# Lecture 18: Recognition IV

Thursday, Nov 15

Prof. Kristen Grauman
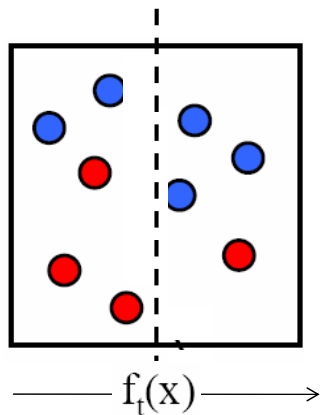
FEMALE

MALE

spade

join

where

# Outline

- Discriminative classifiers
  - SVMs
- Learning categories from weakly supervised images
  - Constellation model
- Shape matching
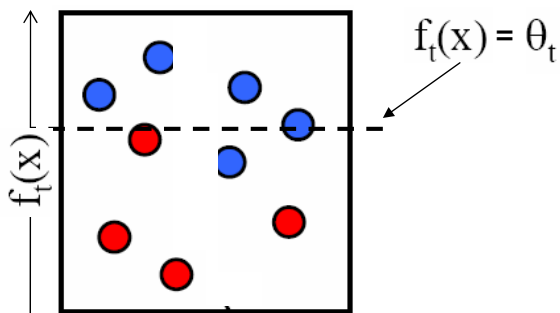  - Shape context, visual CAPTCHA application

# Recall: boosting

- Want to select the single feature that best separates positive and negative examples, in terms of weighted error.



$$\underrightarrow{\qquad f_t(x) \qquad}$$

Each dimension: output of a
possible rectangle feature
on faces and non-faces.

# Recall: boosting

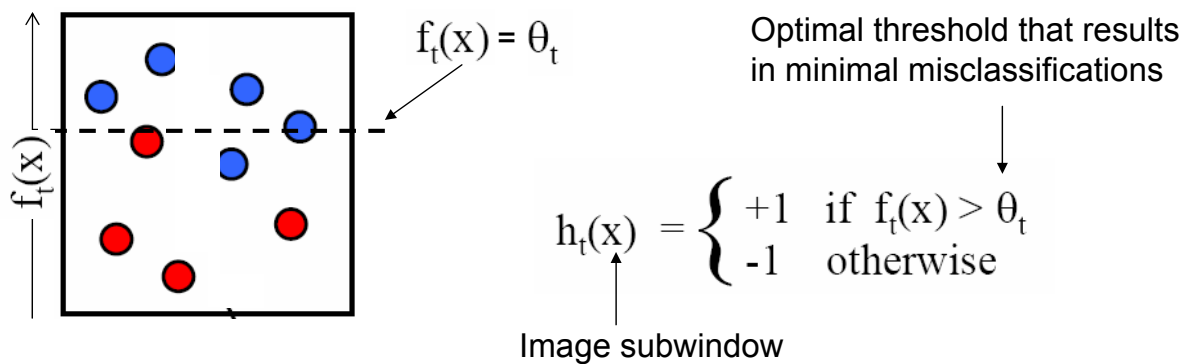- Want to select the single feature that best separates positive and negative examples, in terms of weighted error.



$$f_t(x) = \theta_t$$

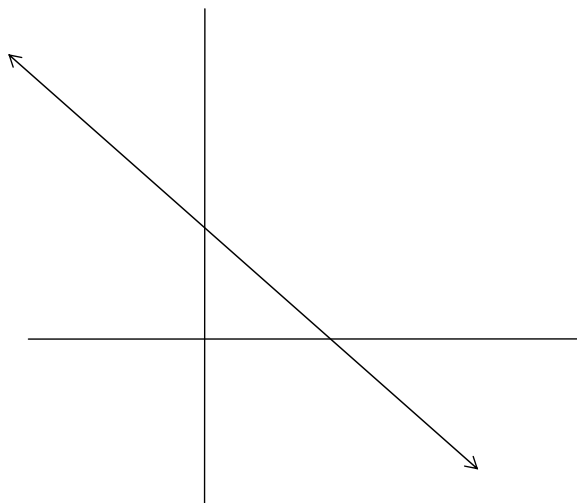Each dimension: output of a possible rectangle feature on faces and non-faces.

# Recall: boosting

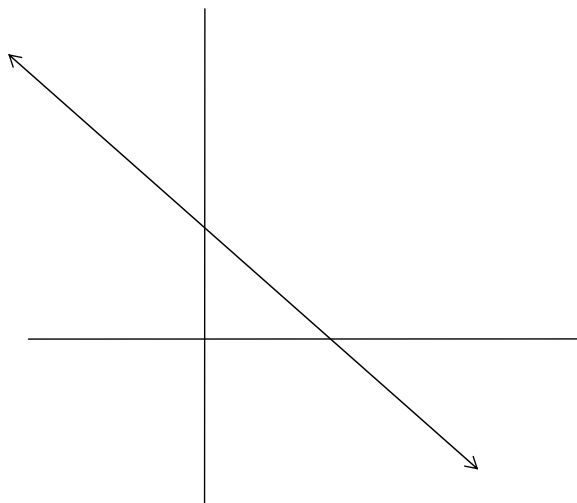- Want to select the single feature that best separates positive and negative examples, in terms of weighted error.



$f_t(x) = \theta_t$

Optimal threshold that results in minimal misclassifications

$$h_t(x) = \begin{cases} +1 & \text{if } f_t(x) > \theta_t \\ -1 & \text{otherwise} \end{cases}$$

Image subwindow

Each dimension: output of a possible rectangle feature on faces and non-faces.

*Notice that any threshold giving same error rate would be equally good here.*

# Lines in R²

$$ax + by + d = 0$$

# Lines in R$^2$

Let $\quad \mathbf{w} = \begin{bmatrix} a \\ b \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$
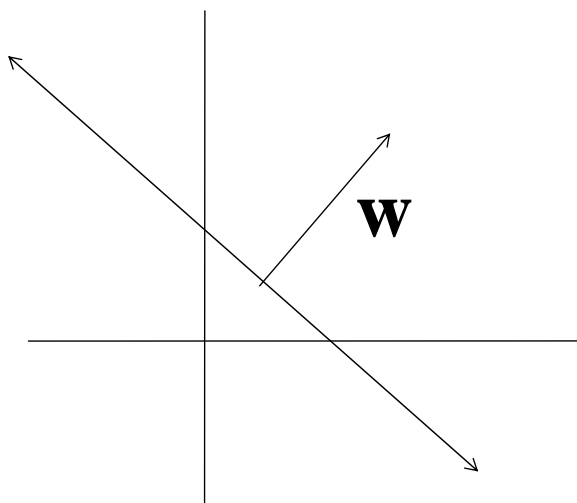
$$ax + by + d = 0$$

# Lines in R$^2$

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

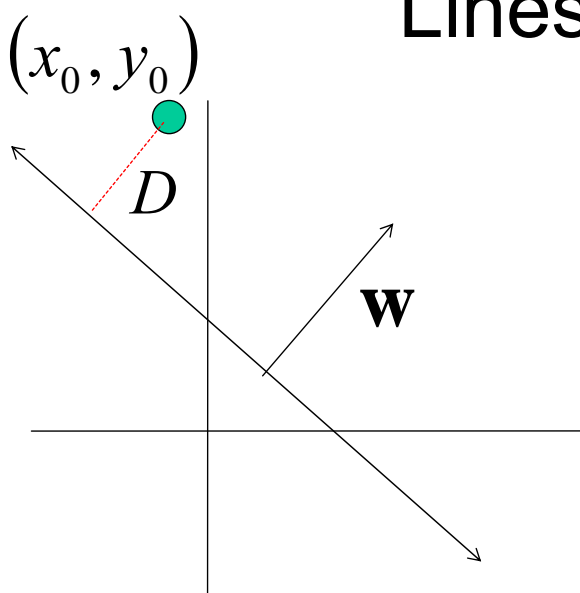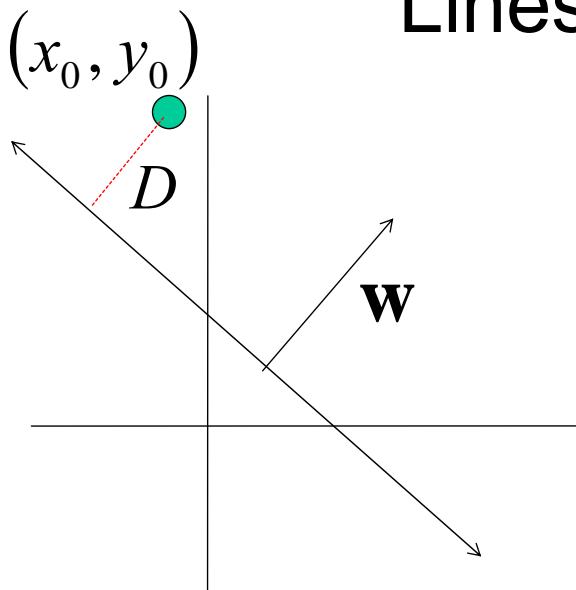$$ax + by + d = 0$$

$$\updownarrow$$

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

$\mathbf{w}$

# Lines in R²

$(x_0, y_0)$

$D$

**w**

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

$$ax + by + d = 0$$

$$\updownarrow$$

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

# Lines in R$^2$

$(x_0, y_0)$

$D$

$\mathbf{W}$

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

$$ax + by + d = 0$$

$$\updownarrow$$

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

$$D = \frac{|ax_0 + by_0 + d|}{\sqrt{a^2 + b^2}}$$

distance from
point to line

# Lines in R$^2$

$(x_0, y_0)$

$D$

$\mathbf{w}$

Let $\mathbf{w} = \begin{bmatrix} a \\ b \end{bmatrix}$ $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$
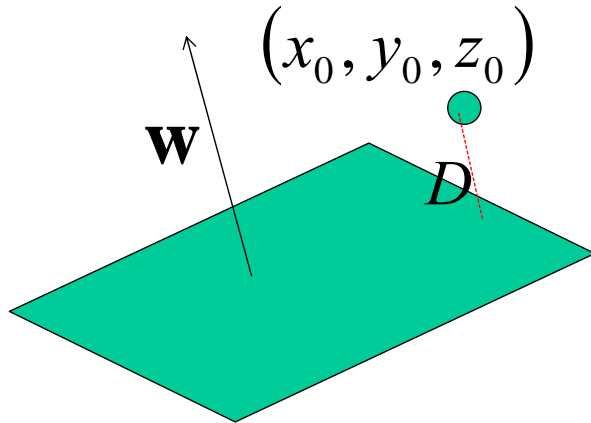
$$ax + by + d = 0$$

$\updownarrow$

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

$$D = \frac{|ax_0 + by_0 + d|}{\sqrt{a^2 + b^2}} = \frac{\mathbf{w}^{\mathrm{T}}\mathbf{x} + d}{|\mathbf{w}|}$$

$\left.\begin{array}{c}\\\\\end{array}\right]$ distance from point to line

# Planes in R$^3$



$(x_0, y_0, z_0)$

**w**

$D$

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

$$ax + by + cz + d = 0$$

$$\updownarrow$$

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

$$D = \frac{\left| ax_0 + by_0 + cz_0 + d \right|}{\sqrt{a^2 + b^2 + c^2}} = \frac{\mathbf{w}^{\mathrm{T}} \mathbf{x} + d}{\left| \mathbf{w} \right|} \left.\right\} \begin{array}{l} \text{distance from} \\ \text{point to plane} \end{array}$$

# Hyperplanes in $R^n$

Hyperplane $H$ is set of all vectors $\mathbf{x} \in R^n$ which satisfy:

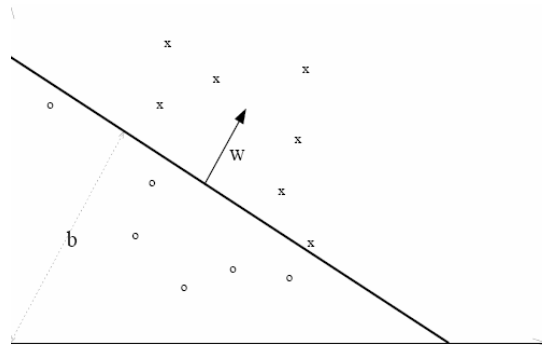$$w_1 x_1 + w_2 x_2 + \ldots + w_n x_n + b = 0$$

$$\updownarrow$$

$$\mathbf{w}^T \mathbf{x} + b = 0$$

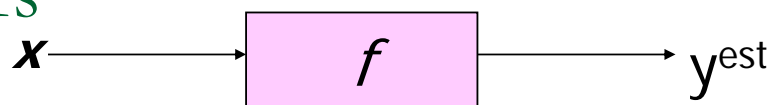$$D(H, \mathbf{x}) = \frac{\mathbf{w}^T \mathbf{x} + b}{|\mathbf{w}|}$$ distance from point to hyperplane
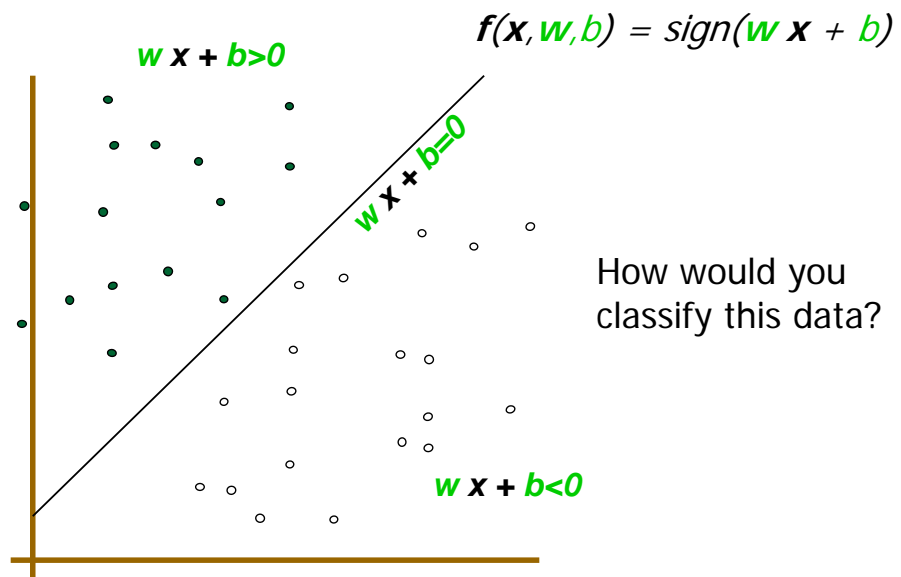
# Support Vector Machines (SVMs)

- Discriminative classifier based on *optimal separating hyperplane*
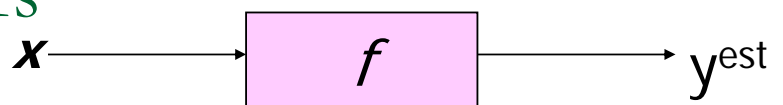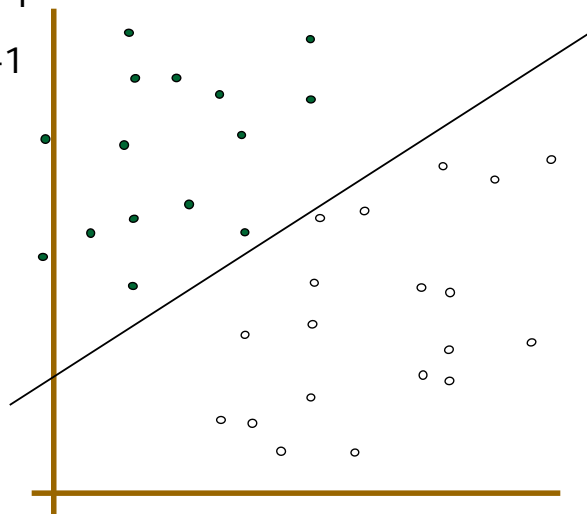


- What hyperplane is optimal?

# Linear Classifiers

$$x \longrightarrow \boxed{f} \longrightarrow y^{est}$$

- denotes +1
- denotes -1

$f(x,w,b) = sign(w\,x + b)$

w x + b>0

w x + b=0

w x + b<0

How would you classify this data?

# Linear Classifiers

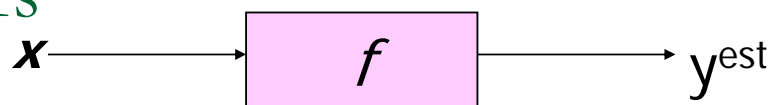$x \longrightarrow \boxed{f} \longrightarrow y^{est}$

$f(x, w, b) = sign(w \ x + b)$

- denotes +1
- denotes -1

How would you classify this data?

# Linear Classifiers

$$x \longrightarrow \boxed{f} \longrightarrow y^{est}$$

$$f(x, w, b) = sign(w\ x + b)$$

- • denotes +1
- ○ denotes -1

How would you
classify this data?

# Linear Classifiers

$x \longrightarrow \boxed{f} \longrightarrow y^{est}$

$f(x, w, b) = sign(w \ x + b)$

- • denotes +1
- ○ denotes -1

Any of these would be fine..

..but which is best?

# Linear Classifiers

$x$ ⟶ $f$ ⟶ $y^{est}$

$f(x,w,b) = sign(w\ x + b)$

- ● denotes +1
- ○ denotes -1

How would you classify this data?

**Misclassified to +1 class**

# Classifier Margin

$x$ → [ $f$ ] → $y^{est}$

$f(x, w, b) = sign(w\ x + b)$

- denotes +1
- denotes -1

Define the margin of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

# Maximum Margin

$x \longrightarrow \boxed{f} \longrightarrow y^{est}$

• denotes +1

○ denotes -1

**Support Vectors** are those datapoints that the margin pushes up against

linear classifier with maximum margin.

This is the simplest kind of SVM (Called an LSVM)

Linear SVM

1. Maximizing the margin is good according to intuition and theory
2. Implies that only support vectors are important; other training examples are ignorable.
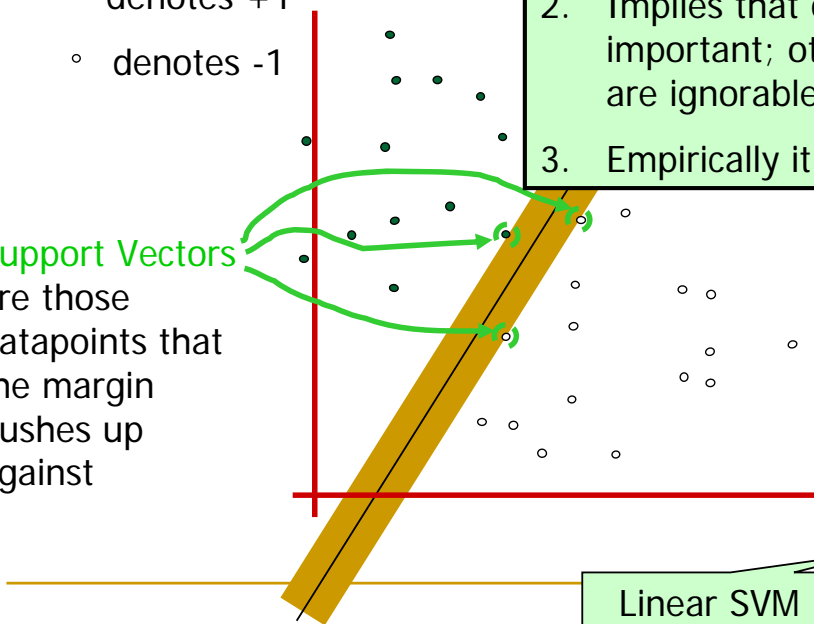3. Empirically it works very very well.

# Linear SVM Mathematically



"Predict Class = +1" zone

**x⁺**

$M$=Margin Width

wx+b=1

wx+b=0

wx+b=-1

**x⁻**

"Predict Class = -1" zone

For the support vectors, distance to hyperplane is 1 for a positives and -1 for negatives.

$$\frac{\mathbf{w}^T\mathbf{x}+b}{|\mathbf{w}|} = \frac{\pm 1}{|\mathbf{w}|}$$

$$M = \left| \frac{1}{|\mathbf{w}|} - \frac{-1}{|\mathbf{w}|} \right| = \frac{2}{|\mathbf{w}|}$$

# Question

- How should we choose values for **w,b?**

  1. want the training data separated by the hyperplane so it classifies them correctly

  2. want the margin width *M* as large as possible

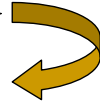# Linear SVM Mathematically

- Goal: **1) Correctly classify all training data**

$$wx_i + b \geq 1 \quad \text{if } y_i = +1$$
$$wx_i + b \leq 1 \quad \text{if } y_i = -1$$
$$y_i(wx_i + b) \geq 1 \quad \text{for all i}$$

**2) Maximize the Margin** $M = \dfrac{2}{|w|}$

**same as minimize** $\dfrac{1}{2}w^t w$

- **Formulated as a Quadratic Optimization Problem, solve for w and b:**

- Minimize $\Phi(w) = \dfrac{1}{2}w^t w$

  subject to $y_i(wx_i + b) \geq 1 \quad \forall i$

# The Optimization Problem Solution

- Solution has the form (omitting derivation):

$$\mathbf{w} = \Sigma \alpha_i y_i \mathbf{x_i} \qquad b = y_k - \mathbf{w^T x_k} \text{ for any } \mathbf{x_k} \text{ such that } \alpha_k \neq 0$$
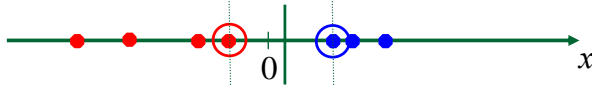
- Each non-zero $\alpha_i$ indicates that corresponding $\mathbf{x_i}$ is a support vector.
- Then the classifying function will have the form:

$$f(\mathbf{x}) = \Sigma \alpha_i y_i \mathbf{x_i^T x} + b$$

- Notice that it relies on an *inner product* between the test point $\mathbf{x}$ and the support vectors $\mathbf{x_i}$
- Solving the optimization problem also involves computing the inner products $\mathbf{x_i^T x_j}$ between all pairs of training points.
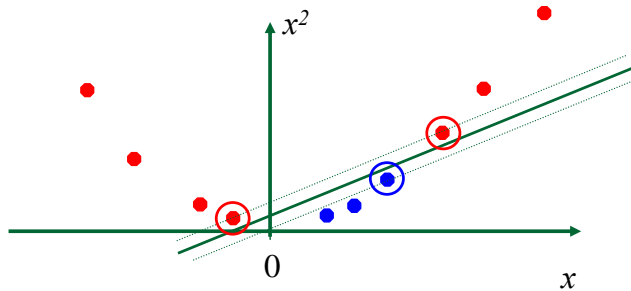
# Non-linear SVMs

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



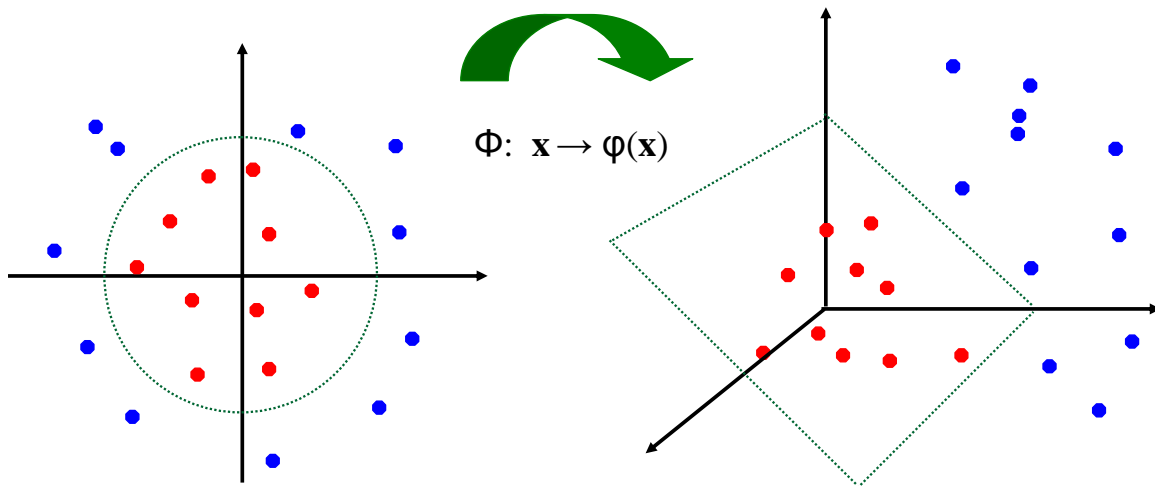- But what are we going to do if the dataset is just too hard?



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

# Non-linear SVMs: Feature spaces

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

$$\Phi: \mathbf{x} \rightarrow \varphi(\mathbf{x})$$

# 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 $\Phi$: $x \rightarrow \varphi(x)$, the dot product becomes:
$$K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j)$$

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

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

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

  $\quad = 1 + x_{i1}^2 x_{j1}^2 + 2\, x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{j1} + 2 x_{i2} x_{j2}$

  $\quad = [1 \ \ x_{i1}^2 \ \ \sqrt{2}\, x_{i1} x_{i2} \ \ x_{i2}^2 \ \ \sqrt{2} x_{i1} \ \ \sqrt{2} x_{i2}]^T \, [1 \ \ x_{j1}^2 \ \ \sqrt{2}\, x_{j1} x_{j2} \ \ x_{j2}^2 \ \ \sqrt{2} x_{j1} \ \ \sqrt{2} x_{j2}]$

  $\quad = \varphi(x_i)^T \varphi(x_j)$,    where $\varphi(x) = [1 \ \ x_1^2 \ \ \sqrt{2}\, x_1 x_2 \ \ x_2^2 \ \ \sqrt{2} x_1 \ \ \sqrt{2} x_2]$
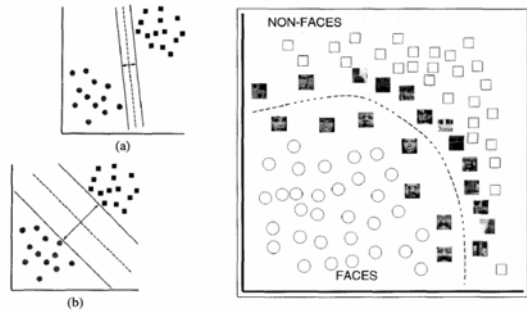
# Examples of General Purpose Kernel Functions

- Linear: $K(\mathbf{x_i}, \mathbf{x_j}) = \mathbf{x_i}^{\mathsf{T}} \mathbf{x_j}$

- Polynomial of power $p$: $K(\mathbf{x_i}, \mathbf{x_j}) = (1 + \mathbf{x_i}^{\mathsf{T}} \mathbf{x_j})^p$

- Gaussian (radial-basis function network):

$$K(\mathbf{x_i}, \mathbf{x_j}) = \exp(-\frac{\left\| \mathbf{x_i} - \mathbf{x_j} \right\|^2}{2\sigma^2})$$
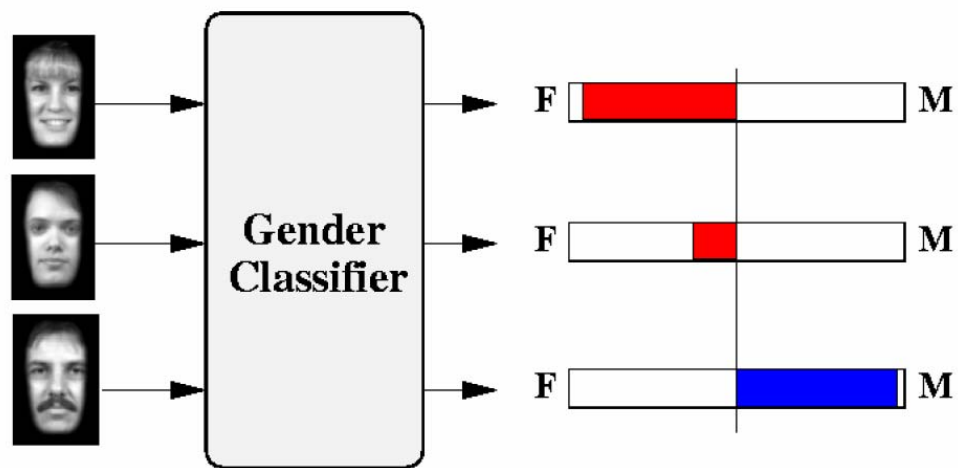
# SVMs for object recognition

1. Define your representation for each example.

2. Select a kernel function.

3. Compute pairwise kernel values between labeled examples, identify support vectors.

4. Compute kernel values between new inputs and support vectors to classify.

# Example: learning gender with SVMs



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

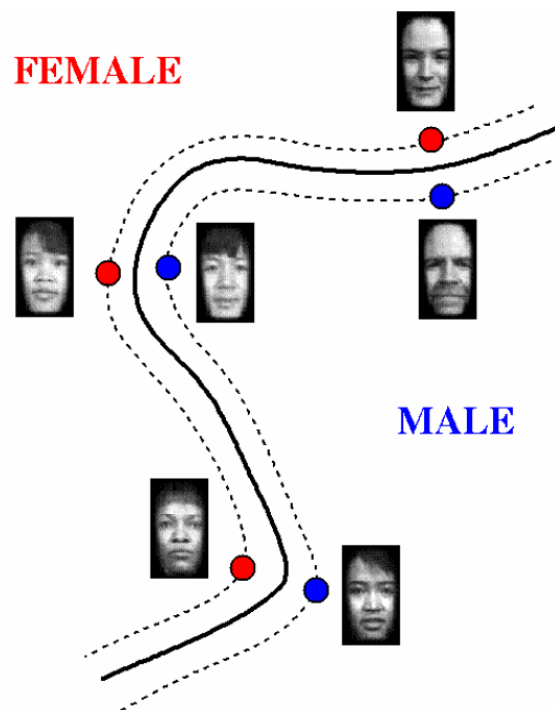Moghaddam and Yang, Face & Gesture 2000.

Face alignment processing

Multiscale Head Search

Feature Search

Scale

Warp

Mask

Processed faces

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

# Support Faces



FEMALE

MALE

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?



Stimuli →

84 x 48          21 x 12

N = 4032          N = 252

Results →

| High-Res | Low-Res |
|----------|---------|
| 6.54% | 30.7% |
| Error | Error |

$\sigma = 3.7\%$

Moghaddam and Yang, Face & Gesture 2000.

# Human vs. Machine



% Error Rates

Figure 6. SVM vs. Human performance

- SVMs perform better than any single human text subject

# Hardest examples for humans



**Top five human misclassifications**

True classification: F, M, M, F, M

# Summary: SVM classifiers

- Discriminative classifier
- Effective for high-dimesional data
- Flexibility/modularity due to kernel
- Very good performance in practice, widely used in vision applications

# Outline

- Discriminative classifiers
  - SVMs
- Learning categories from weakly supervised images
  - Constellation model
- Shape matching
  - Shape context, visual CAPTCHA application

# Weak supervision

- How can we learn object models in the presence of clutter?



Vs.

# Goal

# Weak supervision

- Questions:

  – What about categories where an iconic "template" representation is infeasible?

  – What is the object to be recognized / the part of the image we want to build a model for?

  – For that object, what parts are distinctive or things that can be reliably detected in different instances?

Weber, Welling, Perona. Unsupervised Learning of Models for Recognition, ECCV 2000.

**Fig. 1.** Which objects appear consistently in the left images, but not on the right side? Can a machine learn to recognize instances of the two object classes (*faces* and *cars*) without any further information provided?

What are the features that let us recognize that this is a face?

Slide by Bill Freeman, MIT

A

B

C

D

Slide by Bill Freeman, MIT

# Part-based models

# Part-based models



Main issues:

• measuring the similarity of parts

• representing the

configuration of parts

- Fischler & Elschlager 1973
- Yuille '91
- Brunelli & Poggio '93
- Lades, v.d. Malsburg et al. '93
- Cootes, Lanitis, Taylor et al. '95
- Amit & Geman '95, '99
- Perona et al. '95, '96, '98, '00, '03
- Agarwal & Roth '02

Slide by Fei-Fei Li, 2003.

# One possible constellation model

- Model class with joint probability density function on shape and appearance

Gaussian shape pdf

Gaussian part appearance pdf

mutual positions of the parts, with uncertainty

image patch descriptors, with uncertainty

Figure from Rob Fergus

# Unsupervised learning of part-based models

Main idea:

- Use interest operator to detect small highly textured regions (on both fg and bg)

  - If training objects have similar appearance, these regions will often be similar in different training examples

- Cluster patches: large clusters used to select candidate fg parts

- Choose most informative parts while simultaneously estimating model parameters

  - Iteratively try different combinations of a small number of parts and check model performance on validation set to evaluate quality

Weber, Welling, Perona, ECCV 2000.

# Representation

- Use a scale invariant, scale sensing feature keypoint detector (like the first steps of Lowe's SIFT).



From: Rob Fergus http://www.robots.ox.ac.uk/%7Efergus/

# Features Keys

- A direct appearance model is taken around each located key. This is then normalized to an 11x11 window. PCA further reduces these features.

Unsupervised detector training - 2

"Pattern Space" (100+ dimensions)

Slide by Bill Freeman, MIT

# Candidate parts

For faces         For cars



**Fig. 3.** Points of interest (left) identified on a training image of a human face in cluttered background using Förstner's method. Crosses denote corner-type patterns while circles mark circle-type patterns. A sample of the patterns obtained using k-means clustering of small image patches is shown for faces (center) and cars (right). The car images were high-pass filtered before the part selection process. The total number of patterns selected were 81 for faces and 80 for cars.

At this point, parts appear in both background and foreground of training images.

Weber, Welling, Perona.  Unsupervised Learning of Models for Recognition, 2000.

# Model learning



Which of the candidate parts define the class, and in what configuration?

Let's assume:

• We know number of parts that define the model (and can keep it small).

• Object of interest is only consistent thing somewhere in each training image.

# Model learning



Which of the candidate parts define the class, and in what configuration?

Initialize model parameters randomly.

Iterate while fit improves:

1. Find best assignment in the training images given the parameters

2. Recompute parameters based on current features

# Recognition

- Given a model defining the object class and a model for "background", compute likelihood ratio to make Bayesian decision:

Identified in new image:

X: locations

S: scales

A: appearances

$$R \quad = \quad \frac{p(\text{Object}|\mathbf{X}, \mathbf{S}, \mathbf{A})}{p(\text{No object}|\mathbf{X}, \mathbf{S}, \mathbf{A})}$$

# Recognition

- Given a model defining the object class and a model for "background", compute likelihood ratio to make Bayesian decision:

$$R = \frac{p(\text{Object}|\mathbf{X}, \mathbf{S}, \mathbf{A})}{p(\text{No object}|\mathbf{X}, \mathbf{S}, \mathbf{A})}$$

$$= \frac{p(\mathbf{X}, \mathbf{S}, \mathbf{A}|\text{Object})\, p(\text{Object})}{p(\mathbf{X}, \mathbf{S}, \mathbf{A}|\text{No object})\, p(\text{No object})}$$

X: locations

S: scales

A: appearances

# Example: data from four categories



Slide from Li Fei-Fei http://www.vision.caltech.edu/feifeili/Resume.htm

Face model

Face shape model



Appearance: 10 patches closest to mean for each part

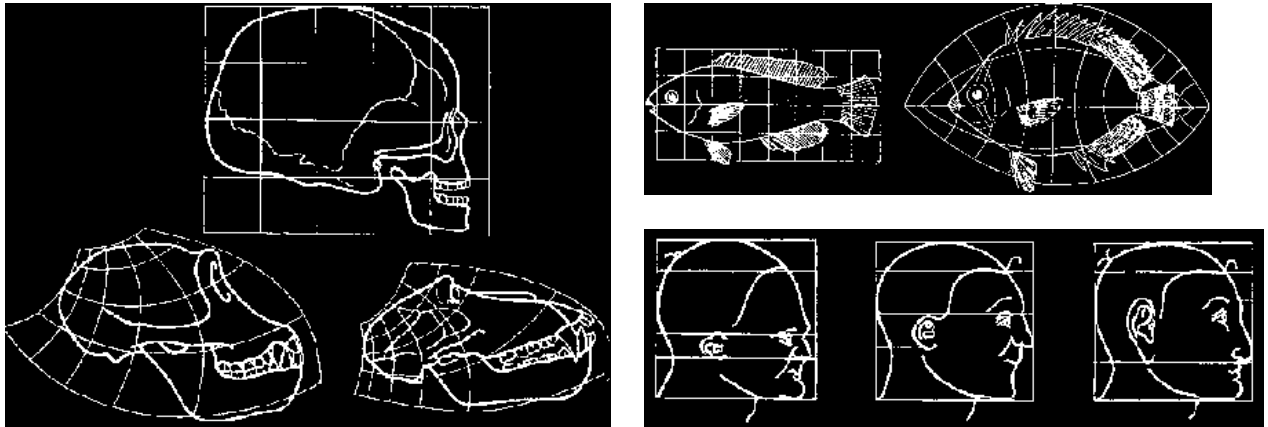Face model

Appearance: 10 patches closest to mean for each part

Recognition results

**Motorbike model**

Motorbike shape model

common features

+ 0.99
+ 0.85
+ 1
+ 1
+ 0.78
+ 0.75

Part 1 – Det:6e–18
Part 2 – Det:6e–22
Part 3 – Det:6e–18
Part 4 – Det:1e–19
Part 5 – Det:3e–17

Background – Det:6e–19

**Appearance: 10 patches closest to mean for each part**

**Recognition results**

Spotted cat model

Spotted cat shape model

Appearance: 10 patches closest to mean for each part

Recognition results

# Outline

- Discriminative classifiers
  - SVMs
- Learning categories from weakly supervised images
  - Constellation model
- Shape matching
  - Shape context, visual CAPTCHA application

# Shape and biology



- D'Arcy Thompson: *On Growth and Form*, 1917
  - studied transformations between shapes of organisms

# Shape matching for recognition



model

target

# Comparing shapes

What points on these two sampled contours are most similar?

# Shape context descriptor
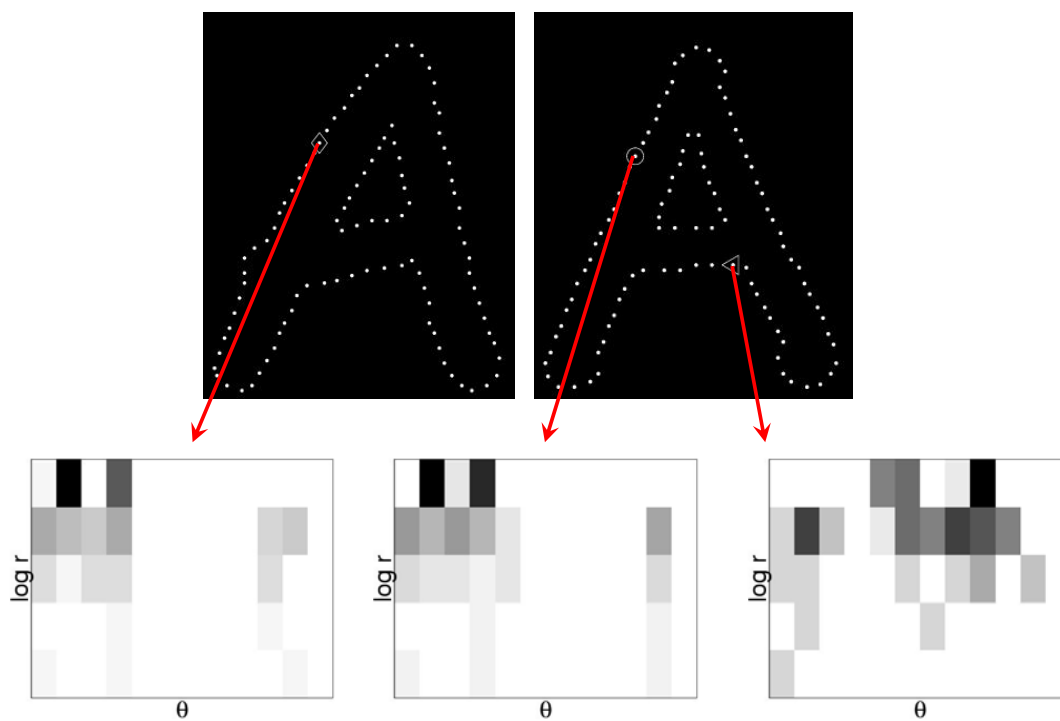


Count the number of points inside each bin, e.g.:
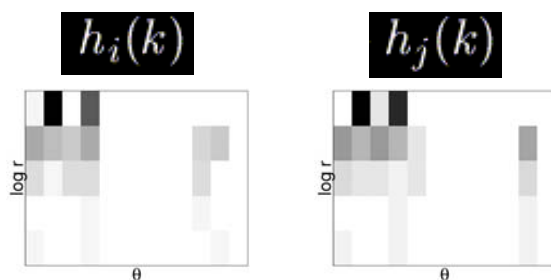
Count = 4

⋮

Count = 10

Compact representation of distribution of points relative to each point
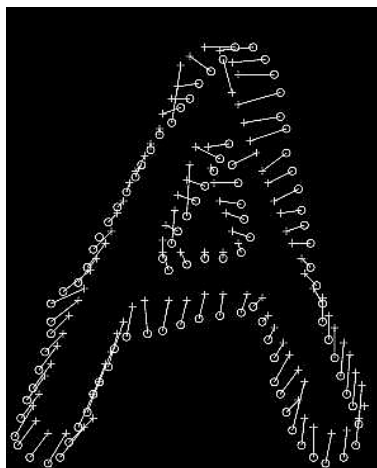
# Shape context descriptor

# Comparing shape contexts

$h_i(k)$      $h_j(k)$



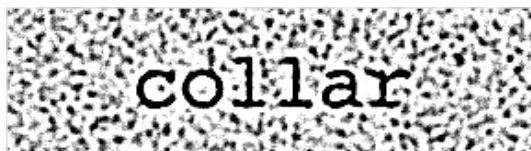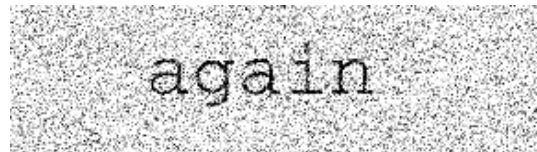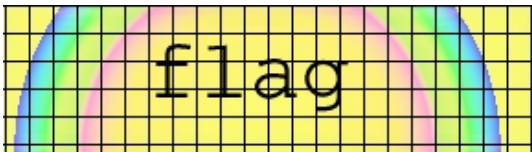$$C_{ij} = \frac{1}{2} \sum_{k=1}^{K} \frac{[h_i(k) - h_j(k)]^2}{h_i(k) + h_j(k)}$$

Recover correspondences by solving for least cost assignment, using costs $C_{ij}$

(Then estimate a parameterized transformation based on these correspondences.)

# CAPTCHA's

- CAPTCHA: Completely Automated Turing Test To Tell Computers and Humans Apart
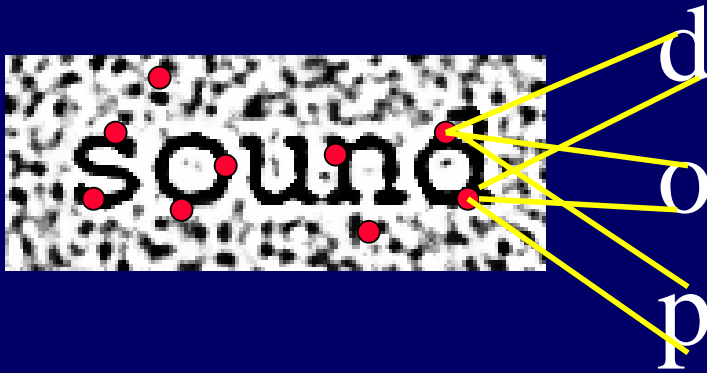- Luis von Ahn, Manuel Blum, Nicholas Hopper and John Langford, CMU, 2000.
- www.captcha.net

# Shape matching application: breaking a visual CAPTCHA

- Use shape matching to recognize characters, words in spite of clutter, warping, etc.

Recognizing Objects in Adversarial Clutter: Breaking a Visual CAPTCHA, by G. Mori and J. Malik, CVPR 2003

# Fast Pruning: Representative Shape Contexts



- Pick k points in the image at random
  - Compare to all shape contexts for all known letters
  - Vote for closely matching letters

- Keep all letters with scores under threshold

# Algorithm A: bottom-up

- Look for letters
  - Representative Shape Contexts

- Find pairs of letters that are "consistent"
  - Letters nearby in space

- Search for valid words

- Give scores to the words

# EZ-Gimpy Results with Algorithm A

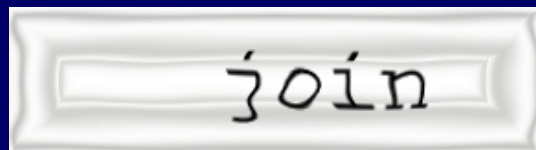- 158 of 191 images correctly identified: 83%
  - Running time: ~10 sec. per image (MATLAB, 1 Ghz P3)


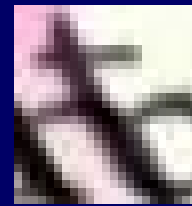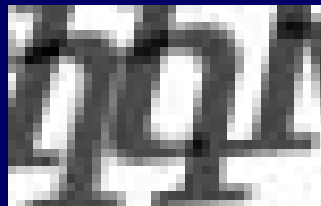
horse
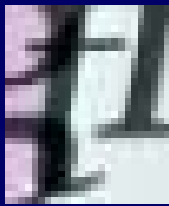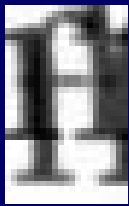


spade



smile



join



canvas



here

# Gimpy



- Multiple words, task is to find 3 words in the image

- Clutter is other objects, not texture
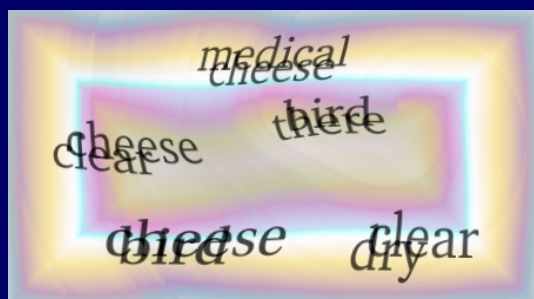
# Algorithm B: Letters are not enough



- Hard to distinguish single letters with so much clutter

- Find words instead of letters
  - Use long range info over entire word
  - Stretch shape contexts into ellipses



- Search problem becomes huge
  - # of words 600 vs. # of letters 26
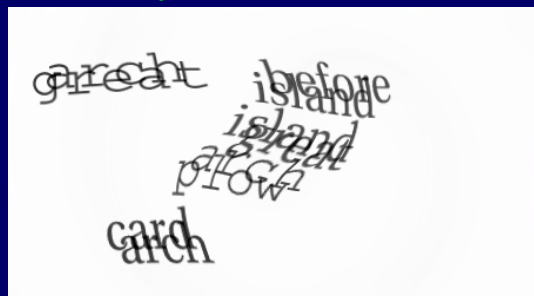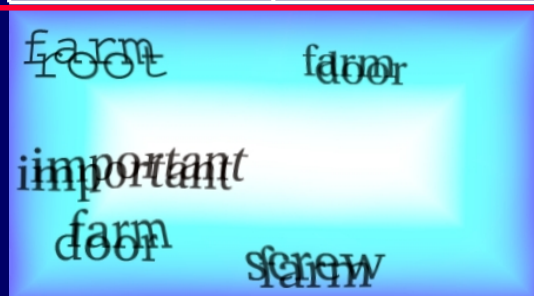  - Prune set of words using opening/closing bigrams

# Results with Algorithm B



dry clear medical



card arch plate

| # Correct words | % tests (of 24) |
|---|---|
| 1 or more | 92% |
| 2 or more | 75% |
| 3 | 33% |
| EZ-Gimpy | 92% |



door farm important

# Coming up

- Face images

- For next week:
  - Read Trucco & Verri handout on Motion

- Problem set 4 due 11/29

# References

- Unsupervised Learning of Models for Recognition, by M. Weber, M. Welling and P. Perona, ECCV 2000.
- Towards Automatic Discovery of Object Categories, by M. Weber, M. Welling and P. Perona, CVPR 2000.
- Object Class Recognition by Unsupervised Scale-Invariant Learning, by Fergus, Perona, and Zisserman, CVPR 2003.
- Matching Shapes, by S. Belongie, J. Malik and J. Puzicha, ICCV 2001.
- Recognizing Objects in Adversarial Clutter: Breaking a Visual CAPTCHA, by G. Mori and J. Malik, CVPR 2003.
- Learning Gender with Support Faces, by Moghaddam and Yang, TPAMI, 2002.
- SVM slides from Andrew Moore, CMU