \\ &= \phi(x_i)^T \phi(x_j), \\ &\quad \text{where } \phi(x) = [1 \ x_1^2 \ \sqrt{2} x_1 x_2 \ x_2^2 \ \sqrt{2} x_1 \ \sqrt{2} x_2] \end{aligned}$$

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

---

---

---

---

---

---

---

---

---

---

## Nonlinear SVMs

- *The kernel trick*: instead of explicitly computing the lifting transformation  $\phi(x)$ , define a kernel function  $K$  such that

$$K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$

66

---

---

---

---

---

---

---

---

---

---

### Finding the maximum margin line

- Solution:  $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$   
 $b = y_i - \mathbf{w} \cdot \mathbf{x}_i$  (for any support vector)  
 $\mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$

67

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

---

---

---

---

---

---

---

---

---

---

### 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:

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

68

---

---

---

---

---

---

---

---

---

---

### Examples of kernel functions

- Linear:  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- Gaussian RBF:  $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$

- Histogram intersection:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \sum_k \min(x_i(k), x_j(k))$$

69

---

---

---

---

---

---

---

---

---

---

## SVMs for recognition

1. Define your representation for each example.
2. Select a kernel function.
3. Compute pairwise kernel values between labeled examples (i.e., training data).
4. Use this “kernel matrix” to solve for SVM support vectors & weights.
5. To classify a new test example: compute kernel values between new input and support vectors, apply weights, check sign of output.

70

---

---

---

---

---

---

---

---

## Example: learning gender with SVMs

Moghaddam and Yang, Learning Gender with Support Faces, TPAMI 2002.

Moghaddam and Yang, Face & Gesture 2000.

71

---

---

---

---

---

---

---

---

### Face alignment processing

Processed faces

Moghaddam and Yang, Learning Gender with Support Faces, TPAMI 2002.

72

---

---

---

---

---

---

---

---



## Learning gender with SVMs

- Training examples:
  - 1044 males
  - 713 females
- Experiment with various kernels, select Gaussian RBF

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

73

---

---

---

---

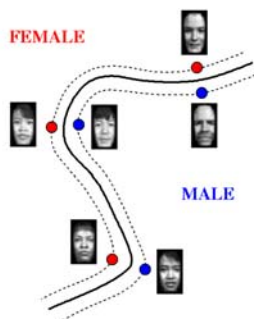
---

---

---

---

## Support Faces



Moghaddam and Yang, Learning Gender with Support Faces, TPAMI 2002.

74

---

---

---

---

---

---

---

---

## Classifier Performance

Classifier	Error Rate		
	Overall	Male	Female
SVM with RBF kernel	3.38%	2.05%	4.79%
SVM with cubic polynomial kernel	4.88%	4.21%	5.59%
Large Ensemble of RBF	5.54%	4.59%	6.55%
Classical RBF	7.79%	6.89%	8.75%
Quadratic classifier	10.63%	9.44%	11.88%
Fisher linear discriminant	13.03%	12.31%	13.78%
Nearest neighbor	27.16%	26.53%	28.04%
Linear classifier	58.95%	58.47%	59.45%

Moghaddam and Yang, Learning Gender with Support Faces, TPAMI 2002.

75

---

---

---

---

---

---

---

---

### Gender perception experiment: How well can humans do?

- Subjects:
  - 30 people (22 male, 8 female)
  - Ages mid-20's to mid-40's
- Test data:
  - 254 face images
  - Low res
  - High res
- Task:
  - Classify as male or female, forced choice
  - No time limit

Moghaddam and Yang, Face & Gesture 2000.

76

---

---

---

---

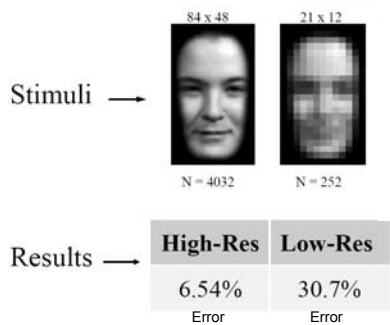
---

---

---

---

### Gender perception experiment: How well can humans do?



Moghaddam and Yang, Face & Gesture 2000.

77

---

---

---

---

---

---

---

---

### Human vs. Machine

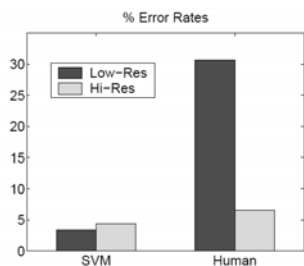


Figure 6. SVM vs. Human performance

78

- SVMs performed better than any single human test subject, at either resolution

---

---

---

---

---

---

---

---

## Hardest examples for humans



Top five human misclassifications

True classification:

Moghaddam and Yang, Face & Gesture 2000.

79

---

---

---

---

---

---

---

---

---

---

## Questions

- What if the features are not 2d?
- What if the data is not linearly separable?
- **What if we have more than just two categories?**

80

---

---

---

---

---

---

---

---

---

---

## Multi-class SVMs

- Achieve multi-class classifier by combining a number of binary classifiers
- **One vs. all**
  - Training: learn an SVM for each class vs. the rest
  - Testing: apply each SVM to test example and assign to it the class of the SVM that returns the highest decision value
- **One vs. one**
  - Training: learn an SVM for each pair of classes
  - Testing: each learned SVM “votes” for a class to assign to the test example

81

---

---

---

---

---

---

---

---

---

---

### SVMs: Pros and cons

- Pros
  - Many publicly available SVM packages:
    - <http://www.kernel-machines.org/software>
    - <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>
  - Kernel-based framework is powerful, flexible
  - Often a sparse set of support vectors – compact at test time
  - Work very well in practice, even with very small training sample sizes
- Cons
  - Can be tricky to select best kernel function for a problem
  - Computation, memory
    - During training time, must compute matrix of kernel values for every pair of examples
    - Learning can take a very long time for large-scale problems

82

Adapted from Leo L. Brezina

---

---

---

---

---

---

---

---

Questions?

See you Tuesday!

83

---

---

---

---

---

---

---

---