

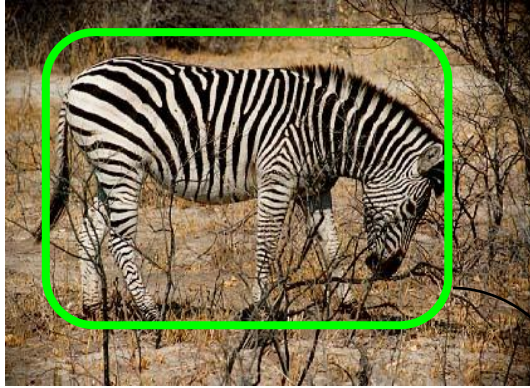
Classifier based methods for Object Recognition

CMP 719– Computer Vision

Pinar Duygulu

(Slide credits:

Kristen Grauman, Fei fei Li, Antonio Torralba, Hames Hays)



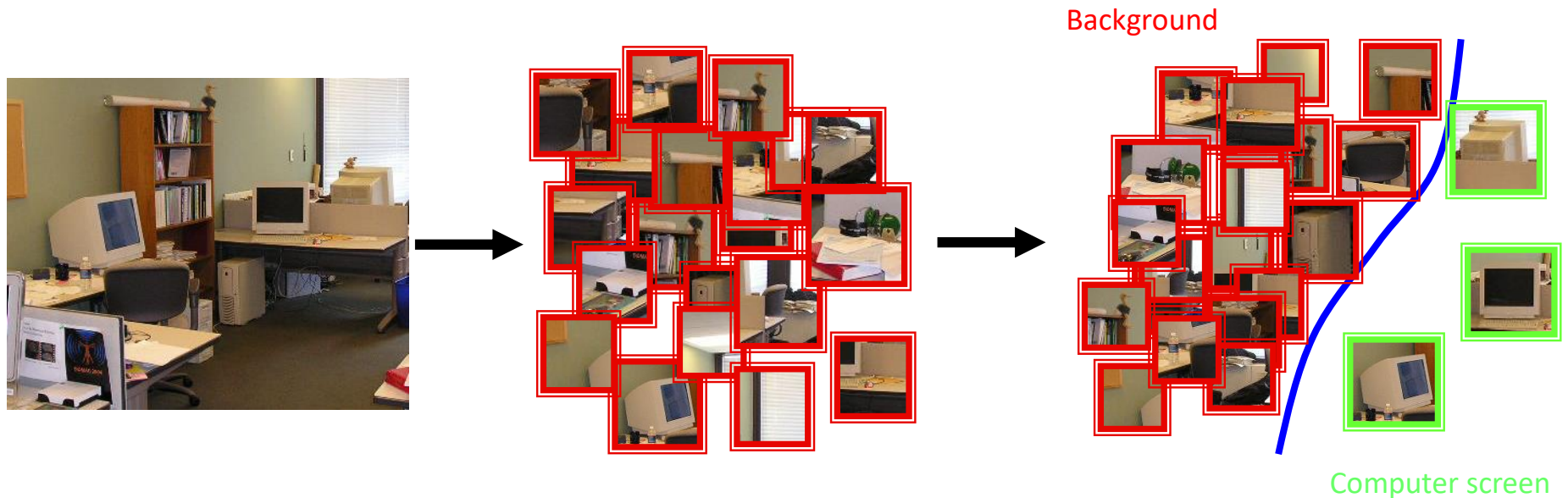
Classifier-based methods

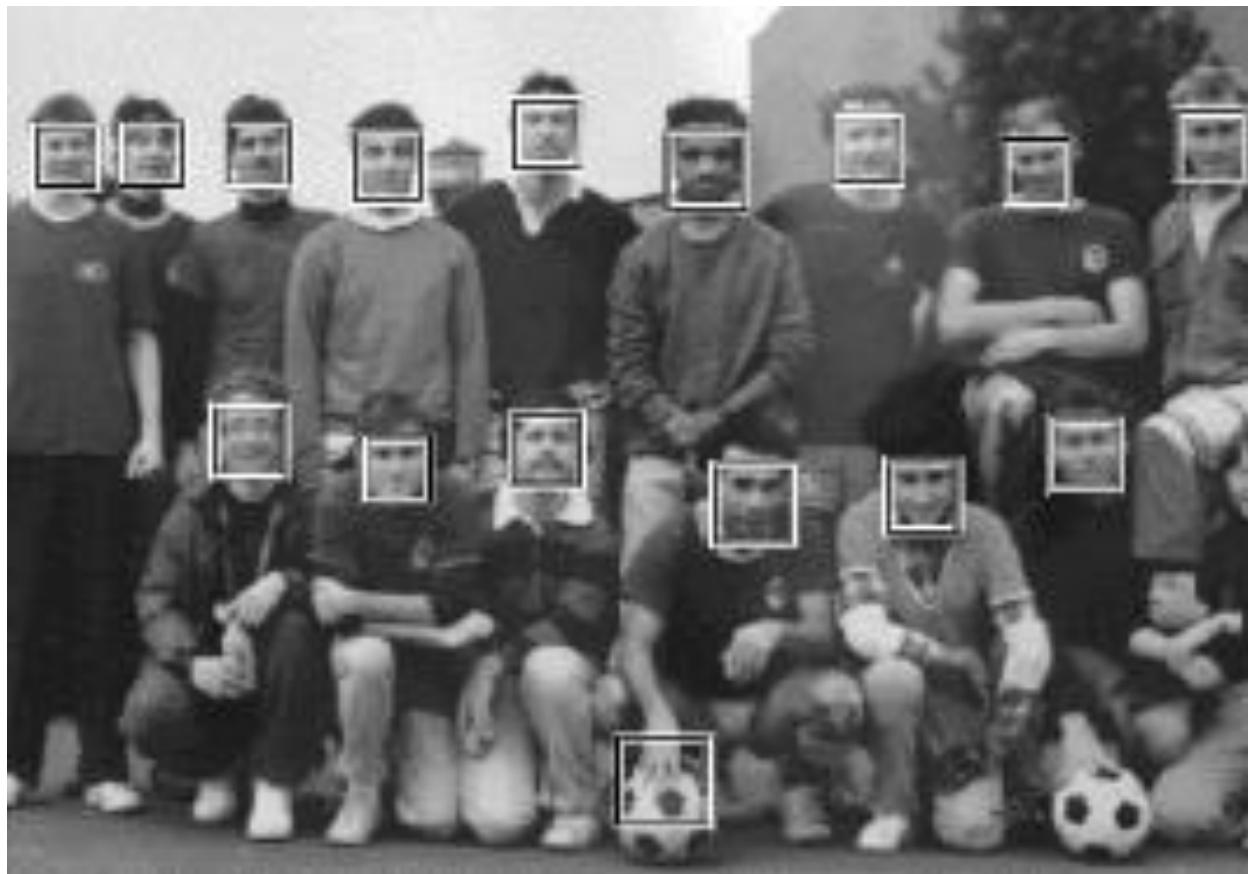
Classifier based methods

Object detection and recognition is formulated as a classification problem.

The image is partitioned into a set of overlapping windows

... and a decision is taken at each window about if it contains a target object or not.





Learning Models

Training

Training Images



Image Features



Training Labels



Training



Learned model

Testing



Test Image



Image Features



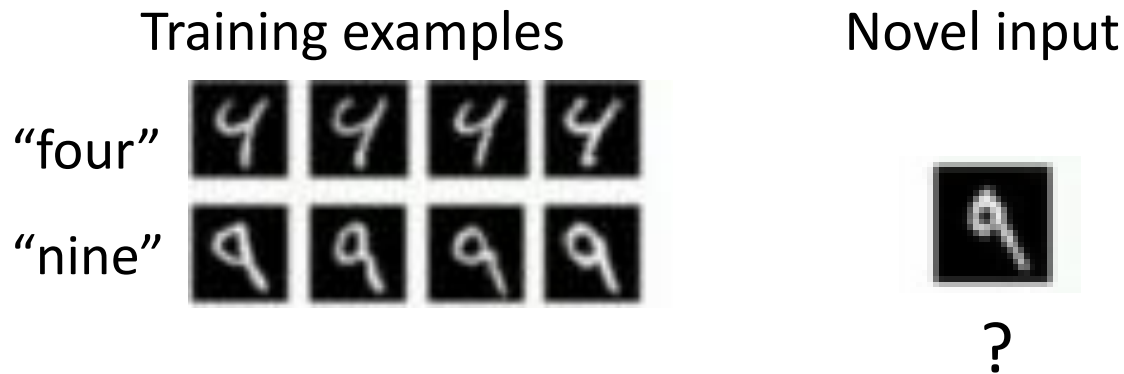
Prediction

Learned model



Supervised classification

- Given a collection of *labeled* examples, come up with a function that will predict the labels of new examples.



- How good is some function we come up with to do the classification?
- Depends on
 - Mistakes made
 - Cost associated with the mistakes

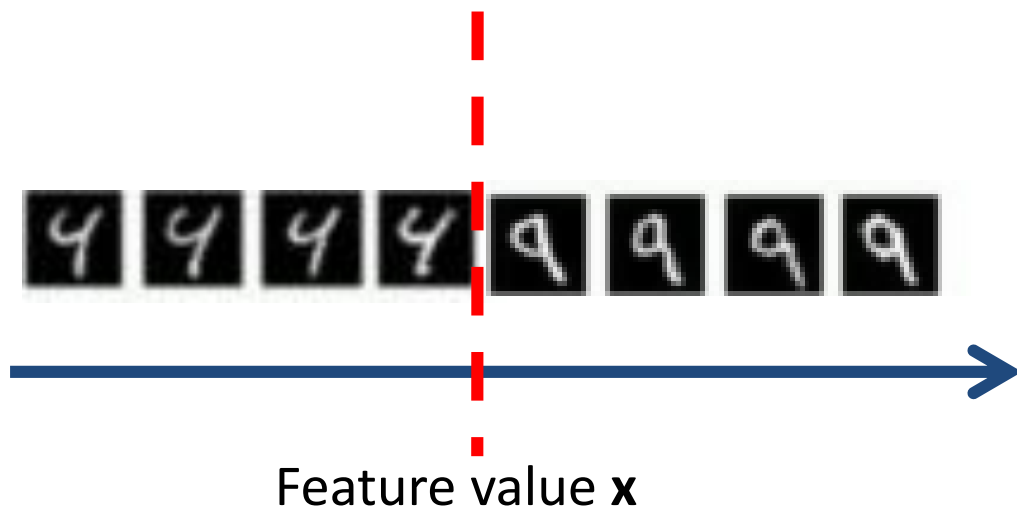
Supervised classification

- Given a collection of *labeled* examples, come up with a function that will predict the labels of new examples.
- Consider the two-class (binary) decision problem
 - $L(4 \rightarrow 9)$: Loss of classifying a 4 as a 9
 - $L(9 \rightarrow 4)$: Loss of classifying a 9 as a 4
- **Risk** of a classifier s is expected loss:

$$R(s) = \Pr(4 \rightarrow 9 \mid \text{using } s)L(4 \rightarrow 9) + \Pr(9 \rightarrow 4 \mid \text{using } s)L(9 \rightarrow 4)$$

- We want to choose a classifier so as to minimize this total risk

Supervised classification



Optimal classifier will minimize total risk.

At decision boundary, either choice of label yields same expected loss.

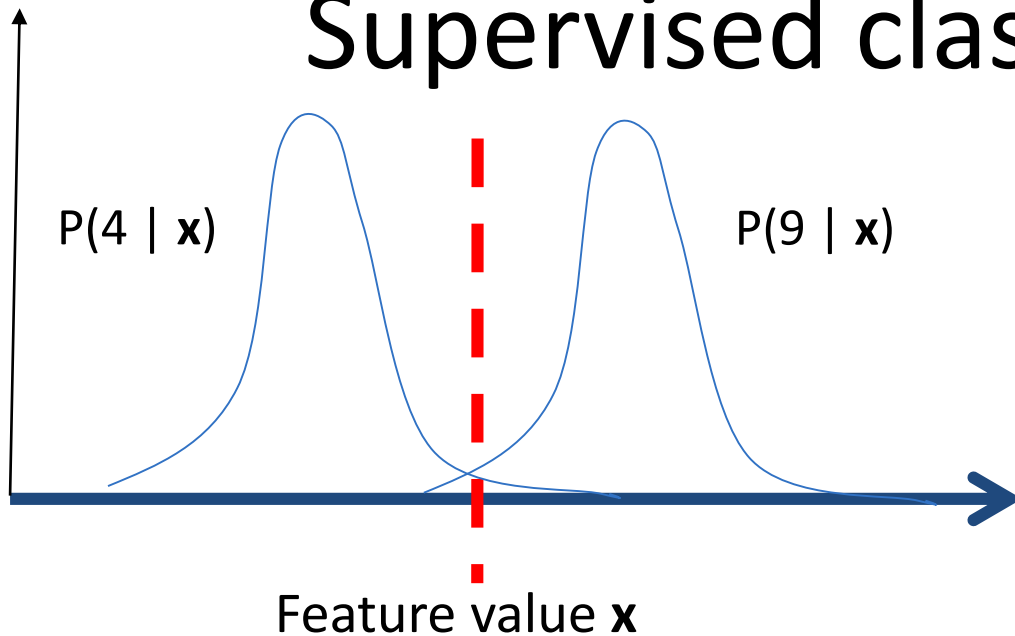
So, best decision boundary is at point \mathbf{x} where

$$P(\text{class is } 9 \mid \mathbf{x}) L(9 \rightarrow 4) = P(\text{class is } 4 \mid \mathbf{x}) L(4 \rightarrow 9)$$

To classify a new point, choose class with lowest expected loss; i.e., choose “four” if

$$P(4 \mid \mathbf{x}) L(4 \rightarrow 9) > P(9 \mid \mathbf{x}) L(9 \rightarrow 4)$$

Supervised classification



Optimal classifier will minimize total risk.

At decision boundary, either choice of label yields same expected loss.

So, best decision boundary is at point \mathbf{x} where

$$P(\text{class is } 9 | \mathbf{x}) L(9 \rightarrow 4) = P(\text{class is } 4 | \mathbf{x}) L(4 \rightarrow 9)$$

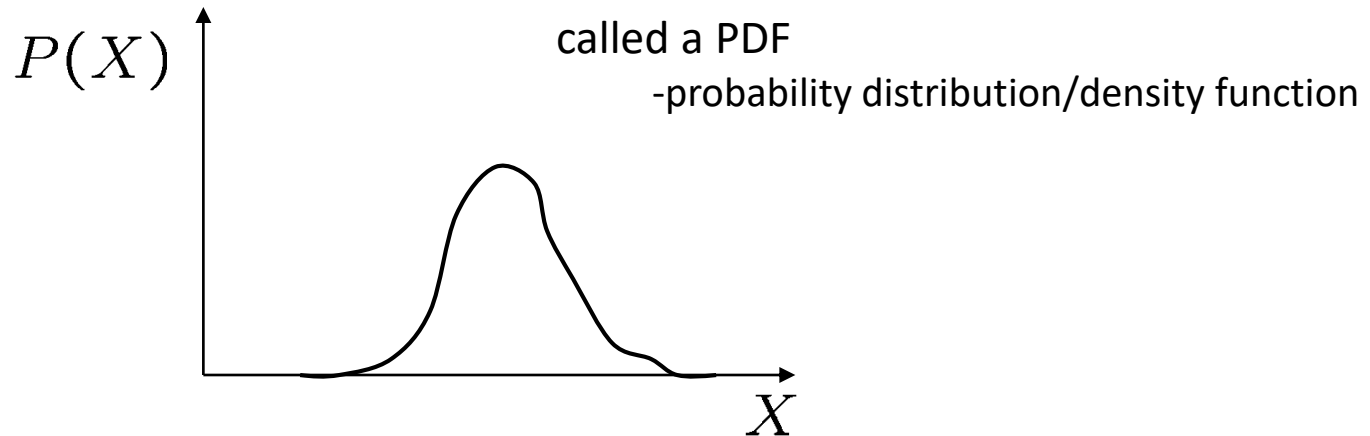
To classify a new point, choose class with lowest expected loss; i.e., choose “four” if

$$P(4 | \mathbf{x}) L(4 \rightarrow 9) > P(9 | \mathbf{x}) L(9 \rightarrow 4)$$

How to evaluate these probabilities?

Probability

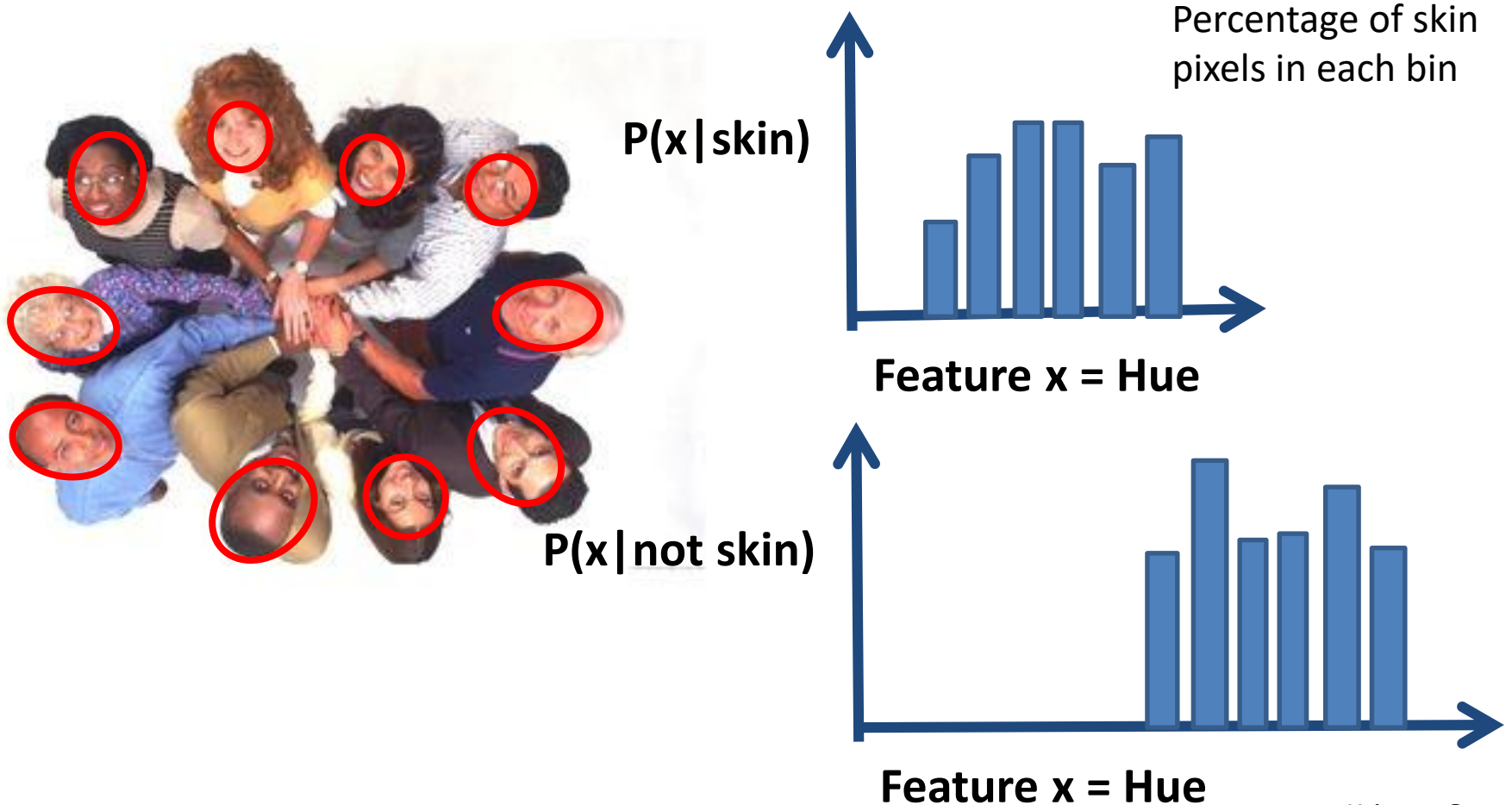
- Basic probability
 - X is a random variable
 - $P(X)$ is the probability that X achieves a certain value



- $0 \leq P(X) \leq 1$
- $\int_{-\infty}^{\infty} P(X) dX = 1$ or $\sum P(X) = 1$
continuous X discrete X
- Conditional probability: $P(X | Y)$
 - probability of X given that we already know Y

Example: learning skin colors

- We can represent a class-conditional density using a histogram (a “non-parametric” distribution)



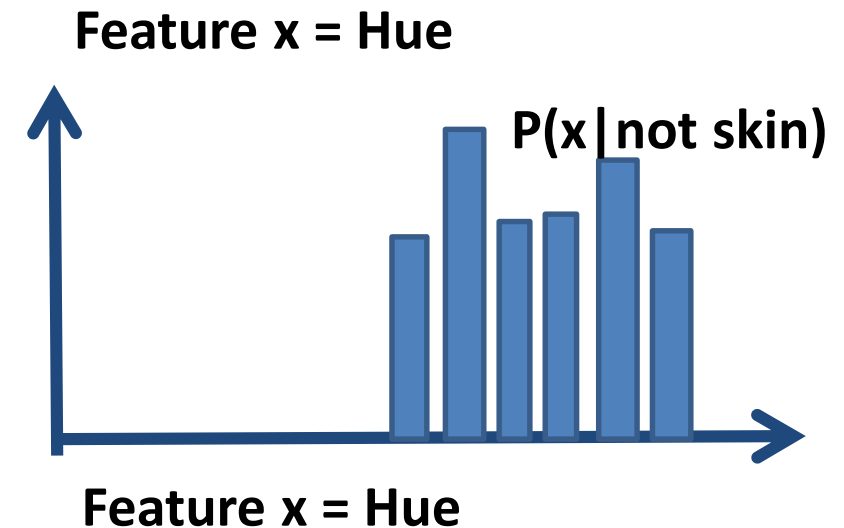
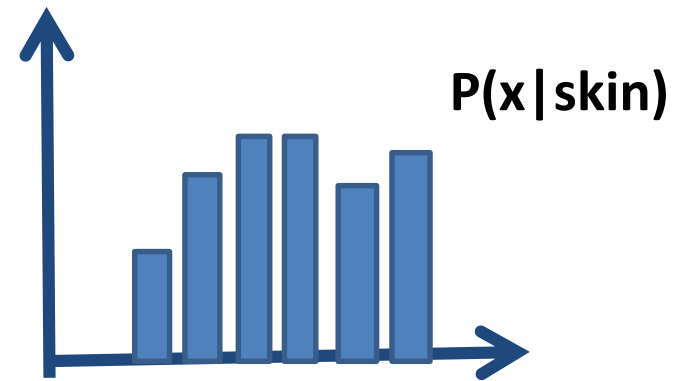
Example: learning skin colors

- We can represent a class-conditional density using a histogram (a “non-parametric” distribution)

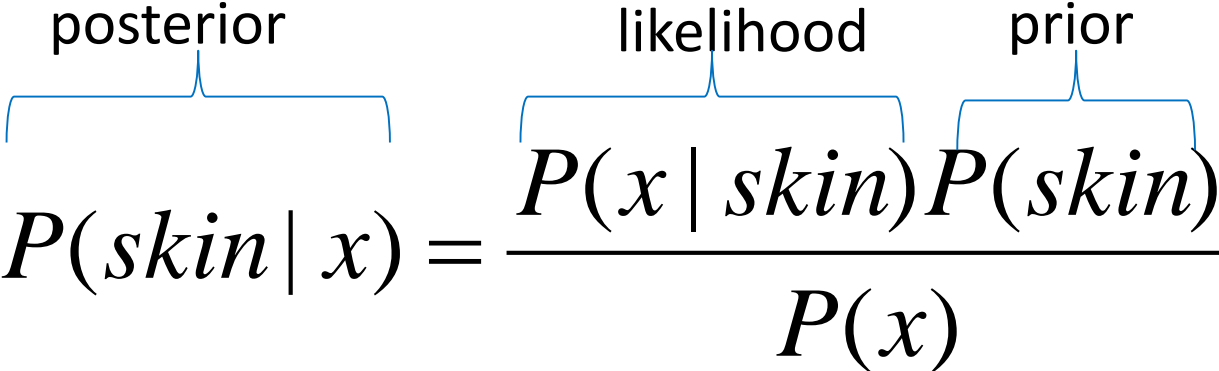


Now we get a new image, and want to label each pixel as skin or non-skin.

What's the probability we care about to do skin detection?



Bayes rule



The diagram illustrates Bayes' rule with three labels above the equation: 'posterior' above the left side, 'likelihood' above the first term of the numerator, and 'prior' above the second term of the numerator. Blue curly braces connect these labels to their respective parts of the equation.

$$\overbrace{P(\textit{skin} \mid x)}^{\text{posterior}} = \frac{\overbrace{P(x \mid \textit{skin})}^{\text{likelihood}} \overbrace{P(\textit{skin})}^{\text{prior}}}{P(x)}$$

$$P(\textit{skin} \mid x) \propto P(x \mid \textit{skin})P(\textit{skin})$$

Example: classifying skin pixels

Now for every pixel in a new image, we can estimate probability that it is generated by skin.



Brighter pixels →
higher probability
of being skin

Classify pixels based on these probabilities

- if $p(\text{skin}|\mathbf{x}) > \theta$, classify as skin
- if $p(\text{skin}|\mathbf{x}) < \theta$, classify as not skin

Example: classifying skin pixels

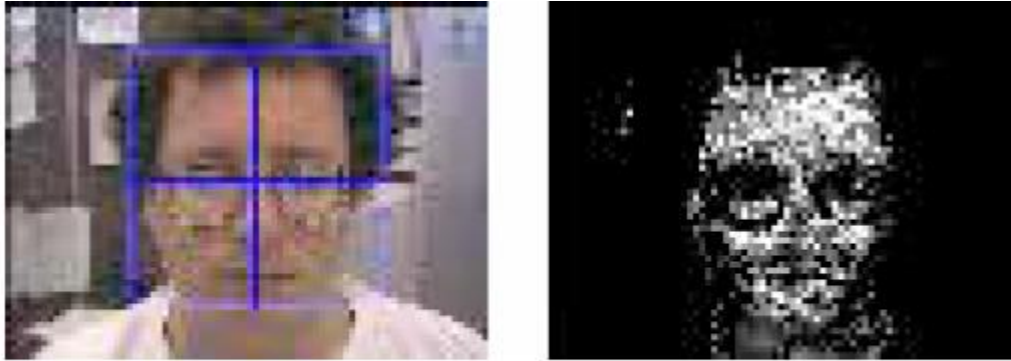


Figure 6: A video image and its flesh probability image



Figure 7: Orientation of the flesh probability distribution marked on the source video image

Example: classifying skin pixels

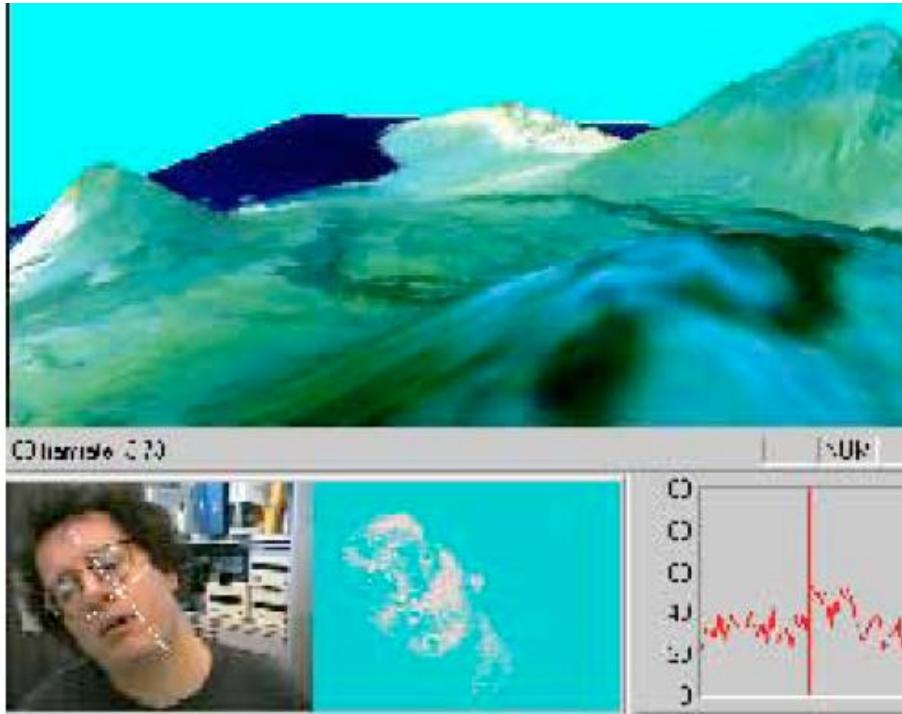


Figure 13: CAMSHIFT-based face tracker used to play Quake 2 hands free by inserting control variables into the mouse queue

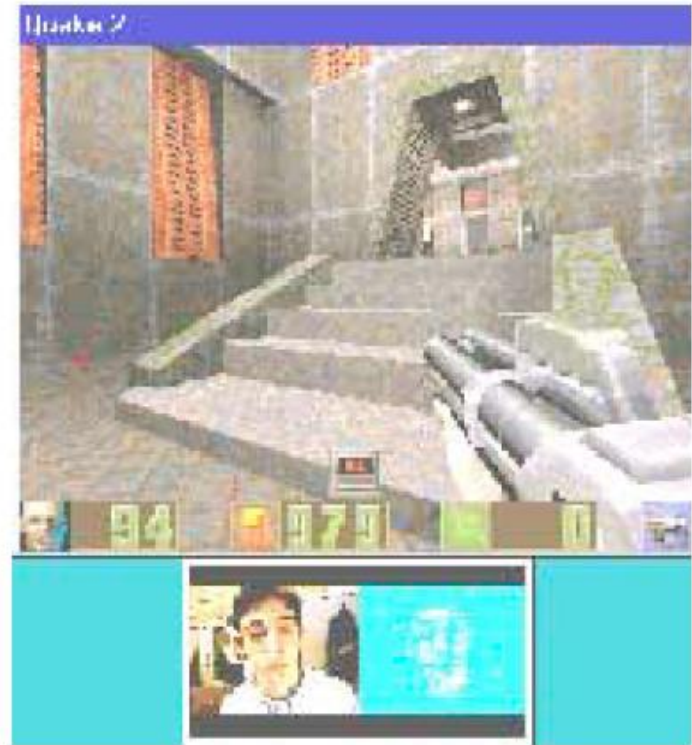


Figure 12: CAMSHIFT-based face tracker used to play Quake 2 hands free by inserting control variables into the mouse queue

Using skin color-based face detection and pose estimation as a video-based interface

Supervised classification

- Want to minimize the expected misclassification
- Two general strategies
 - Use the training data to build representative probability model; separately model class-conditional densities and priors (*generative*)
 - Directly construct a good decision boundary, model the posterior (*discriminative*)

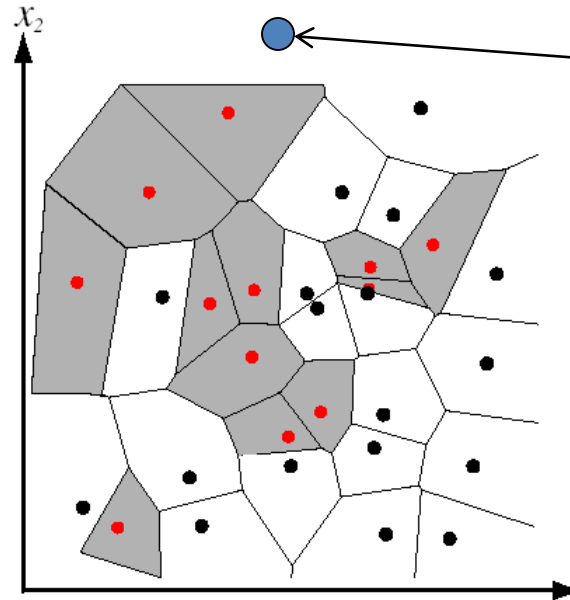
Discriminative classifiers for image recognition

- 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

Black = negative
Red = positive



Novel test example

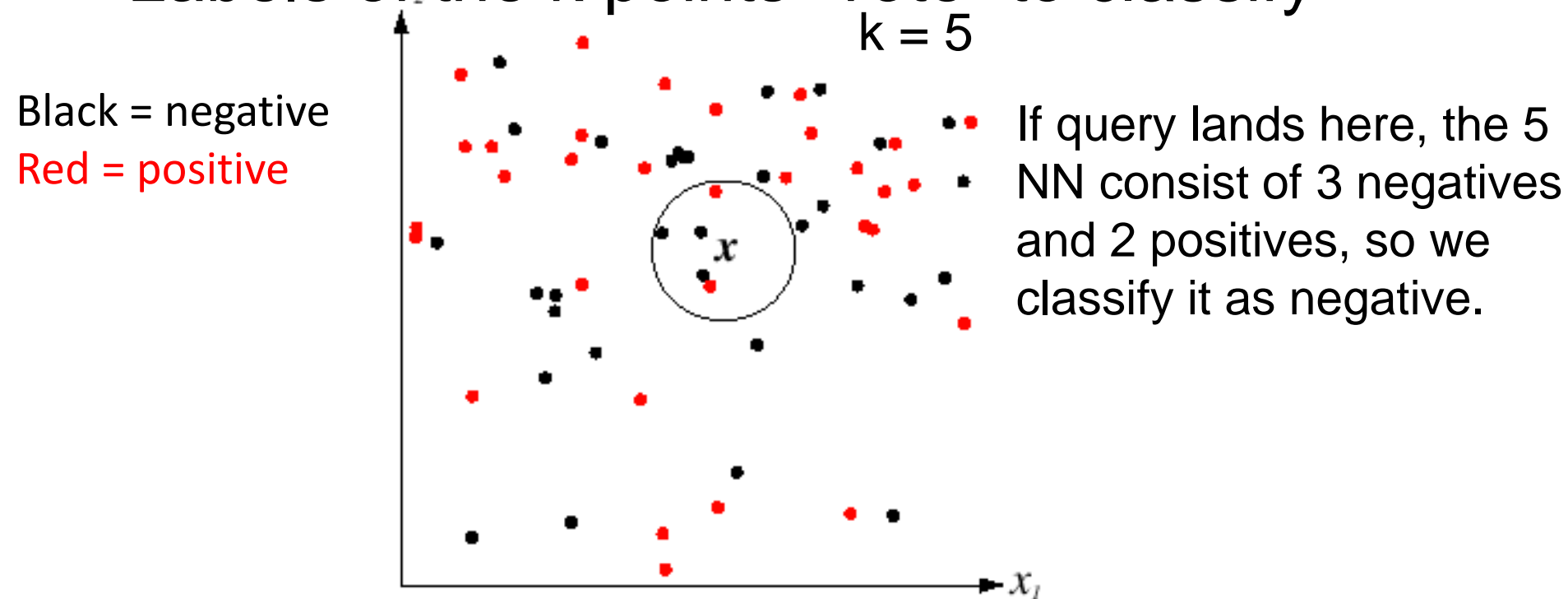
Closest to a
positive example
from the training
set, so classify it
as positive.

from Duda *et al.*

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

K-Nearest Neighbors classification

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



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

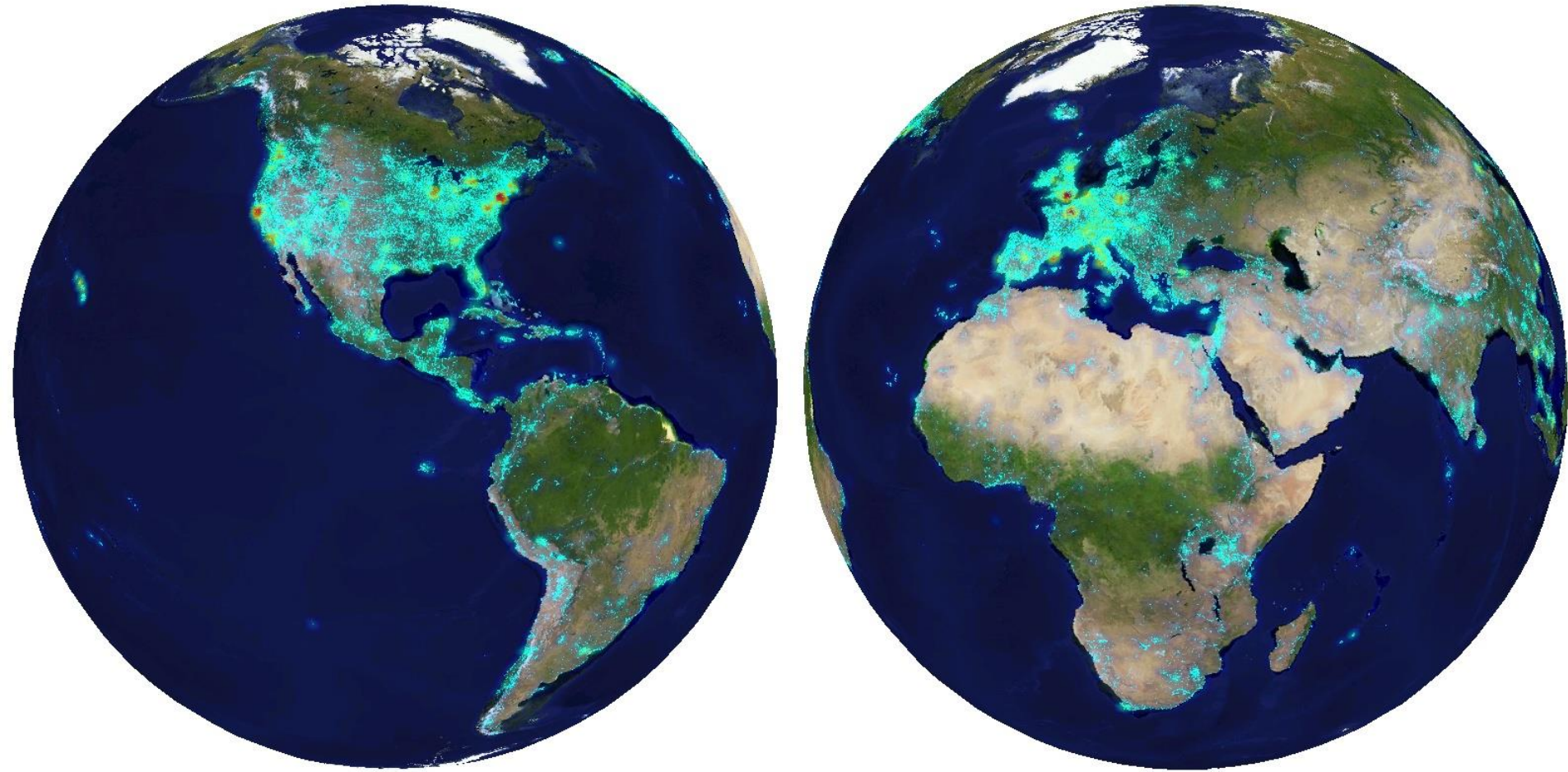
Where in the World?



Where in the World?



6+ million geotagged photos
by 109,788 photographers



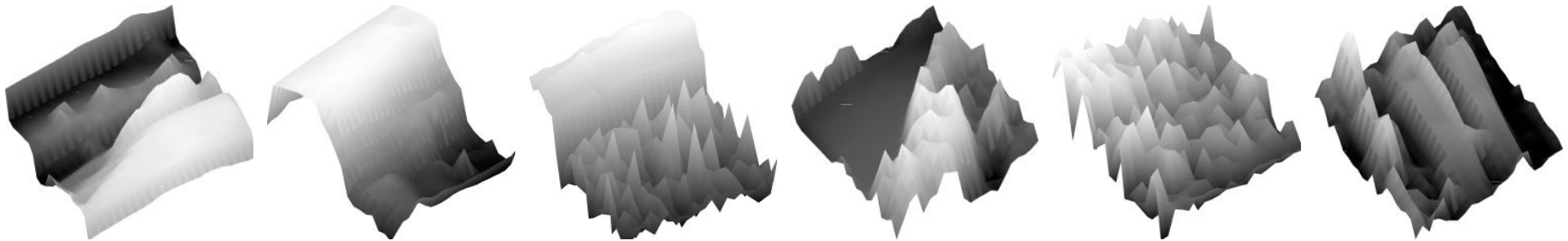
Annotated by Flickr users

Slides: James Hays

Which scene properties are relevant?

Spatial Envelope Theory of Scene Representation

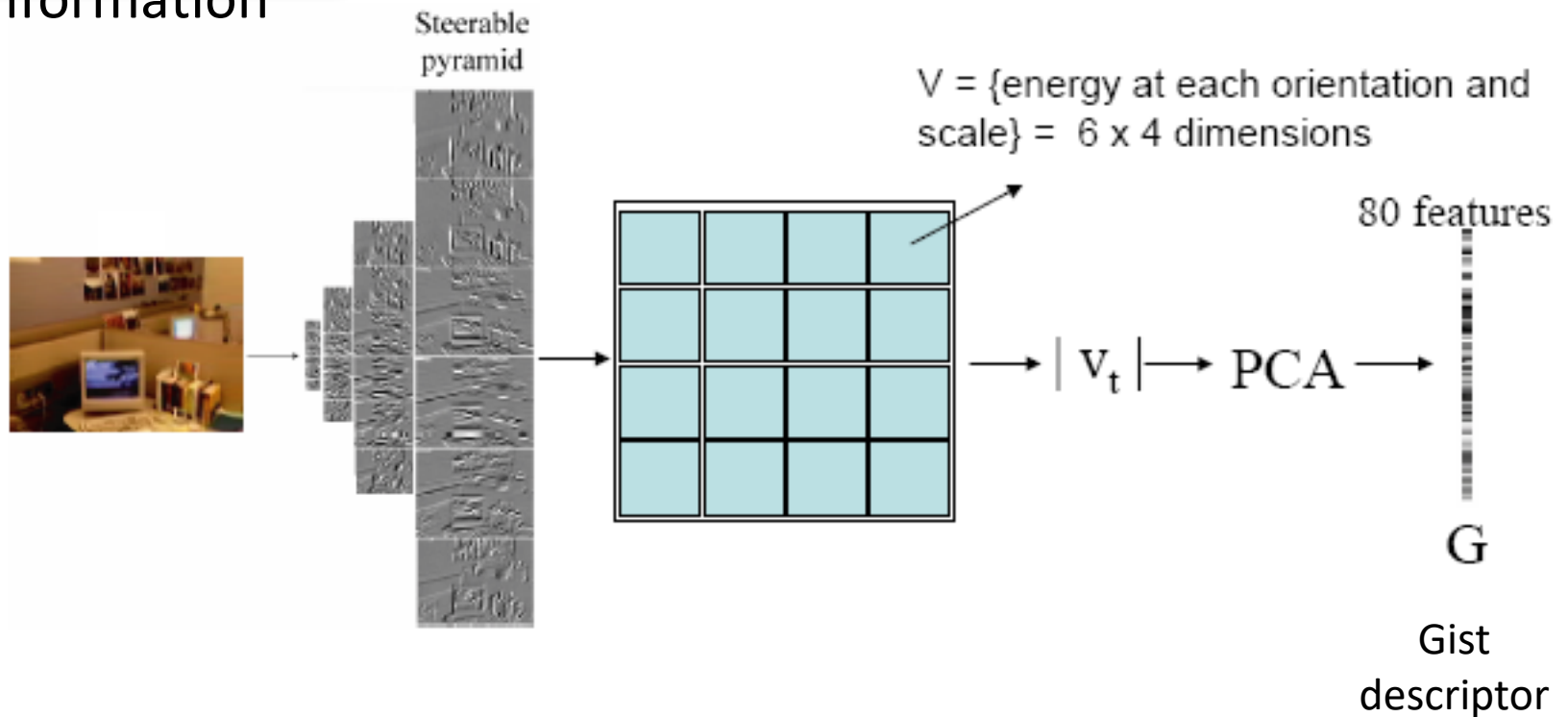
Oliva & Torralba (2001)



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

Global texture: capturing the “Gist” of the scene

Capture global image properties while keeping some spatial information



Which scene properties are relevant?

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

Scene Matches



Madrid



england



France



Paris



Croatia



heidelberg



Macau



Malta



Cairo



Italy



Italy



Italy



Latvia



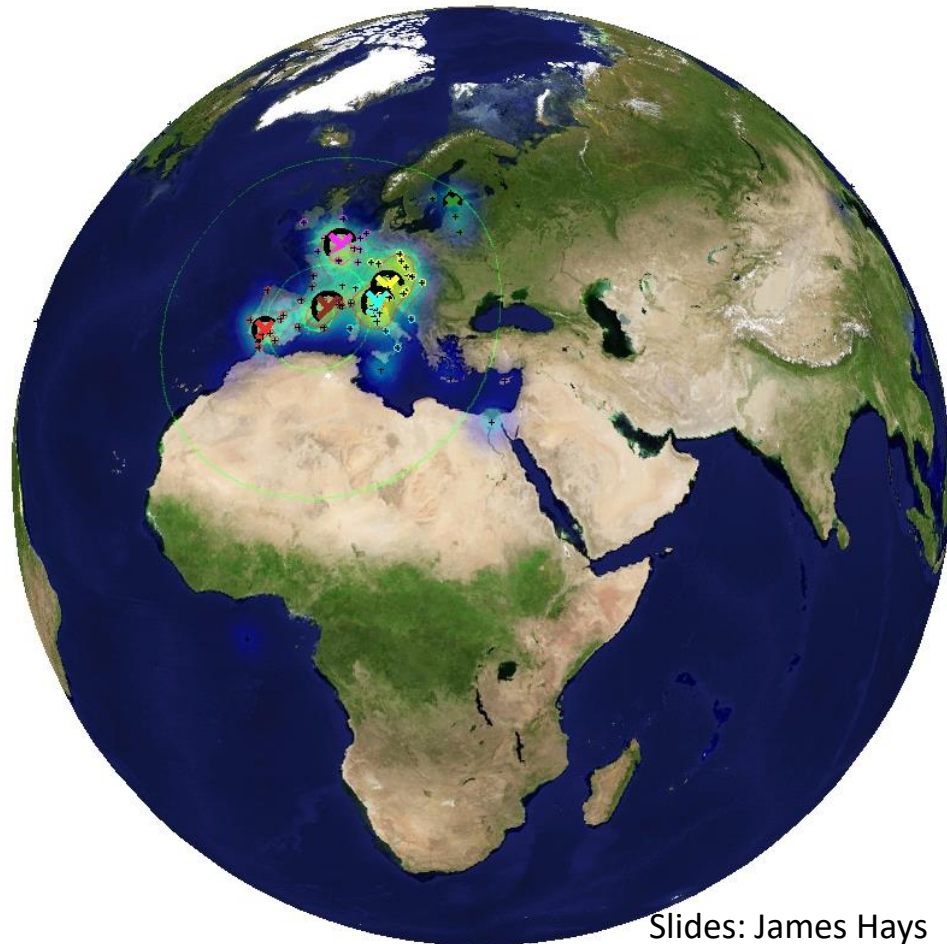
europe



Barcelona



Austria



Scene Matches



Paris



Paris



Paris



Paris



Paris



Paris



Paris



Madrid



Rome



Paris



Cuba



Paris



Paris



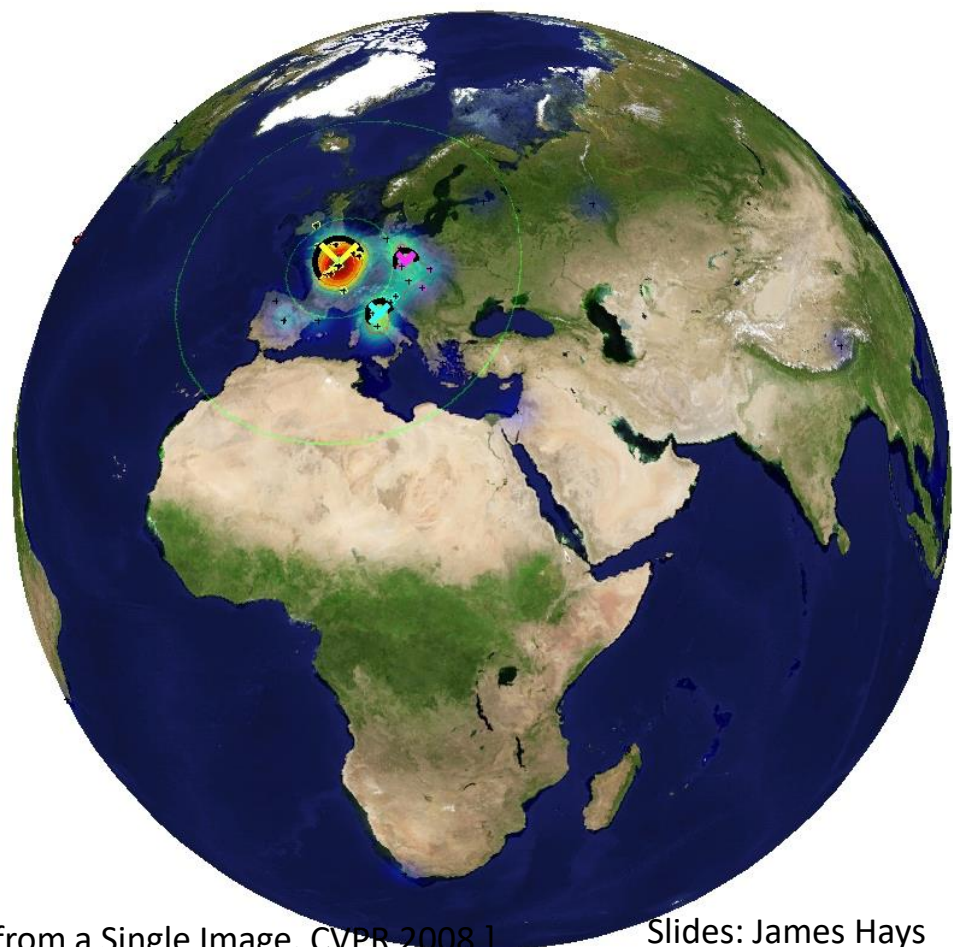
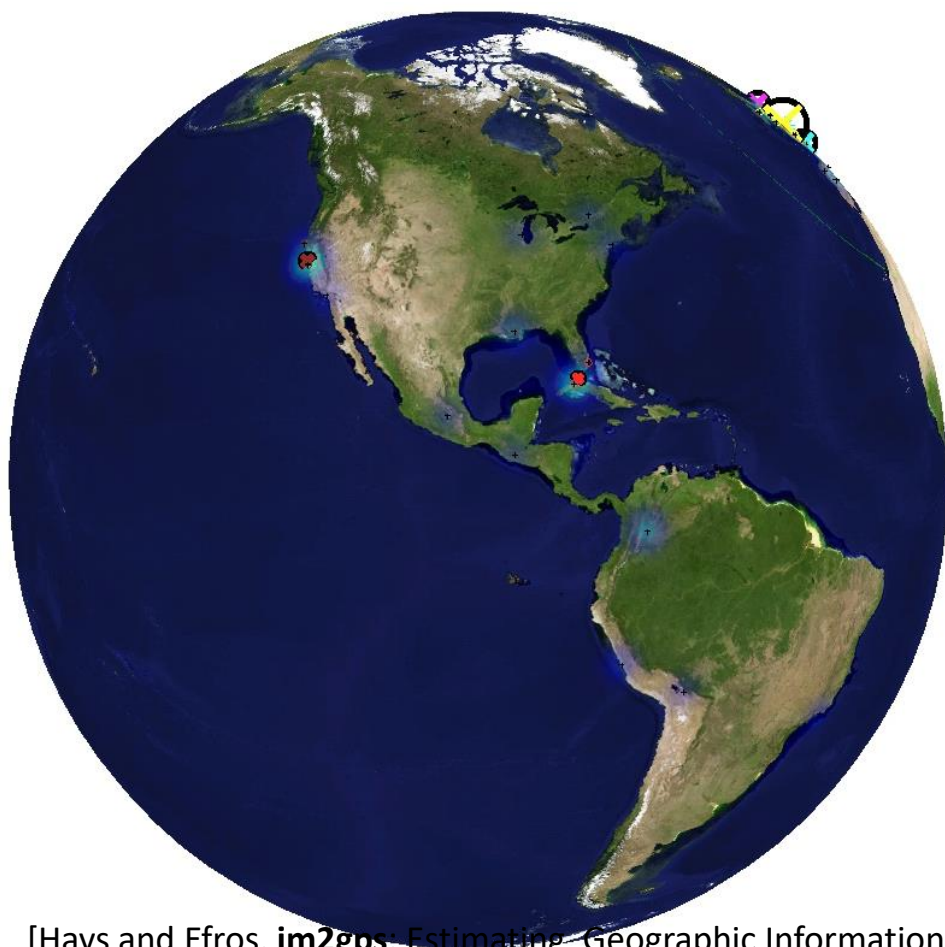
Poland



Paris



Paris



Scene Matches



Philippines



Houston



Thailand



Houston



Maldives



Philippines



NewZealand



Bermuda



Palau



Mexico2



Brazil



Mendoza



Brazil



Thailand



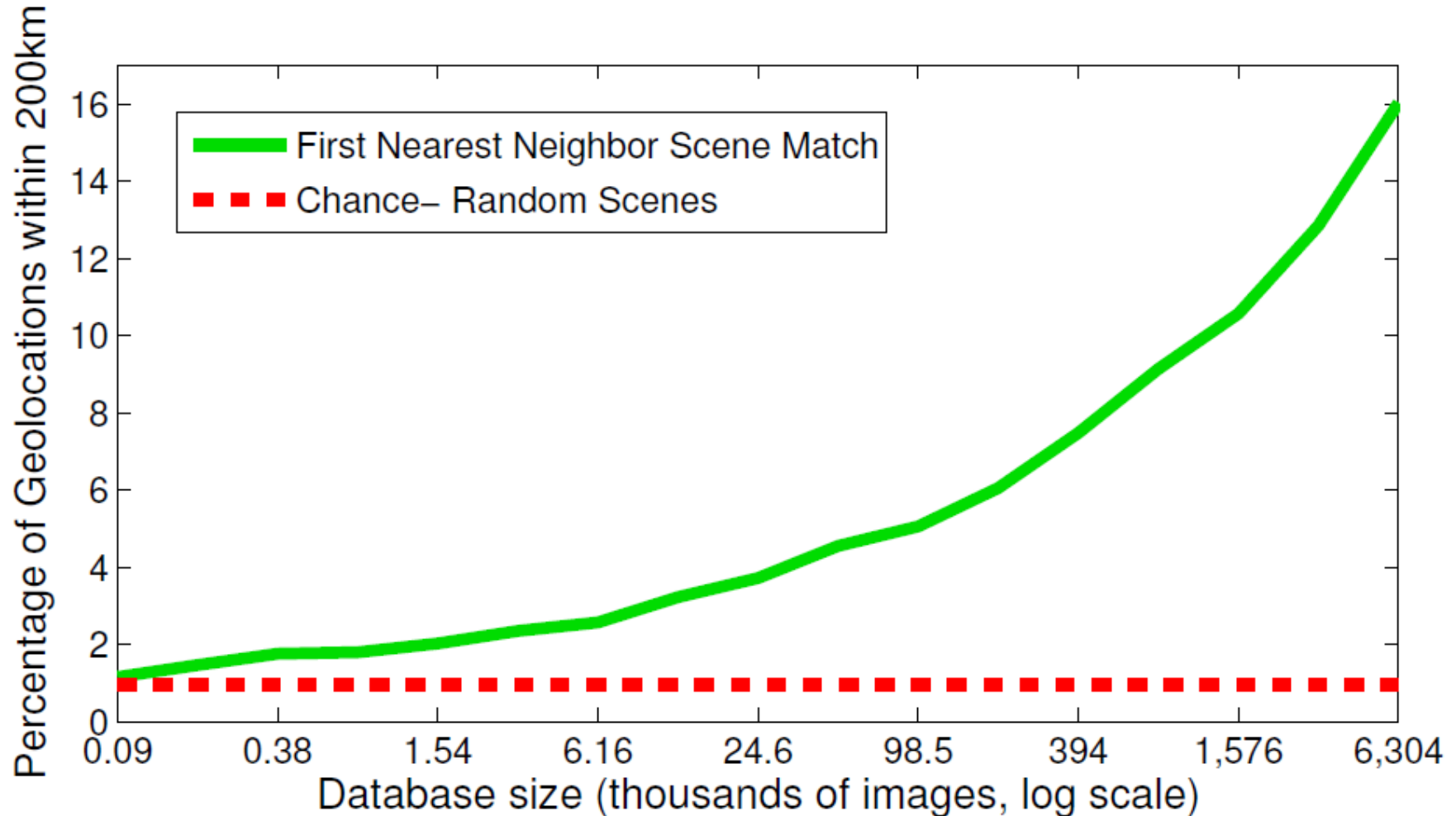
Arkansas



Hawaii



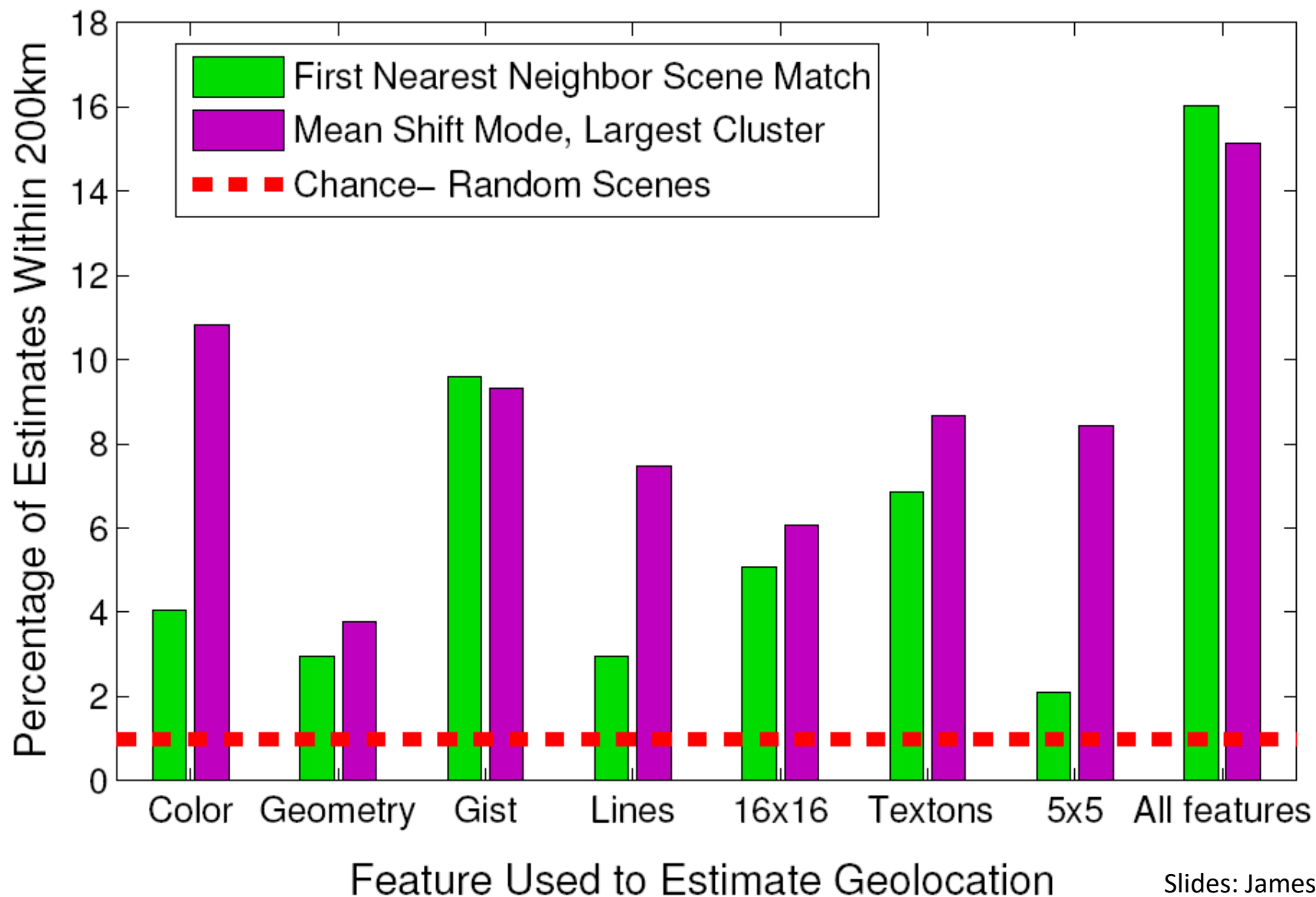
The Importance of Data



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

Slides: James Hays

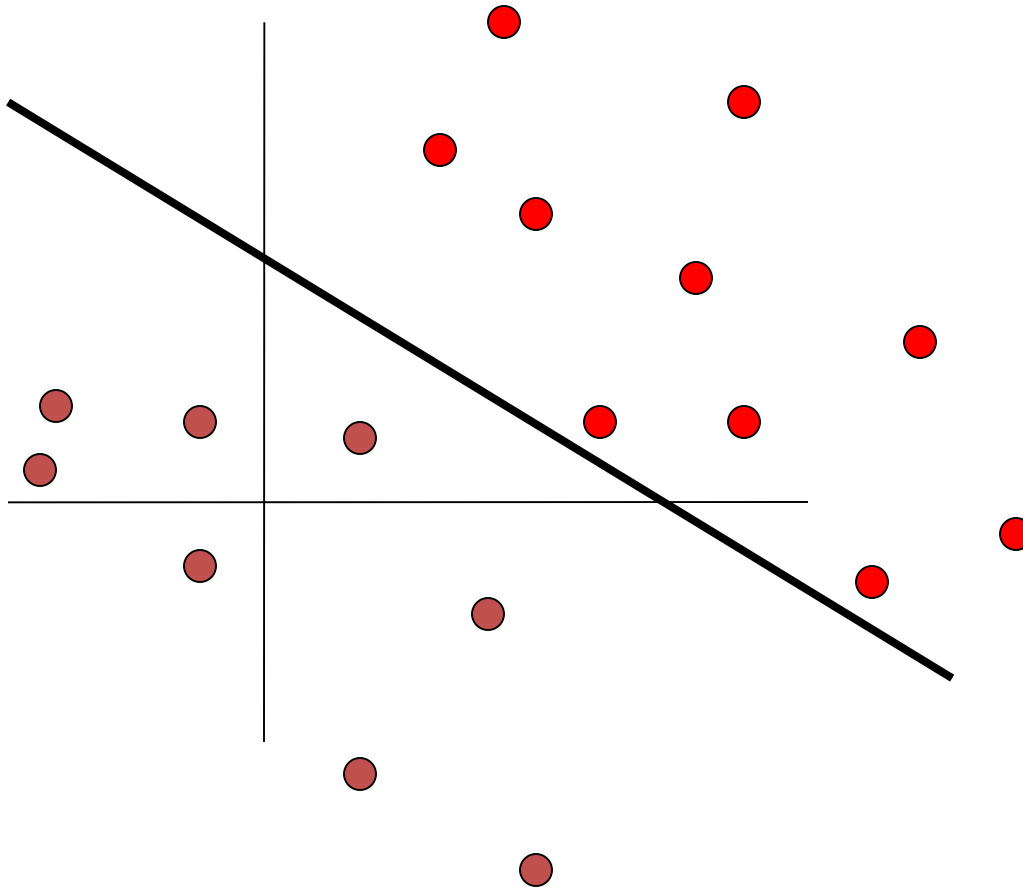
Feature Performance



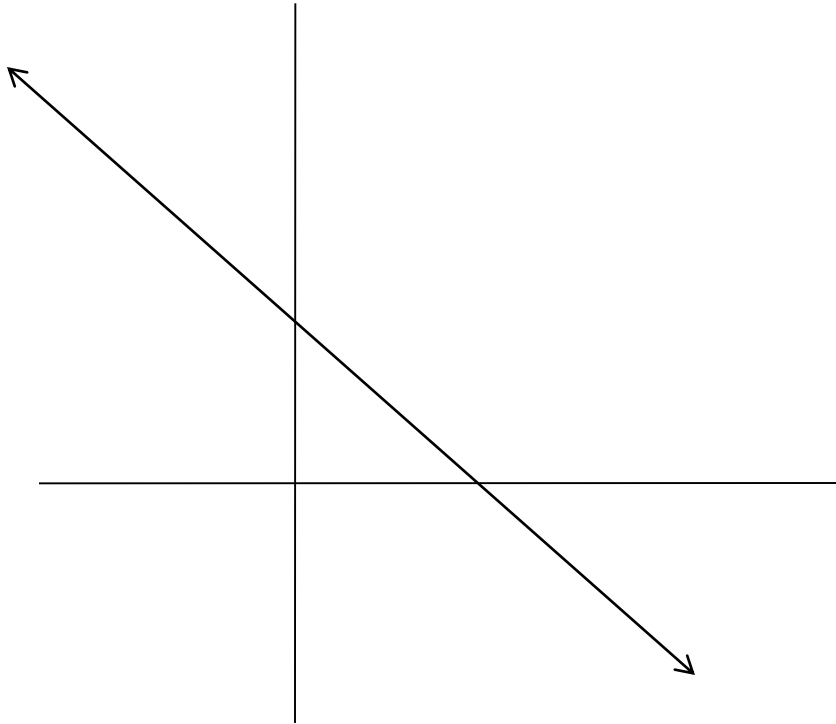
Nearest neighbors: pros and cons

- **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

Linear classifiers



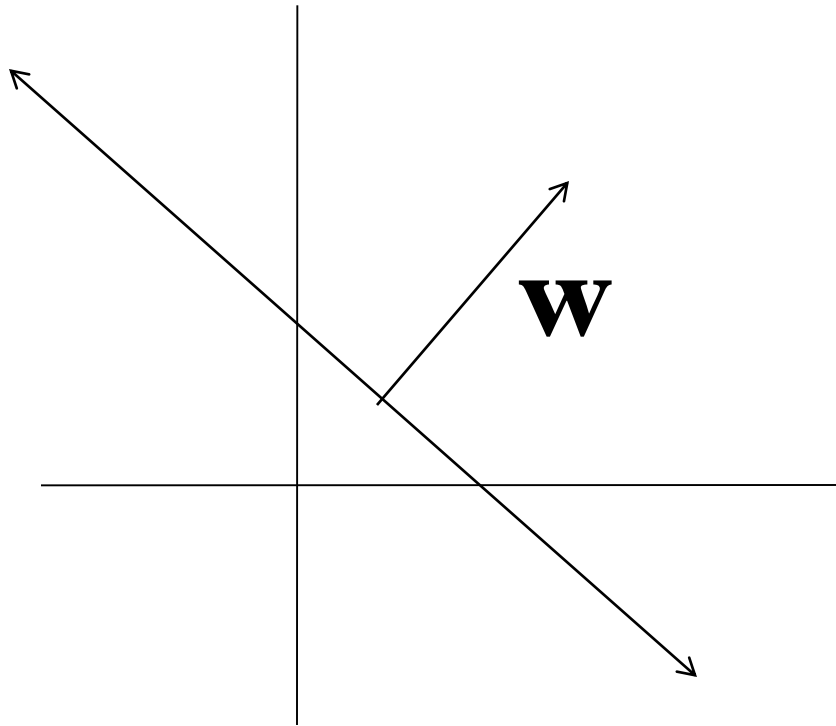
Lines in \mathbb{R}^2



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



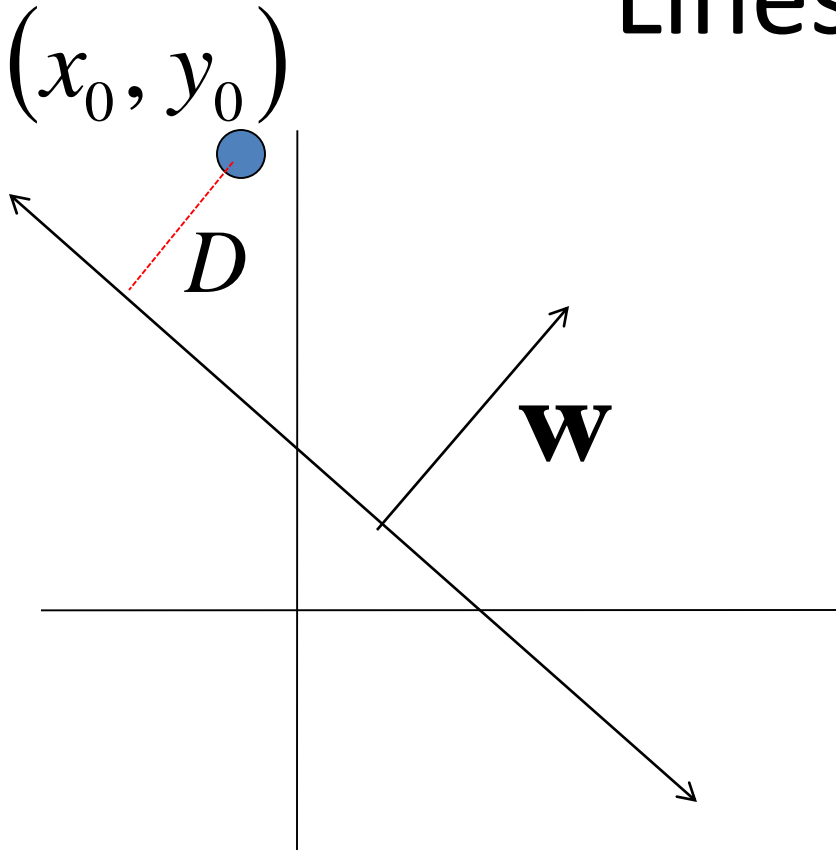
Let $\mathbf{w} = \begin{bmatrix} a \\ c \end{bmatrix}$ $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$

$$ax + cy + b = 0$$



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

Lines in \mathbb{R}^2



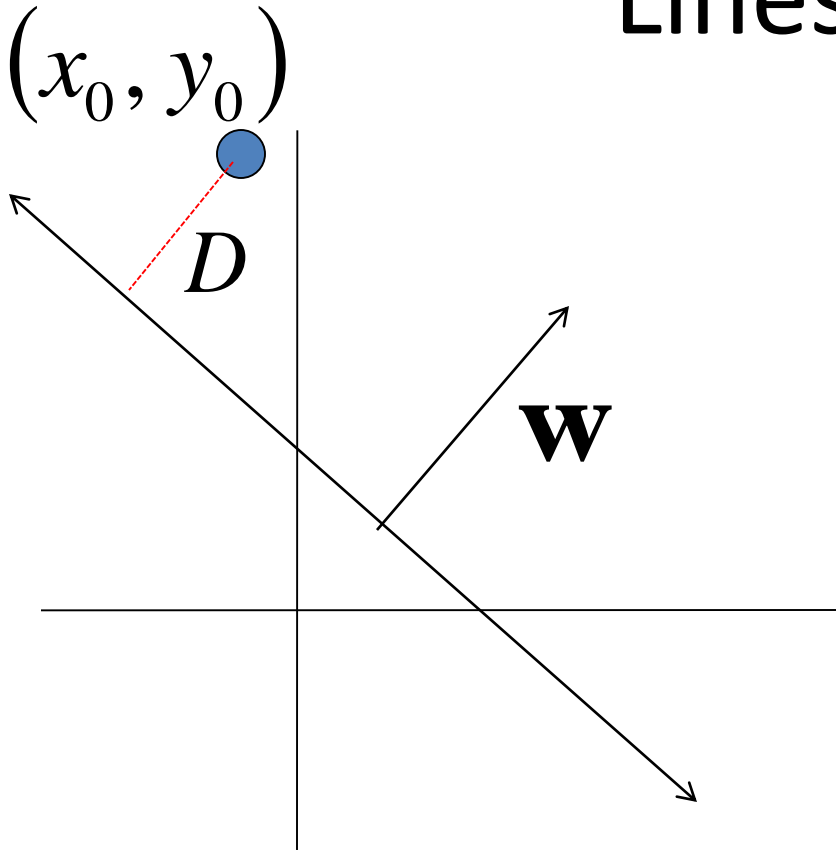
Let $\mathbf{w} = \begin{bmatrix} a \\ c \end{bmatrix}$ $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$

$$ax + cy + b = 0$$



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

Lines in \mathbb{R}^2



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

$$ax + cy + b = 0$$

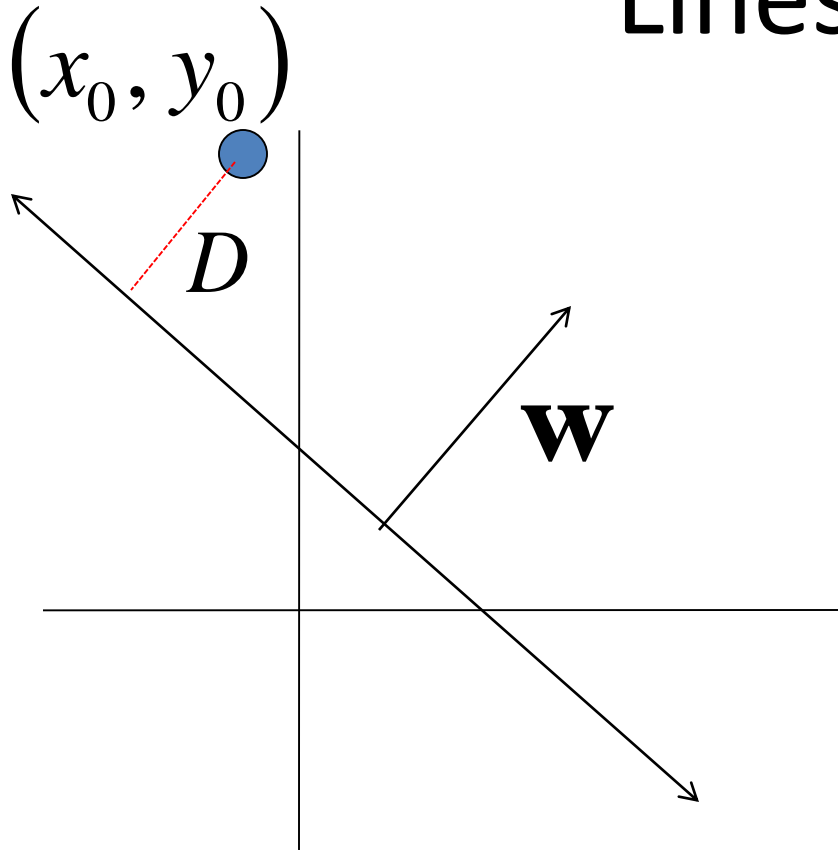


$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

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

distance from
point to line

Lines in \mathbb{R}^2



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

$$ax + cy + b = 0$$

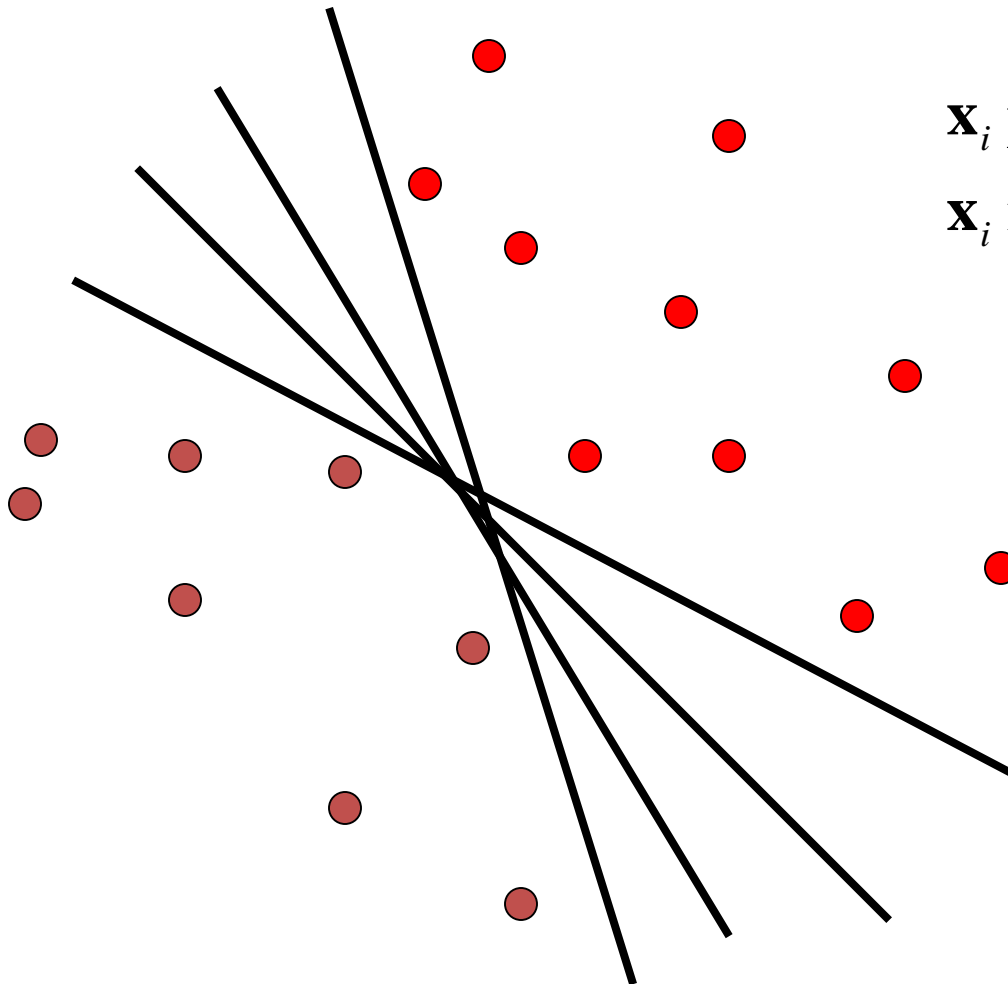


$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

$$D = \frac{|ax_0 + cy_0 + b|}{\sqrt{a^2 + c^2}} = \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|} \quad \left. \vphantom{\frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}} \right\} \begin{array}{l} \text{distance from} \\ \text{point to line} \end{array}$$

Linear classifiers

- Find linear function to separate positive and negative examples

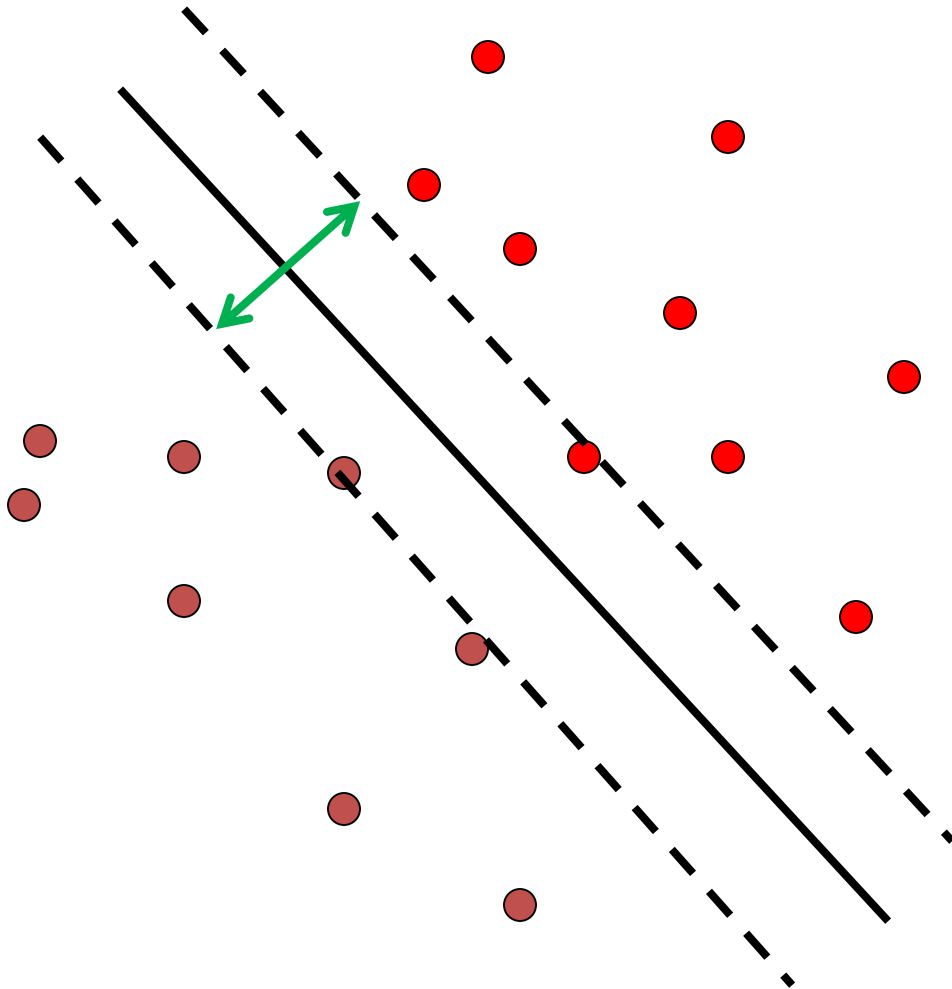


$$\mathbf{x}_i \text{ positive : } \mathbf{x}_i \cdot \mathbf{w} + b \geq 0$$

$$\mathbf{x}_i \text{ negative : } \mathbf{x}_i \cdot \mathbf{w} + b < 0$$

Which line
is best?

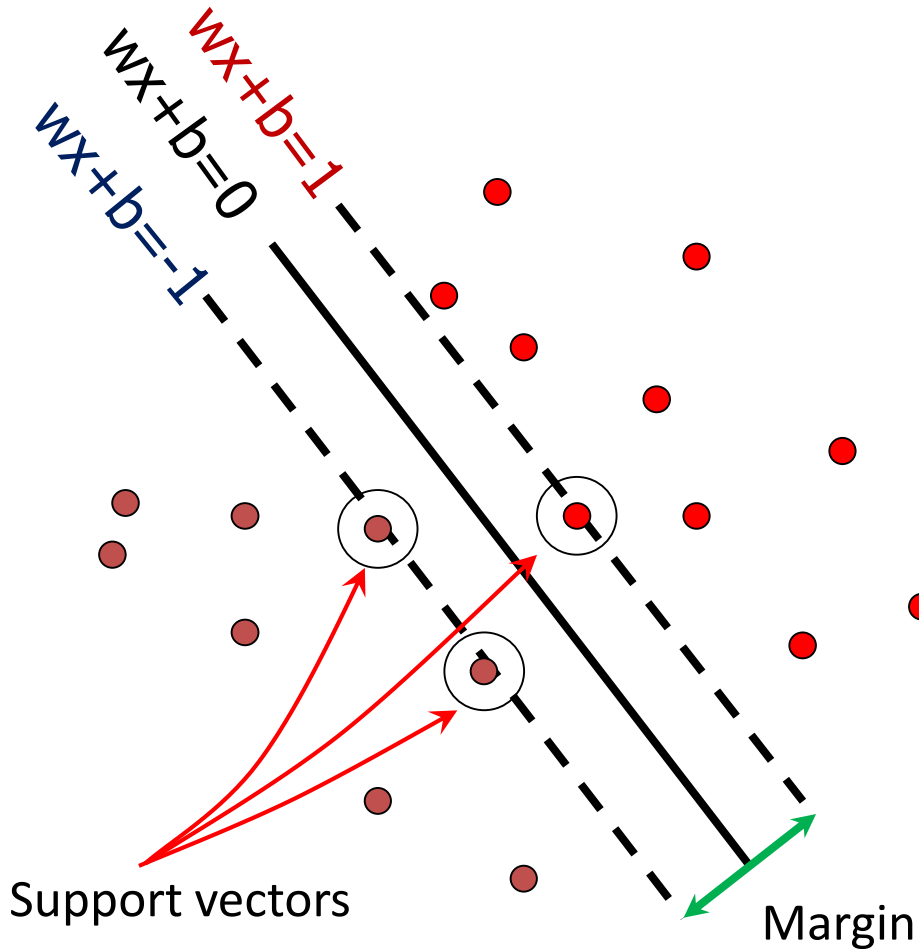
Support Vector Machines (SVMs)



- 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.



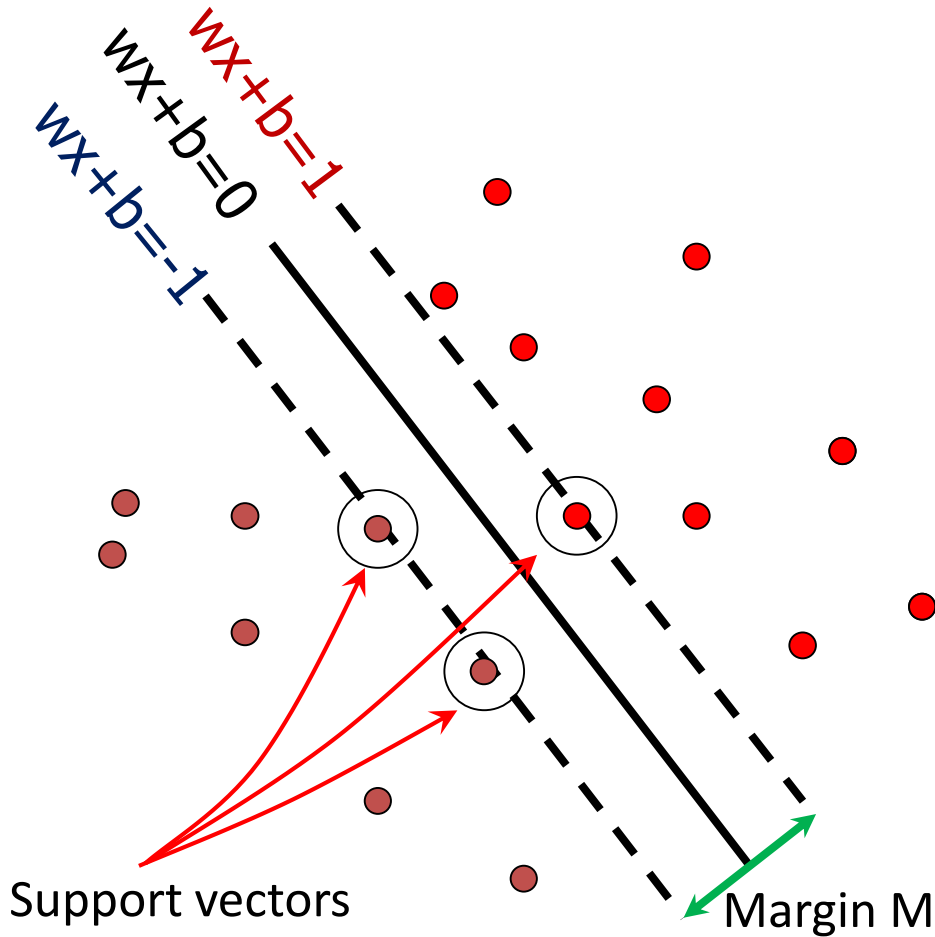
$$\mathbf{x}_i \text{ positive } (y_i = 1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

$$\text{For support, vectors,} \quad \mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$$

Support vector machines

- Want line that maximizes the margin.



$$\mathbf{x}_i \text{ positive } (y_i = 1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

$$\text{For support, vectors,} \quad \mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$$

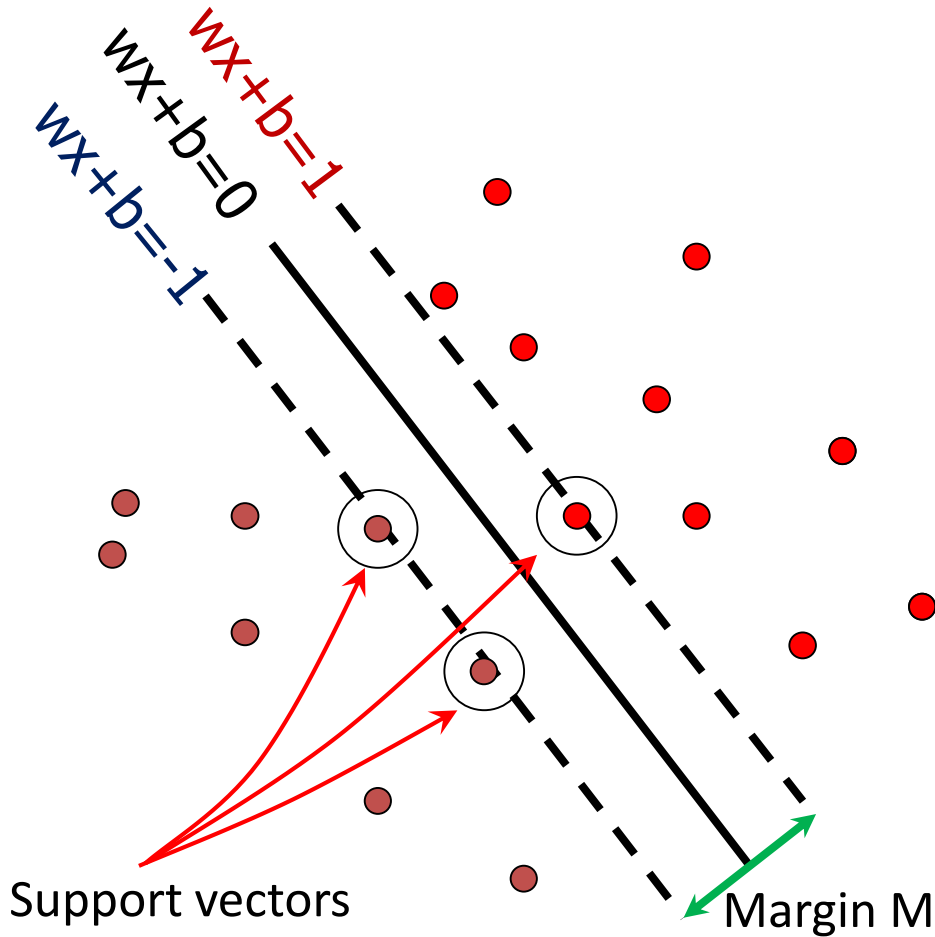
$$\text{Distance between point and line:} \quad \frac{|\mathbf{x}_i \cdot \mathbf{w} + b|}{\|\mathbf{w}\|}$$

For support vectors:

$$\frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|} = \frac{\pm 1}{\|\mathbf{w}\|} \quad M = \left| \frac{1}{\|\mathbf{w}\|} - \frac{-1}{\|\mathbf{w}\|} \right| = \frac{2}{\|\mathbf{w}\|}$$

Support vector machines

- Want line that maximizes the margin.



$$\mathbf{x}_i \text{ positive } (y_i = 1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

$$\text{For support, vectors,} \quad \mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$$

$$\text{Distance between point and line:} \quad \frac{|\mathbf{x}_i \cdot \mathbf{w} + b|}{\|\mathbf{w}\|}$$

$$\text{Therefore, the margin is } 2 / \|\mathbf{w}\|$$

Finding the maximum margin line

1. Maximize margin $2/\|\mathbf{w}\|$
2. Correctly classify all training data points:

$$\mathbf{x}_i \text{ positive } (y_i = 1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

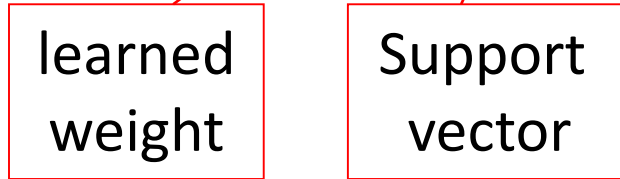
$$\mathbf{x}_i \text{ negative } (y_i = -1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

- *Quadratic optimization problem:*

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

Finding the maximum margin line

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




Finding the maximum margin line

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

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

- Classification function:


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

*If $f(x) < 0$, classify as negative,
if $f(x) > 0$, classify as positive*

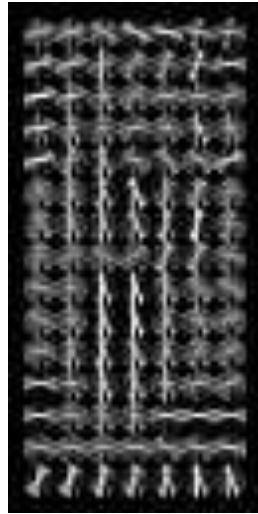
Questions

- **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.

Person detection with HoG's & linear SVM's

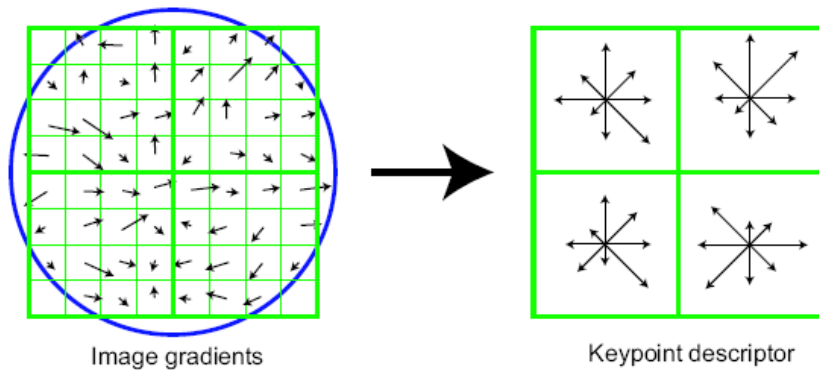


- Histograms of Oriented Gradients for Human Detection, [Navneet Dalal](#), [Bill Triggs](#), International Conference on Computer Vision & Pattern Recognition - June 2005
- <http://lear.inrialpes.fr/pubs/2005/DT05/>

Histograms of oriented gradients

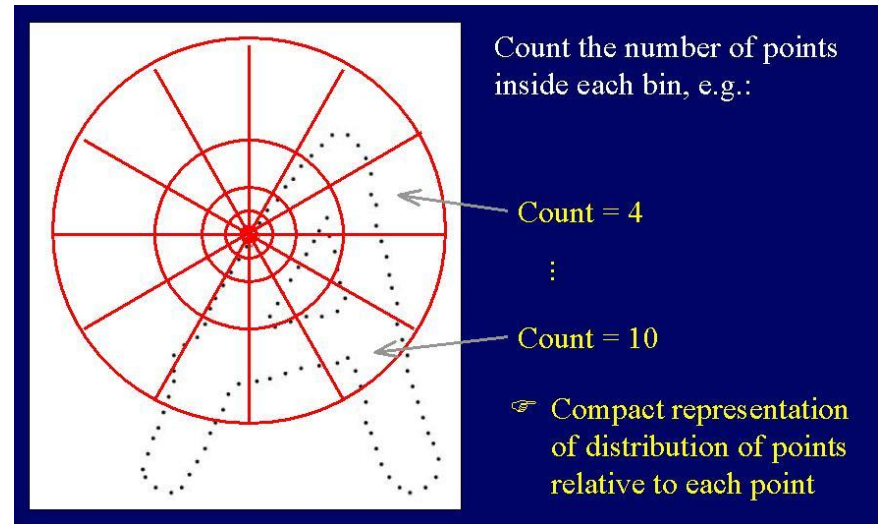
Histograms of oriented gradients

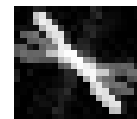
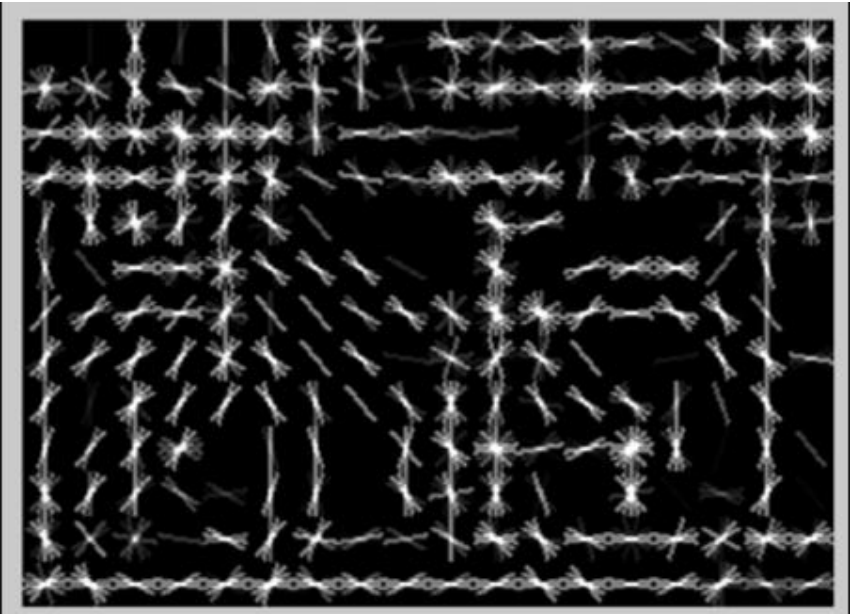
SIFT, D. Lowe, ICCV 1999



Shape context

Belongie, Malik, Puzicha, NIPS 2000





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>

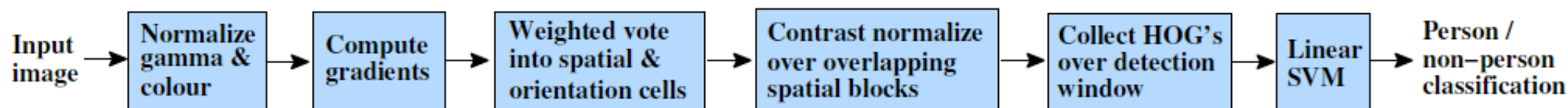


Figure 1. An overview of our feature extraction and object detection chain. The detector window is tiled with a grid of overlapping blocks in which Histogram of Oriented Gradient feature vectors are extracted. The combined vectors are fed to a linear SVM for object/non-object classification. The detection window is scanned across the image at all positions and scales, and conventional non-maximum suppression is run on the output pyramid to detect object instances, but this paper concentrates on the feature extraction process.

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>

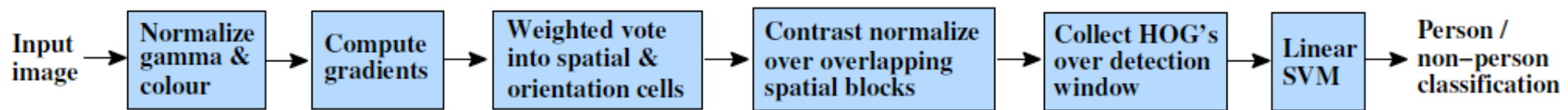
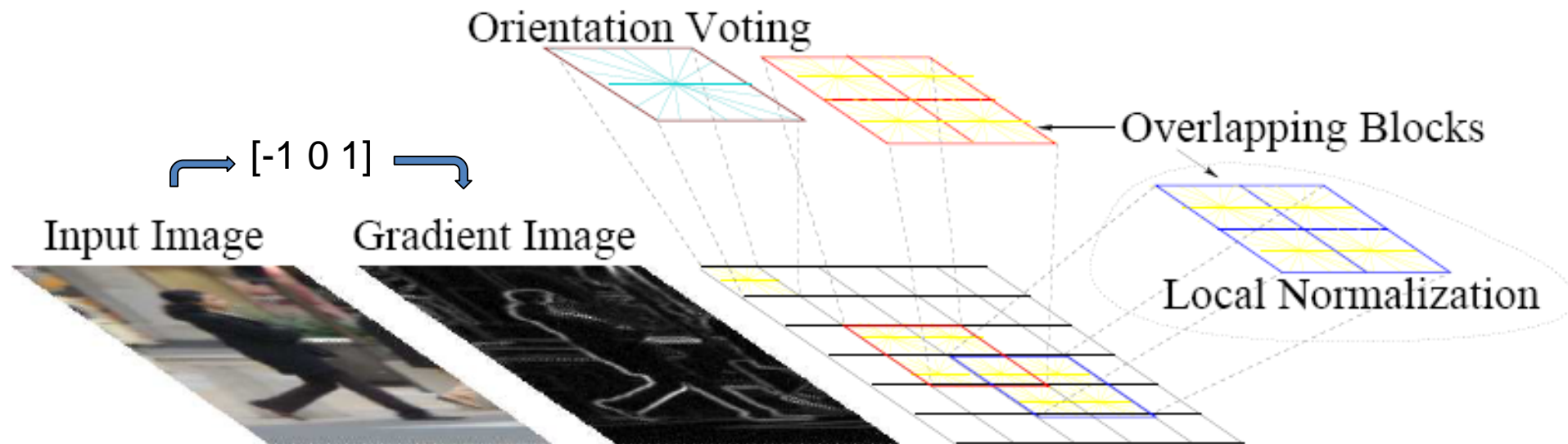


Figure 1. An overview of our feature extraction and object detection chain. The detector window is tiled with a grid of overlapping blocks in which Histogram of Oriented Gradient feature vectors are extracted. The combined vectors are fed to a linear SVM for object/non-object classification. The detection window is scanned across the image at all positions and scales, and conventional non-maximum suppression is run on the output pyramid to detect object instances, but this paper concentrates on the feature extraction process.



SVM

A Support Vector Machine (SVM) learns a classifier with the form:

$$H(x) = \sum_{m=1}^M a_m y_m k(x, x_m)$$

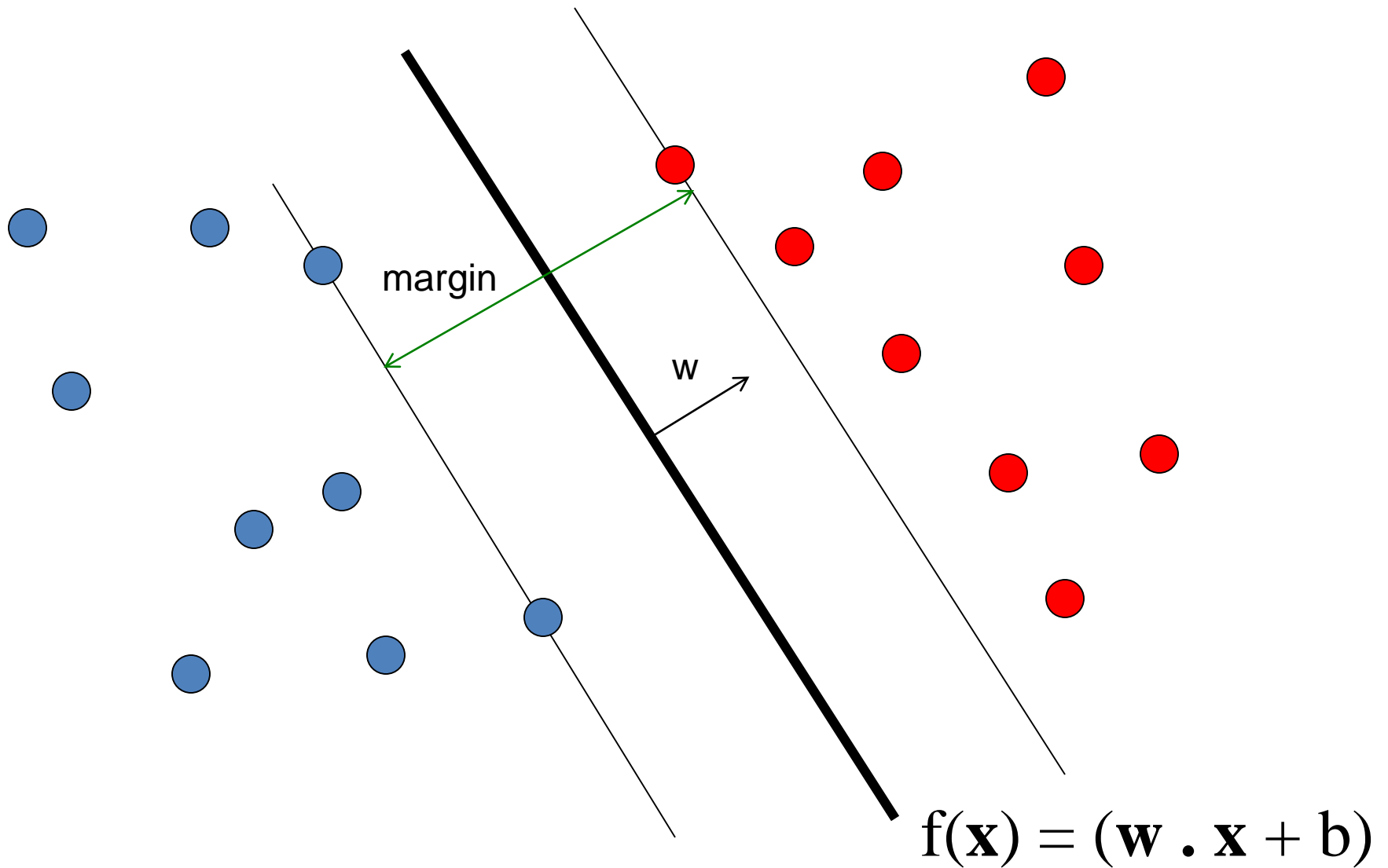
Where $\{x_m, y_m\}$, for $m = 1 \dots M$, are the training data with x_m being the input feature vector and $y_m = +1, -1$ the class label. $k(x, x_m)$ is the kernel and it can be any symmetric function satisfying the Mercer Theorem.

The classification is obtained by thresholding the value of $H(x)$.

There is a large number of possible kernels, each yielding a different family of decision boundaries:

- Linear kernel: $k(x, x_m) = x^T x_m$
- Radial basis function: $k(x, x_m) = \exp(-|x - x_m|^2/\sigma^2)$.
- Histogram intersection: $k(x, x_m) = \sum_i (\min(x(i), x_m(i)))$

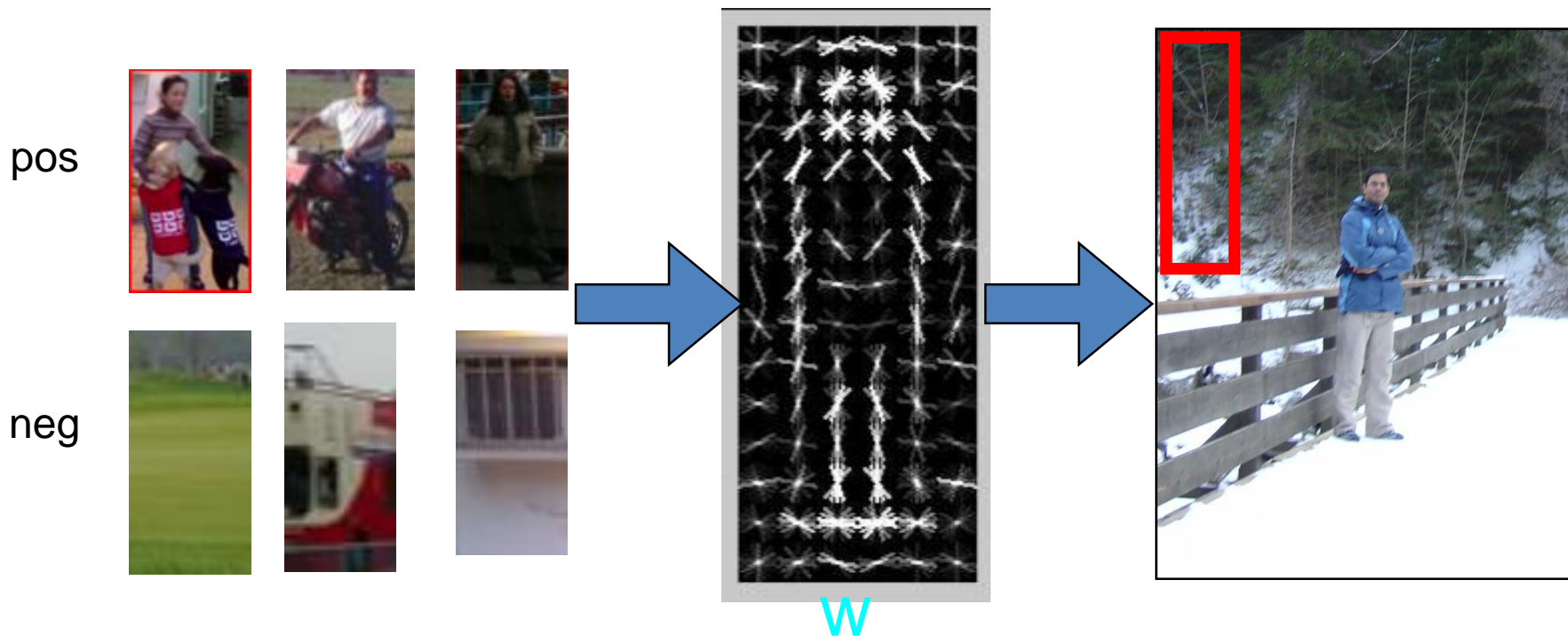
Linear SVM



Scanning-window templates

Dalal and Triggs CVPR05 (HOG)

Papageorgiou and Poggio ICIP99 (wavelets)



w = weights for orientation and spatial bins

$$w \cdot x > 0$$



Train with a linear classifier (perceptron, logistic regression, SVMs...)

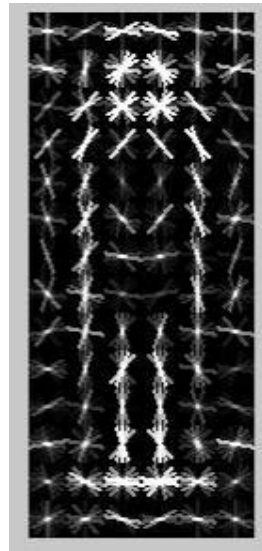
How to interpret positive and negative weights?

$$w \cdot x > 0$$

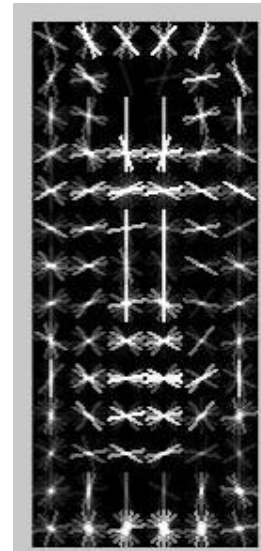
$$(w_{\text{pos}} - w_{\text{neg}}) \cdot x > 0$$

$$w_{\text{pos}} \cdot x > w_{\text{neg}} \cdot x$$

Pedestrian
template



>



Pedestrian
background
template

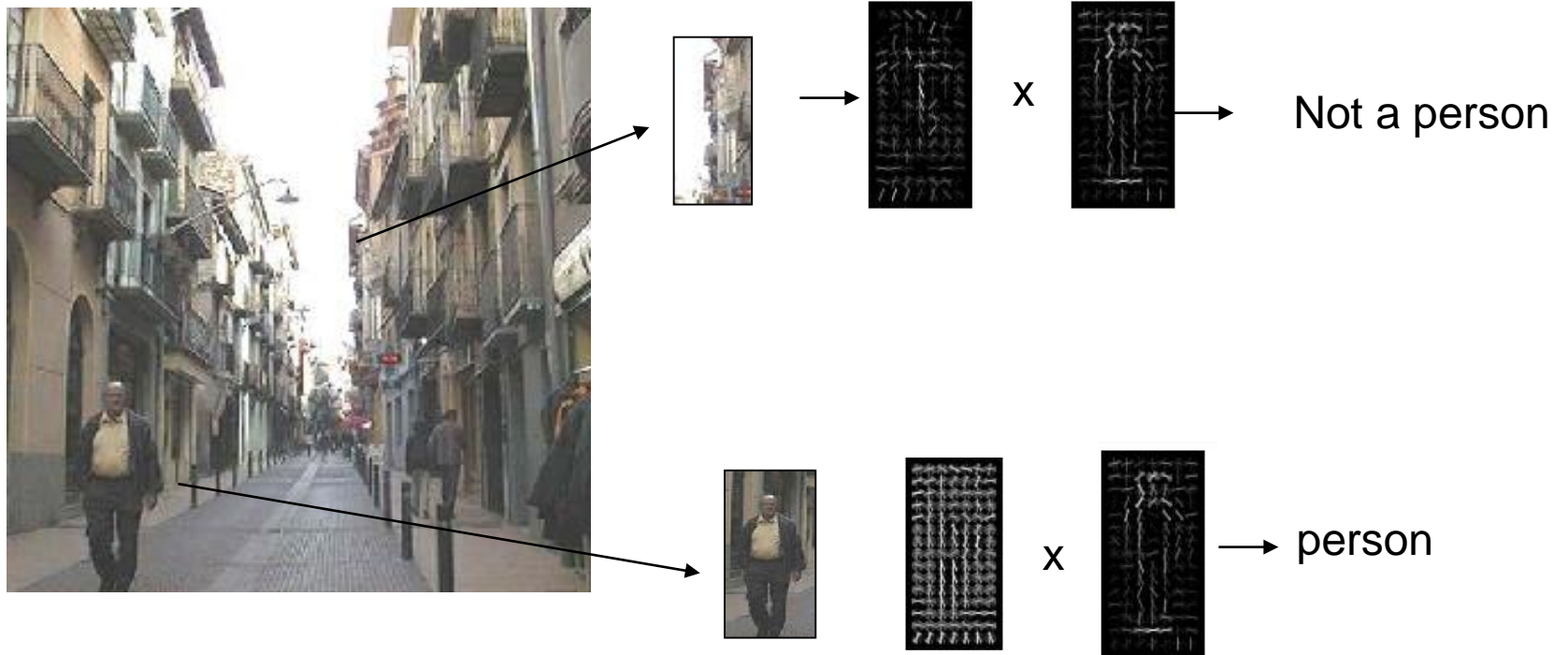
$w_{\text{pos}}, w_{\text{neg}}$ = weighted average of positive, negative support vectors

Right approach is to **compete** pedestrian, pillar, doorway... models

Background class is hard to model - easier to penalize particular vertical edges

Histograms of oriented gradients

Dalal & Trigs, 2006



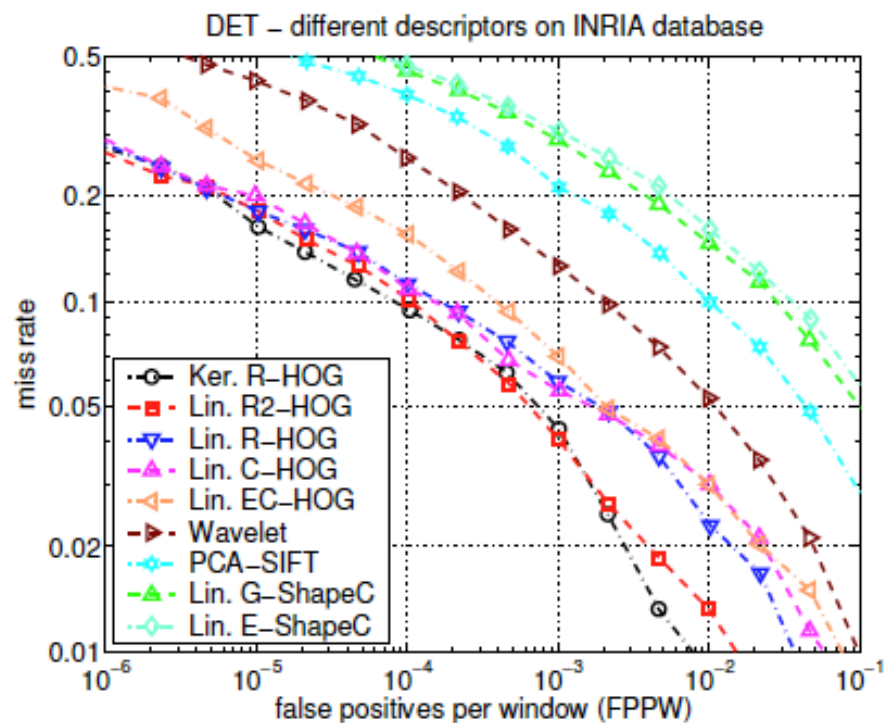
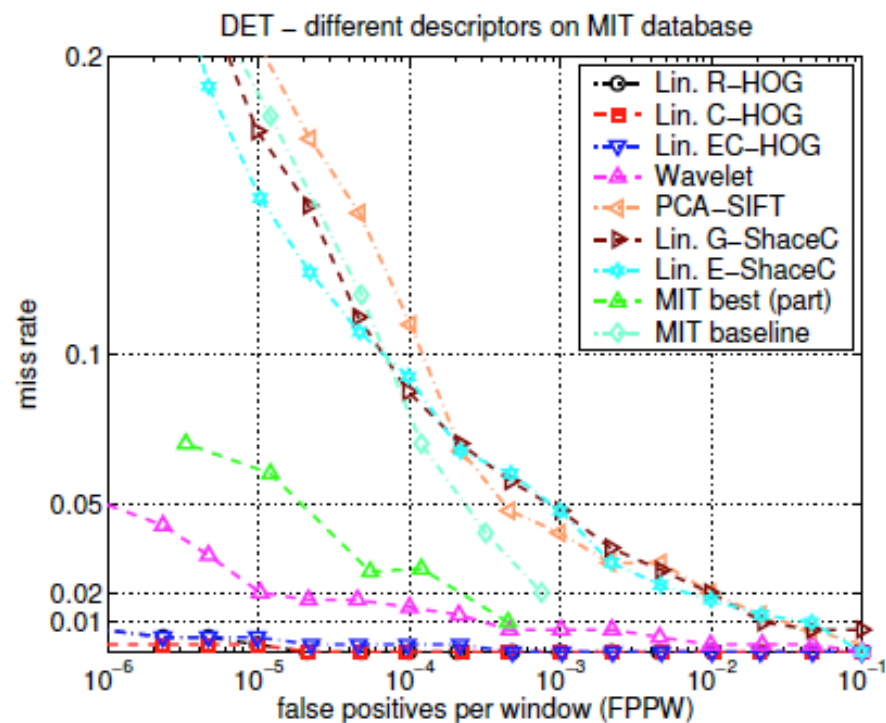


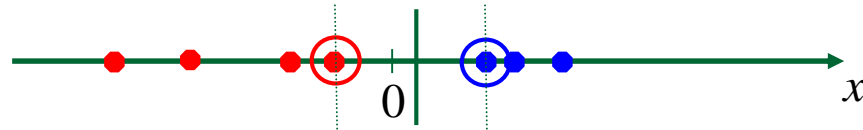
Figure 3. The performance of selected detectors on (left) MIT and (right) INRIA data sets. See the text for details.

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

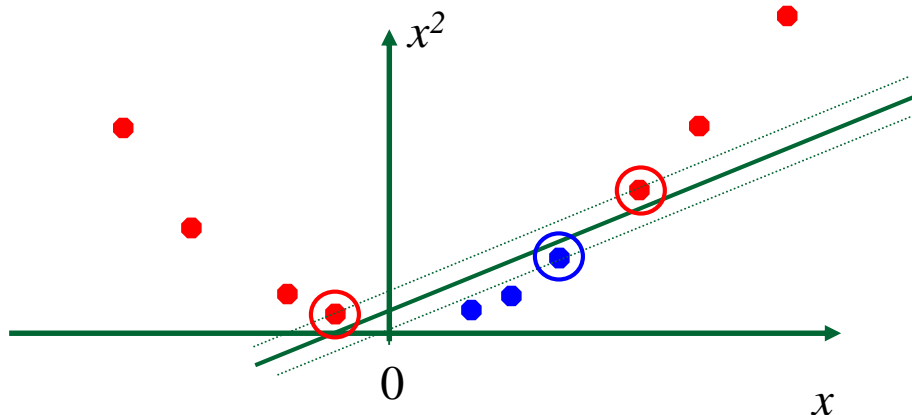
- Datasets that are linearly separable with some noise work out great:



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

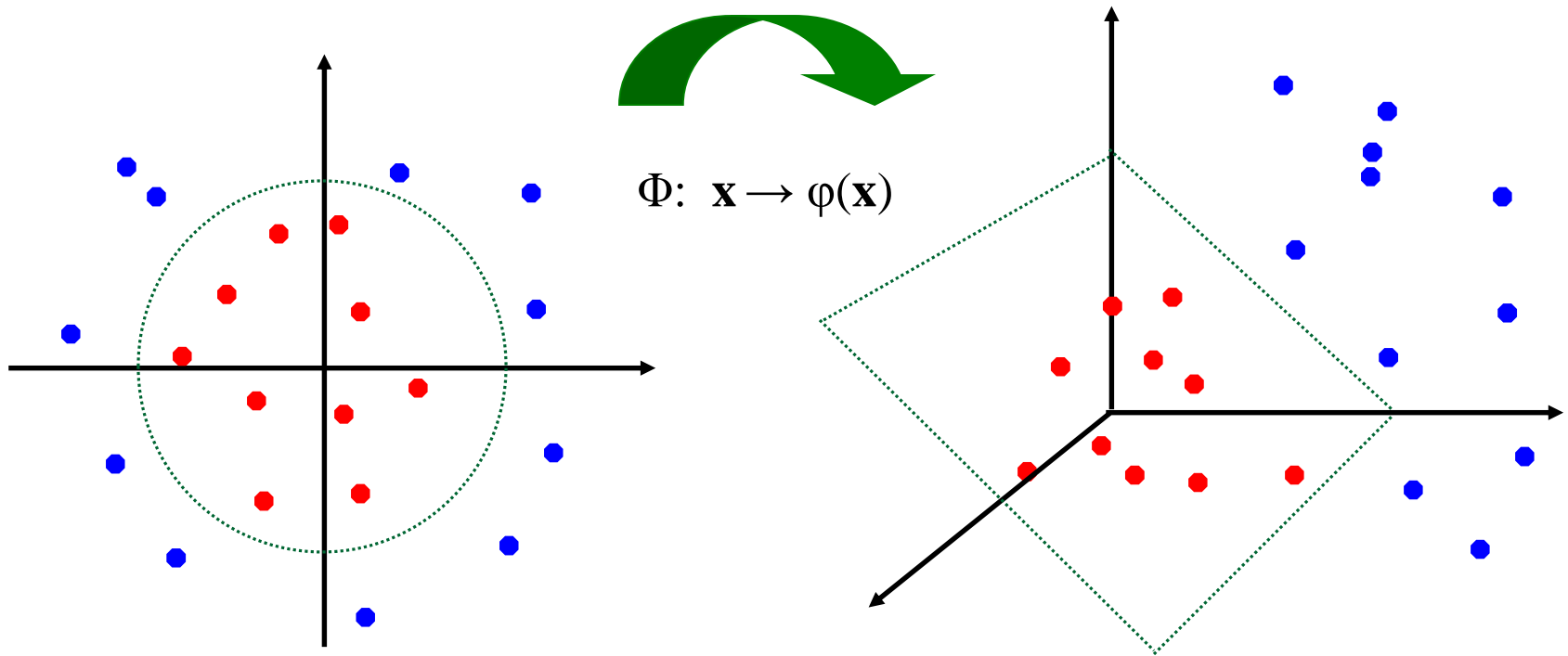


- How about... mapping data to a higher-dimensional space:



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:



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.

Example

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

let $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^2$

Need to show that $K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$:

$$\begin{aligned} K(\mathbf{x}_i, \mathbf{x}_j) &= (1 + \mathbf{x}_i^T \mathbf{x}_j)^2, \\ &= 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{j1} + 2 x_{i2} x_{j2} \\ &= [1 \ x_{i1}^2 \ \sqrt{2} x_{i1} x_{i2} \ x_{i2}^2 \ \sqrt{2} x_{i1} \ \sqrt{2} x_{i2}]^T \\ &\quad [1 \ x_{j1}^2 \ \sqrt{2} x_{j1} x_{j2} \ x_{j2}^2 \ \sqrt{2} x_{j1} \ \sqrt{2} x_{j2}] \\ &= \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j), \end{aligned}$$

where $\varphi(\mathbf{x}) = [1 \ x_1^2 \ \sqrt{2} x_1 x_2 \ x_2^2 \ \sqrt{2} x_1 \ \sqrt{2} x_2]$

Nonlinear SVMs

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

$$K(\mathbf{x}_i, \mathbf{x}_j) = \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}_j)$$

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

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

Examples of kernel functions

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

- Linear:

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

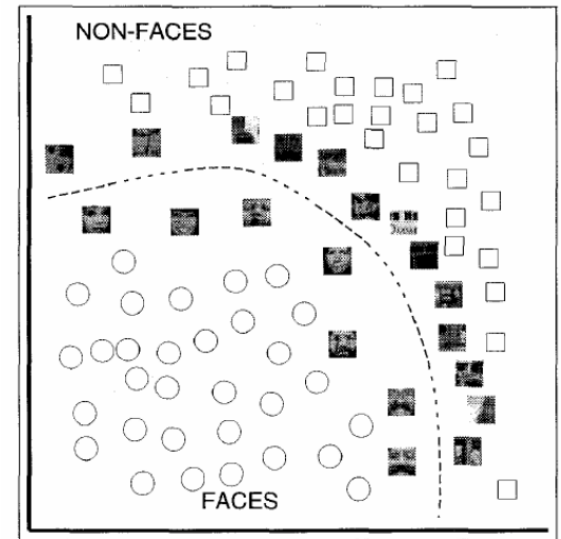
- Gaussian RBF:

$$K(x_i, x_j) = \sum_k \min(x_i(k), x_j(k))$$

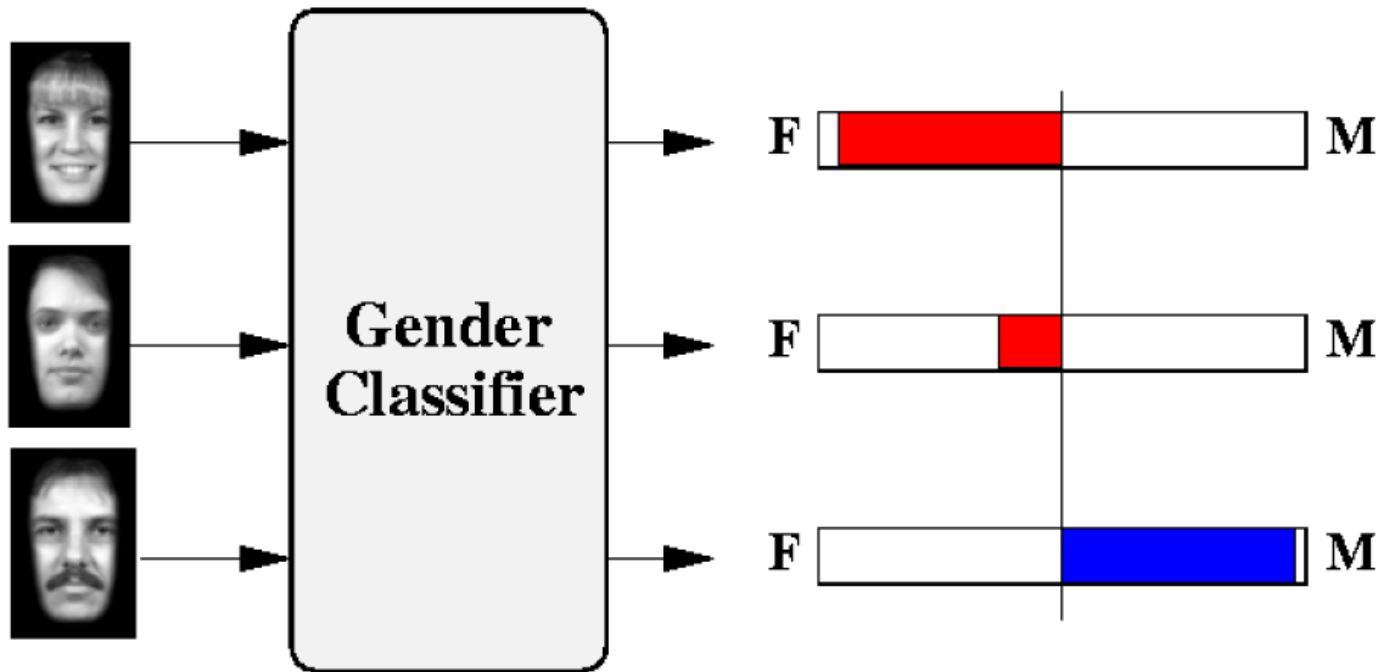
- Histogram intersection:

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.



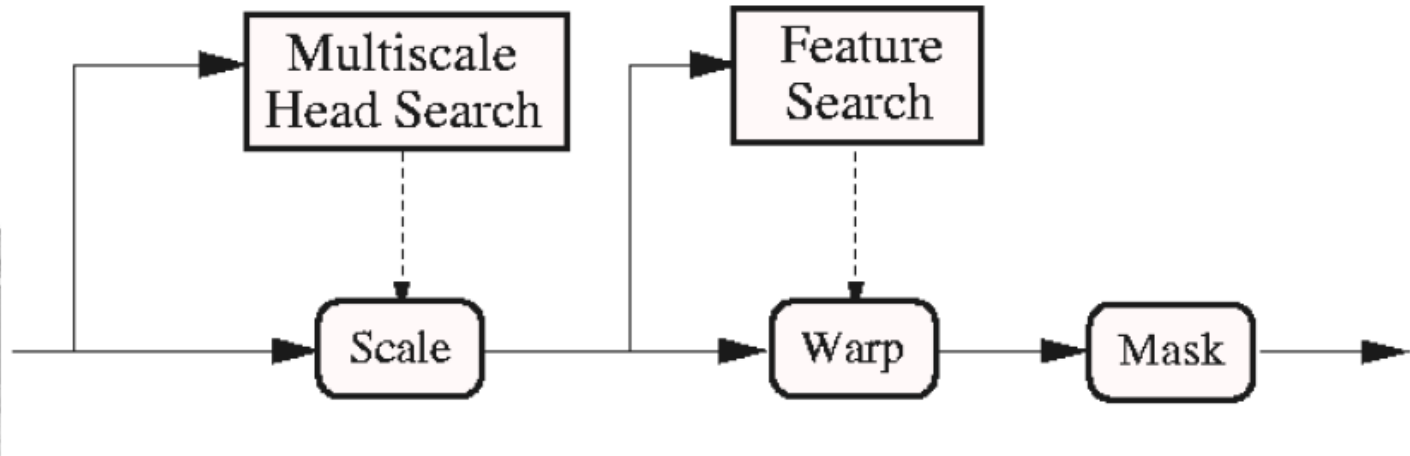
Example: learning gender with SVMs



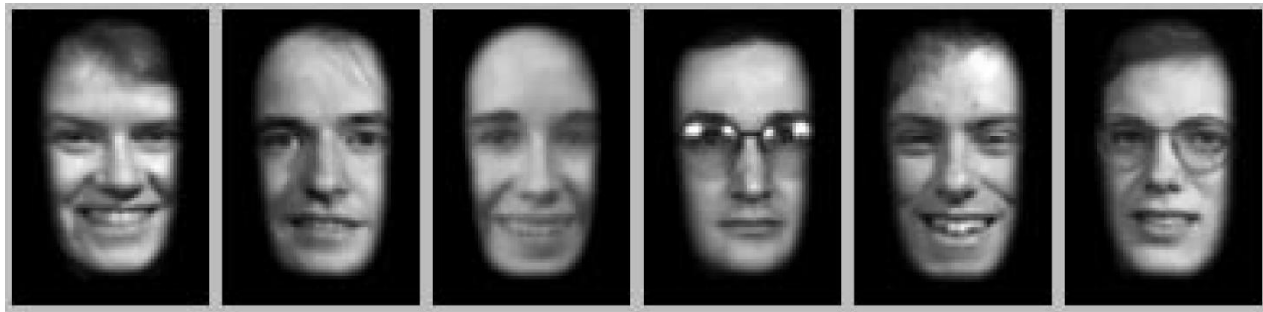
Moghaddam and Yang, Learning Gender with Support Faces, TPAMI 2002.

Moghaddam and Yang, Face & Gesture 2000.

Face alignment
processing



Processed faces

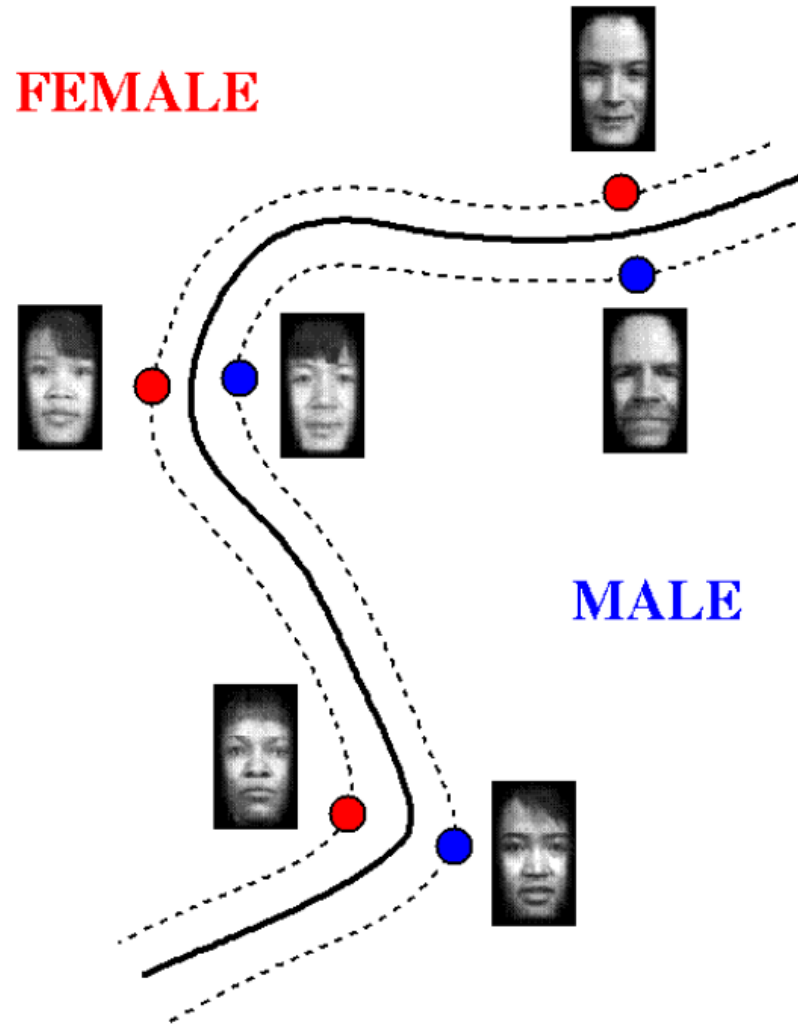


Learning gender with SVMs

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

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

Support Faces



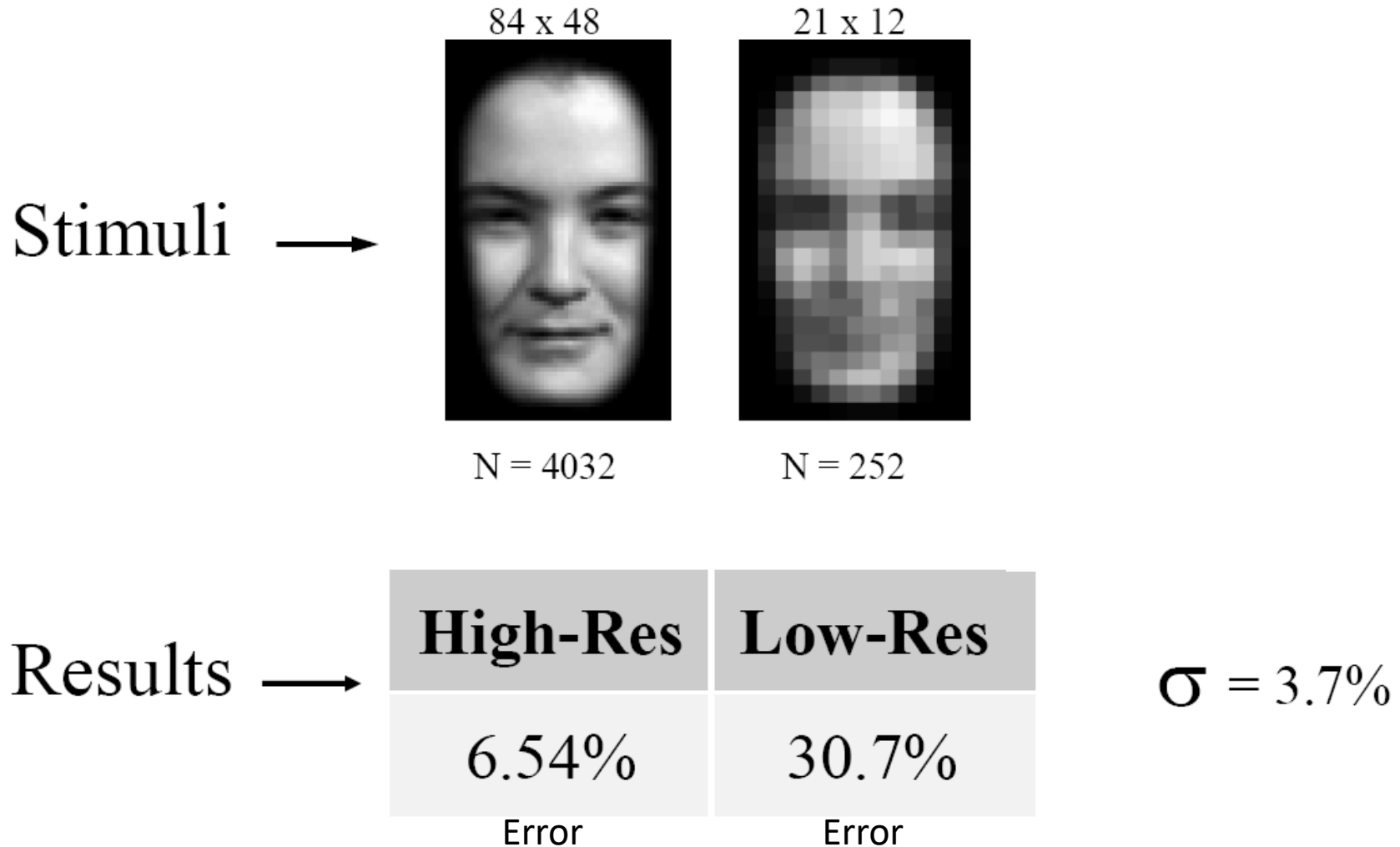
Classifier Performance

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%

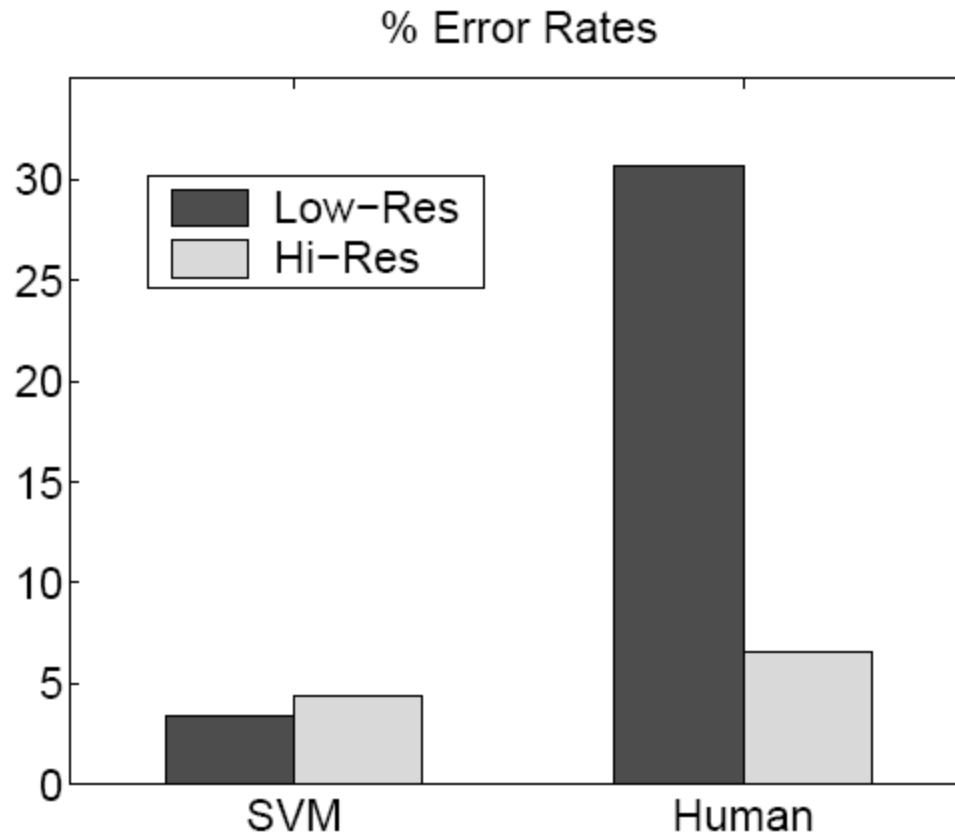
Gender perception experiment: How well can humans do?

- 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

Gender perception experiment: How well can humans do?



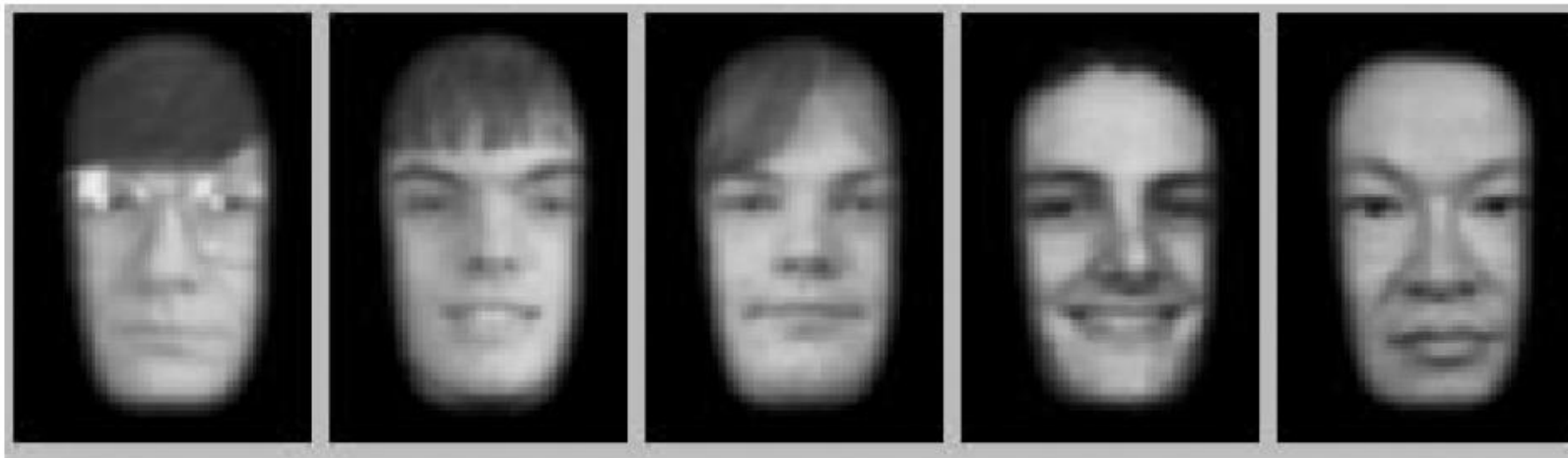
Human vs. Machine



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

Figure 6. SVM vs. Human performance

Hardest examples for humans



Top five human misclassifications

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

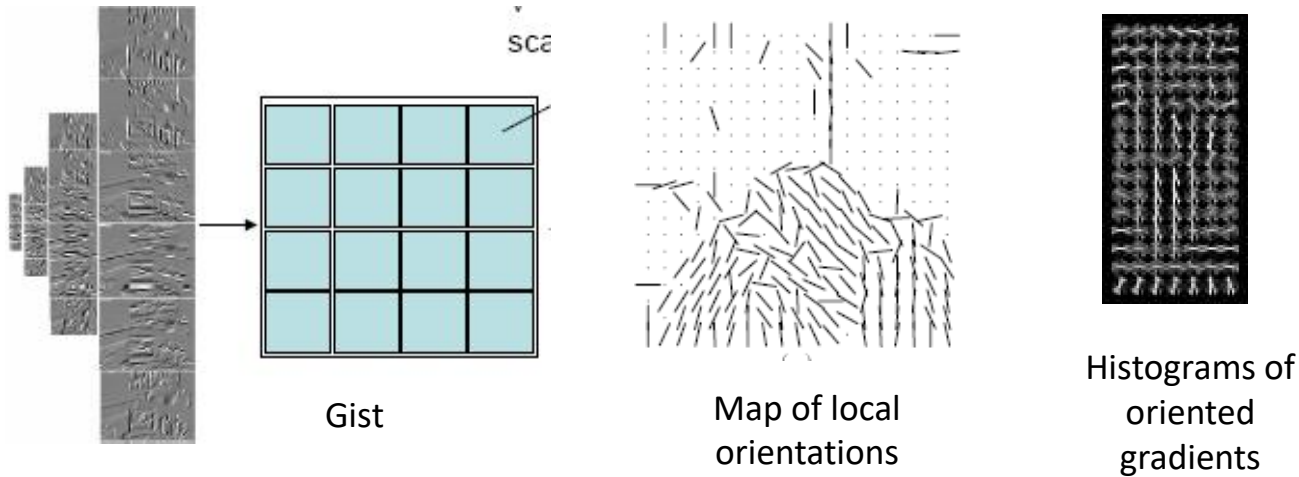
SVMs: Pros and cons

- 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

Summary

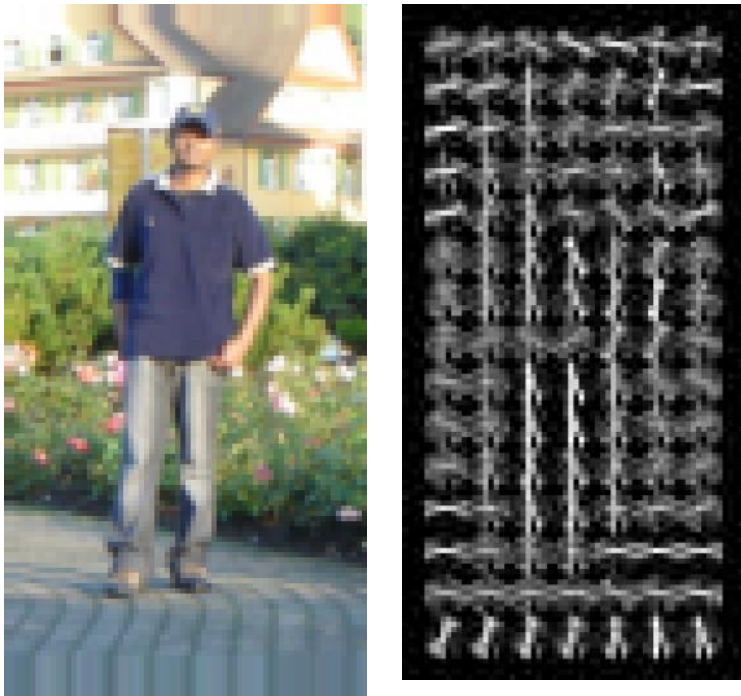
- Discriminative classifiers
 - Boosting
 - Nearest neighbors
 - Support vector machines
- Useful for object recognition when combined with “window-based” or holistic appearance descriptors

Global window-based appearance representations

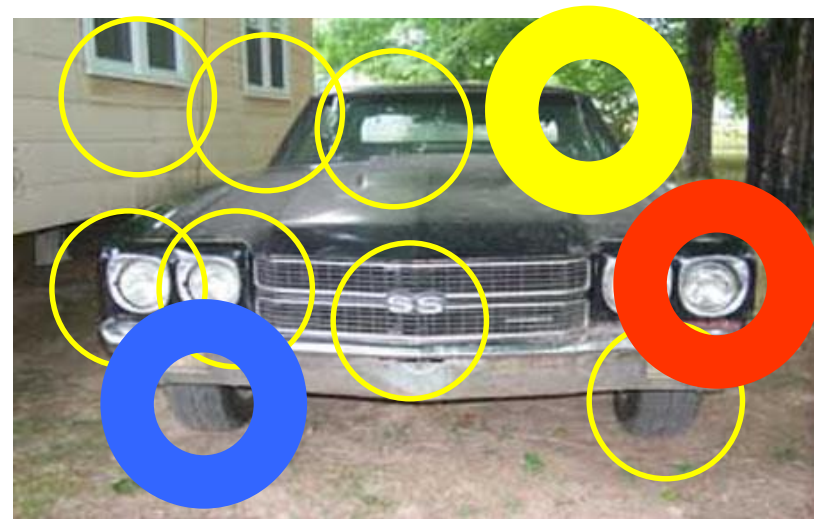


- These examples are truly global; each pixel in the window contributes to the representation.
- Classifier can account for relative relevance...
- *When might this not be ideal?*

Generic category recognition: representation choice



Window-based



Part-based