

Discriminative classifiers for image recognition

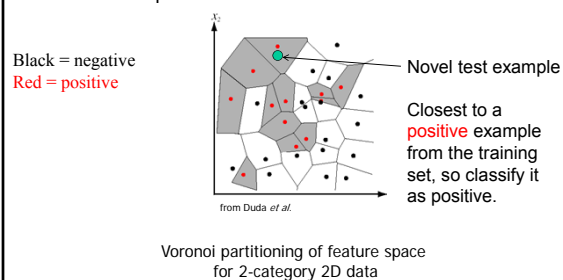
Wednesday, April 13
Kristen Grauman
UT-Austin

Outline

- **Last time:** window-based generic object detection
 - basic pipeline
 - face detection with boosting as case study
- **Today:** discriminative classifiers for image recognition
 - nearest neighbors (+ scene match app)
 - support vector machines (+ gender, person app)

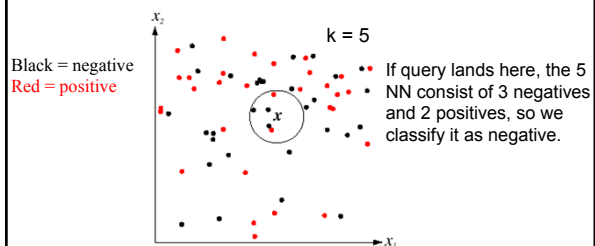
Nearest Neighbor classification

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



K-Nearest Neighbors classification

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



Source: D. Lowe

A nearest neighbor recognition example

Where in the World?



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

Slides: James Hays

Where in the World?



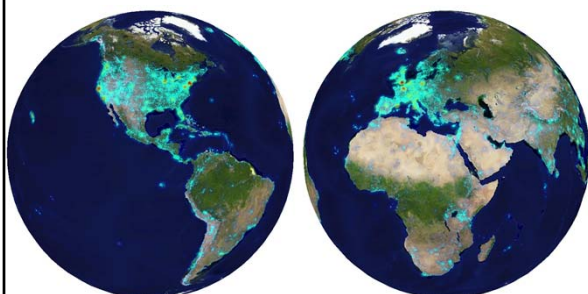
Slides: James Hays

Where in the World?



Slides: James Hays

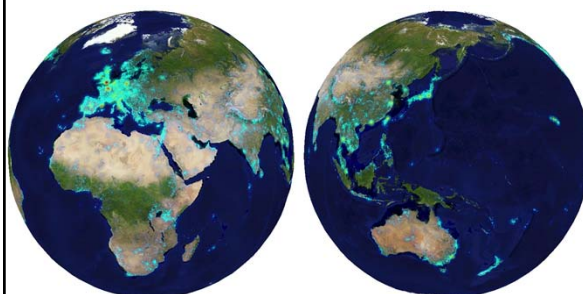
6+ million geotagged photos
by 109,788 photographers



Annotated by Flickr users

Slides: James Hays

6+ million geotagged photos
by 109,788 photographers

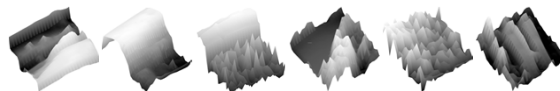


Annotated by Flickr users

Slides: James Hays

Which scene properties are relevant?

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

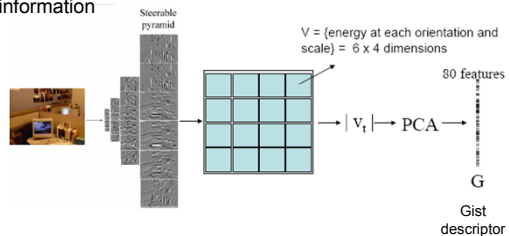


A scene is a single surface that can be
represented by global (statistical) descriptors

Slide Credit: Aude Oliva

Global texture: capturing the “Gist” of the scene

Capture global image properties while keeping some spatial information



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

Which scene properties are relevant?

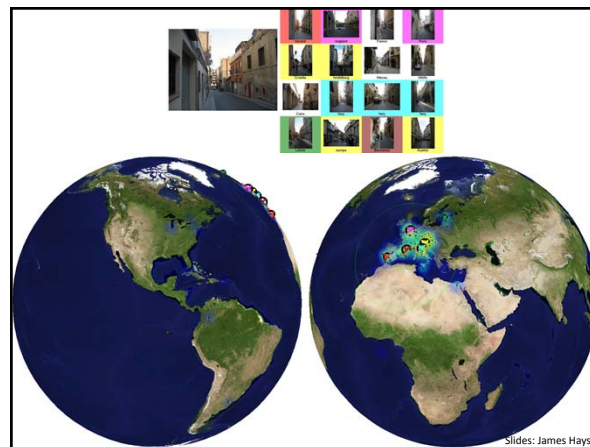
- **Gist scene descriptor**
- **Color Histograms** - $L \times A \times B$ 4x14x14 histograms
- **Texton Histograms** – 512 entry, filter bank based
- **Line Features** – Histograms of straight line stats

Scene Matches



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

Slides: James Hays



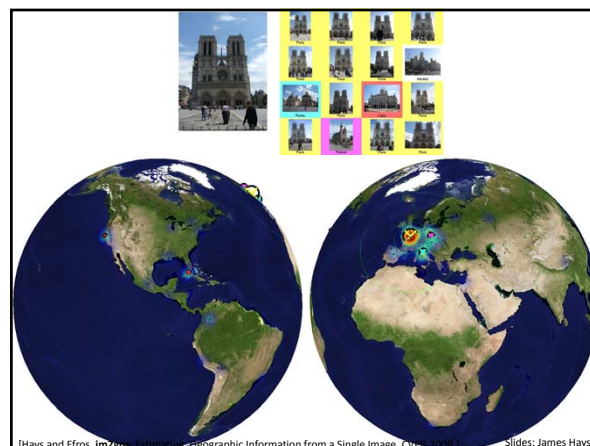
Slides: James Hays

Scene Matches



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

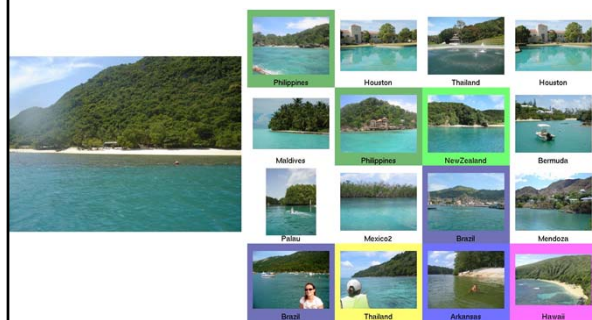
Slides: James Hays



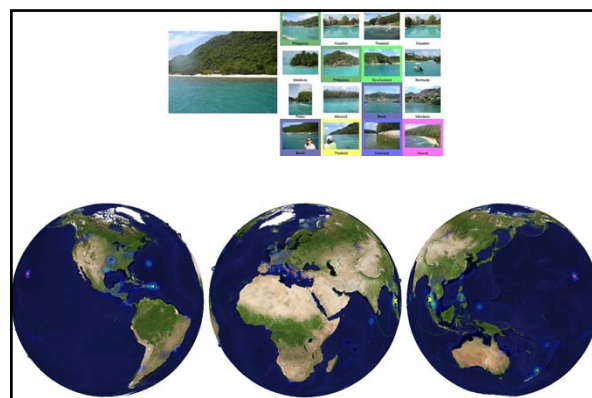
[Hays and Efros, *im2gps*: Estimating Geographic Information from a Single Image, CVPR 2008.]

Slides: James Hays

Scene Matches

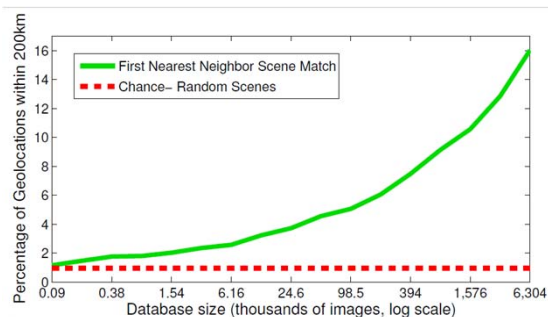


[Hays and Efros. **im2gps**: Estimating Geographic Information from a Single Image. CVPR 2008.] Slides: James Hays



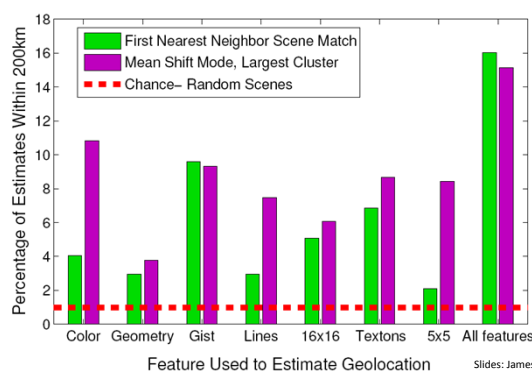
[Hays and Efros. **im2gps**: Estimating Geographic Information from a Single Image. CVPR 2008.] Slides: James Hays

The Importance of Data



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

Feature Performance



Slides: James Hays

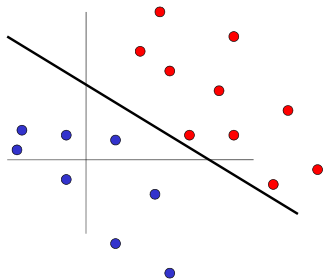
Nearest neighbors: pros and cons

- **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
 - Storage of data
 - Must know we have a meaningful distance function

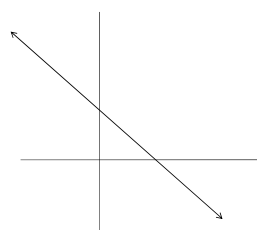
Outline

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

Linear classifiers



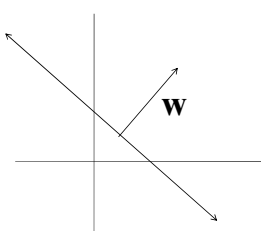
Lines in \mathbb{R}^2



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

$$ax + cy + b = 0$$

Lines in \mathbb{R}^2



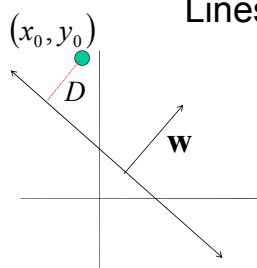
Let $\mathbf{w} = \begin{bmatrix} a \\ c \end{bmatrix}$ $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$

$$ax + cy + b = 0$$

$$\updownarrow$$

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

Lines in \mathbb{R}^2

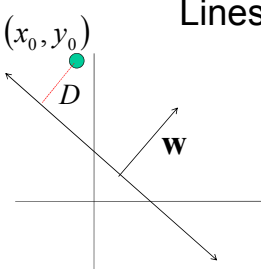


Let $\mathbf{w} = \begin{bmatrix} a \\ c \end{bmatrix}$ $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$

$$ax + cy + b = 0$$

$$\updownarrow$$

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



Lines in \mathbb{R}^2

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

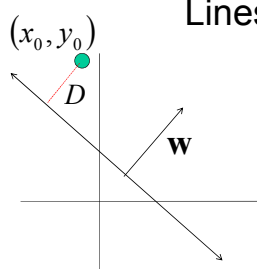
$$ax + cy + b = 0$$

$$\updownarrow$$

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

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

distance from
point to line



Lines in \mathbb{R}^2

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

$$ax + cy + b = 0$$

$$\updownarrow$$

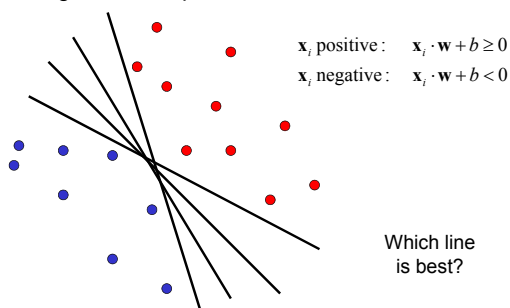
$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

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

distance from
point to line

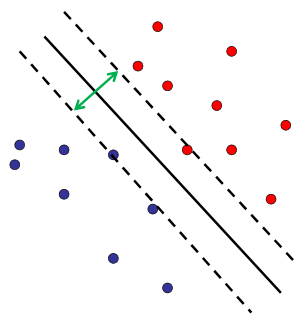
Linear classifiers

- Find linear function to separate positive and negative examples



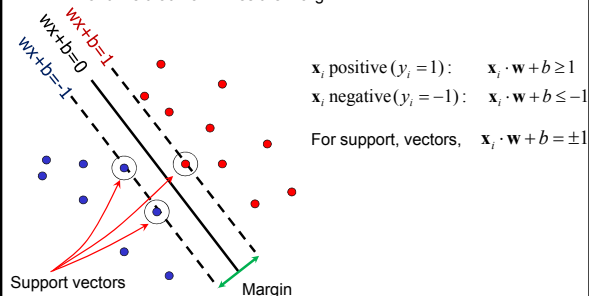
Support Vector Machines (SVMs)

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



Support vector machines

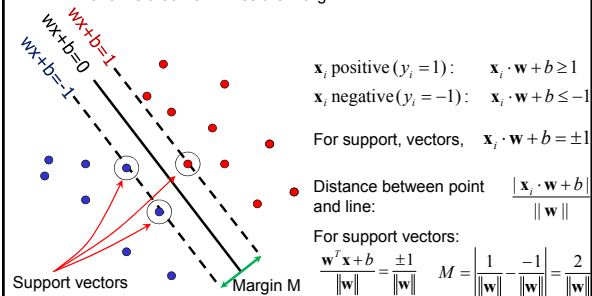
- Want line that maximizes the margin.



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

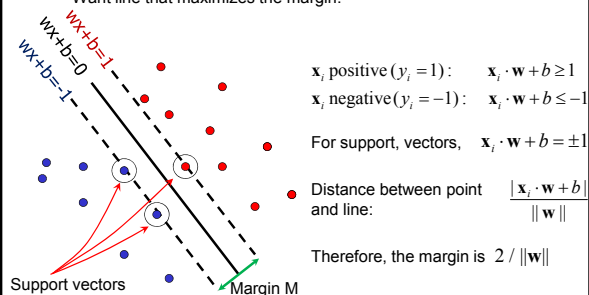
Support vector machines

- Want line that maximizes the margin.



Support vector machines

- Want line that maximizes the margin.



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:

$$\begin{aligned} &\text{Minimize } \frac{1}{2} w^T w \\ &\text{Subject to } y_i(w \cdot x_i + b) \geq 1 \end{aligned}$$

Finding the maximum margin line

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

learned
weight

Support
vector

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:

$$f(x) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b) \quad \begin{array}{l} \text{If } f(x) < 0, \text{ classify} \\ \text{as negative,} \\ \text{if } f(x) > 0, \text{ classify} \\ \text{as positive} \end{array}$$

$$= \text{sign}\left(\sum_i \alpha_i \mathbf{x}_i \cdot \mathbf{x} + b\right)$$

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

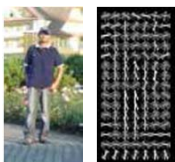
Questions

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

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?

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

Code available:
<http://pascal.inrialpes.fr/soft/ol/>

Person detection with HoG's & linear SVM's



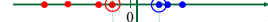
- Histograms of Oriented Gradients for Human Detection, [Navneet Dalal](#), [Bill Triggs](#), International Conference on Computer Vision & Pattern Recognition - June 2005
- <http://lear.inrialpes.fr/pubs/2005/Dalal/>

Questions

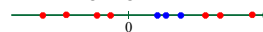
- What if the features are not 2d?
- **What if the data is not linearly separable?**
- What if we have more than just two categories?

Non-linear SVMs

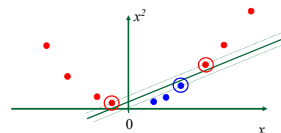
- Datasets that are linearly separable with some noise work out great:



- But what are we going to do if the dataset is just too hard?

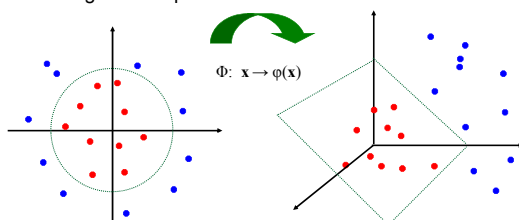


- How about... mapping data to a higher-dimensional space:



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 separable:



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

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>

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

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

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

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

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

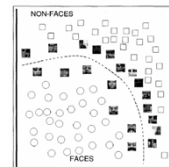
$$\sum_i \alpha_i y_i K(x_i, x) + b$$

Examples of kernel functions

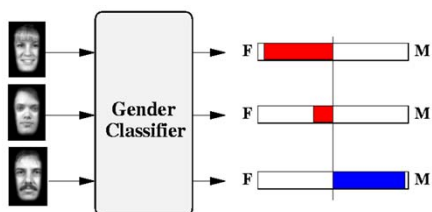
- Linear: $K(x_i, x_j) = x_i^T x_j$
- Gaussian RBF: $K(x_i, x_j) = \exp(-\frac{\|x_i - x_j\|^2}{2\sigma^2})$
- Histogram intersection: $K(x_i, x_j) = \sum_k \min(x_i(k), x_j(k))$

SVMs for recognition

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

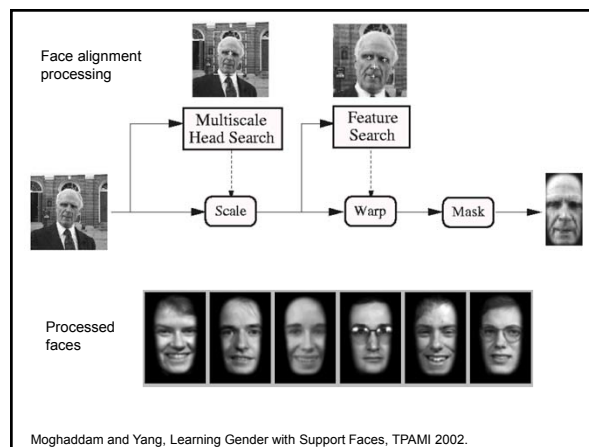


Example: learning gender with SVMs



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

Moghaddam and Yang, Face & Gesture 2000.



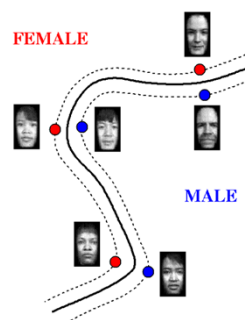
Moghaddam and Yang, Learning Gender with Support Faces, TPAMI 2002.

Learning gender with SVMs

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

$$K(x_i, x_j) = \exp(-\frac{\|x_i - x_j\|^2}{2\sigma^2})$$

Support Faces



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

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%

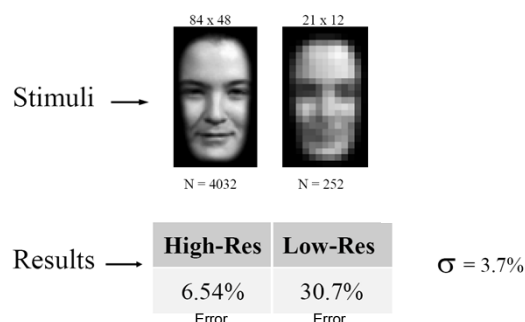
Moghaddam and Yang, Learning Gender with Support Faces, TPAMI 2002.

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 (6 males, 4 females)
 - Low res and high res versions
- Task:
 - Classify as male or female, forced choice
 - No time limit

Moghaddam and Yang, Face & Gesture 2000.

Gender perception experiment: How well can humans do?



Moghaddam and Yang, Face & Gesture 2000.

Human vs. Machine

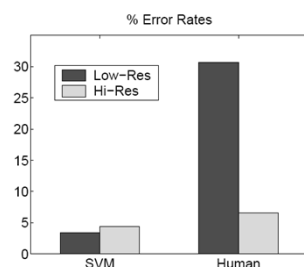
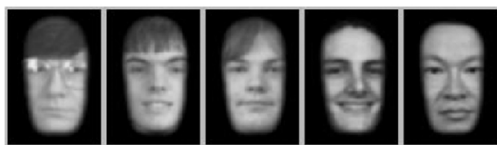


Figure 6. SVM vs. Human performance

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

Hardest examples for humans



Top five human misclassifications

Moghaddam and Yang, Face & Gesture 2000.

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?

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

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 very 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
 - No “direct” multi-class SVM, must combine two-class SVMs
 - 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

Adapted from Leon Lagadec

Coming up

- Part-based models
- Video processing: motion, tracking, activity