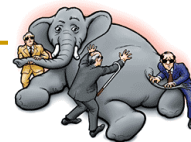

Representing Images

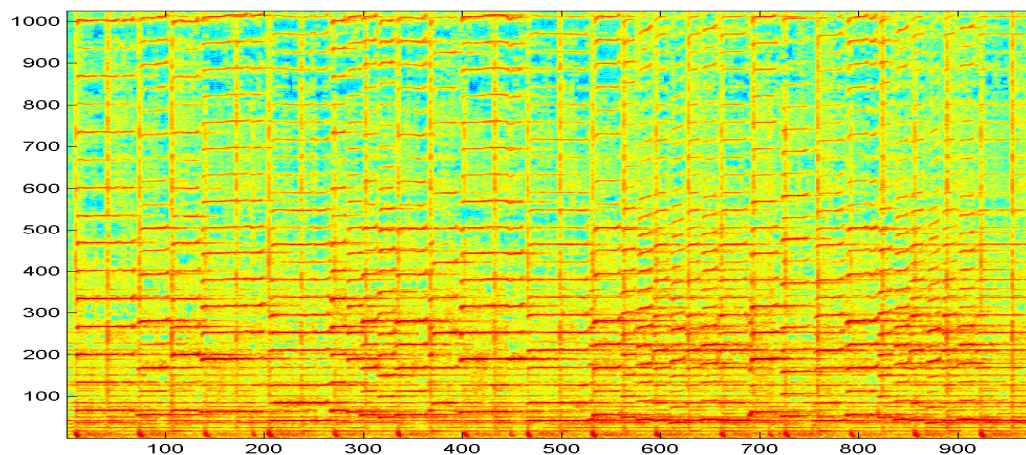
Detecting faces in images

Class 5. 15 Sep 2009

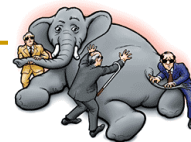
Instructor: Bhiksha Raj



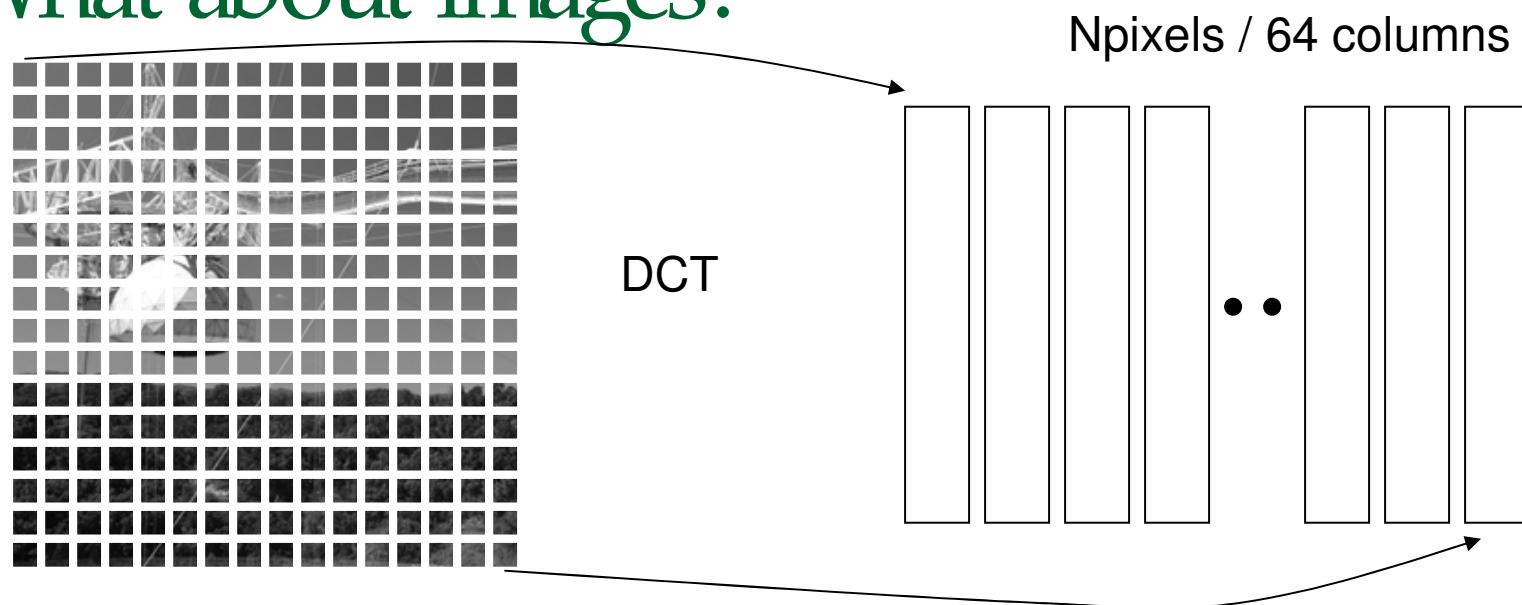
Last Class: Representing Audio



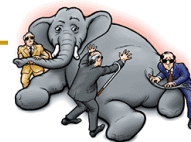
- n Basic DFT
- n Computing a Spectrogram
- n Computing additional features from a spectrogram



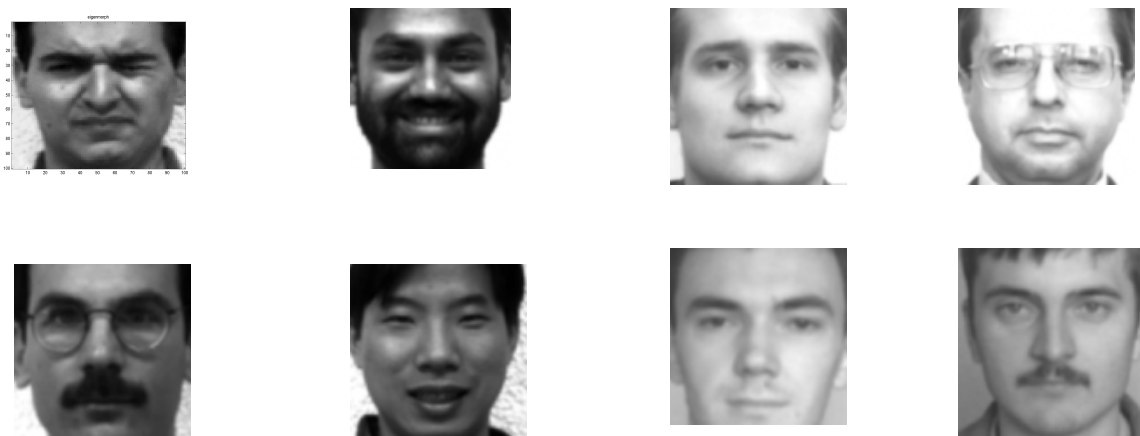
What about images?



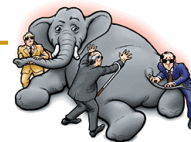
- n DCT of small segments
 - q 8x8
 - q Each image becomes a matrix of DCT vectors
- n DCT of the image
- n Haar transform (checkerboard)
- n ***Or data-driven representations..***



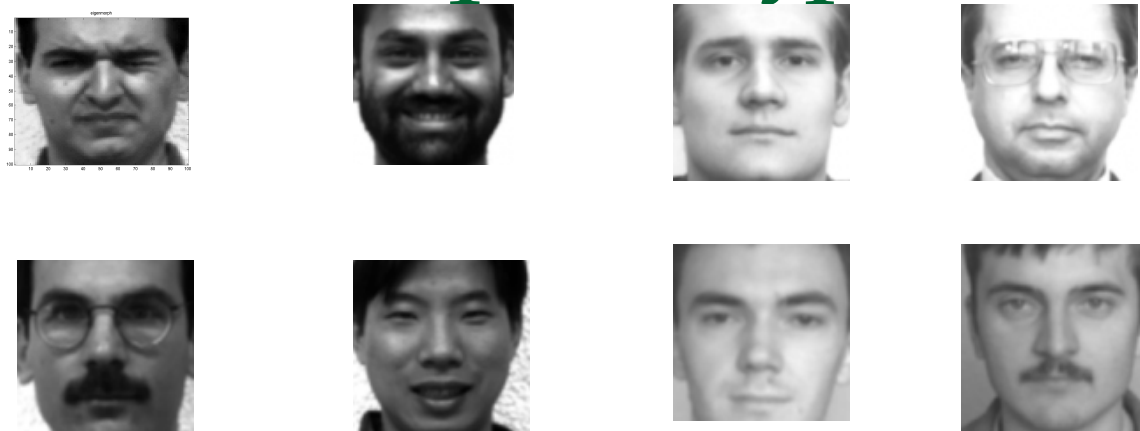
Returning to Eigen Computation



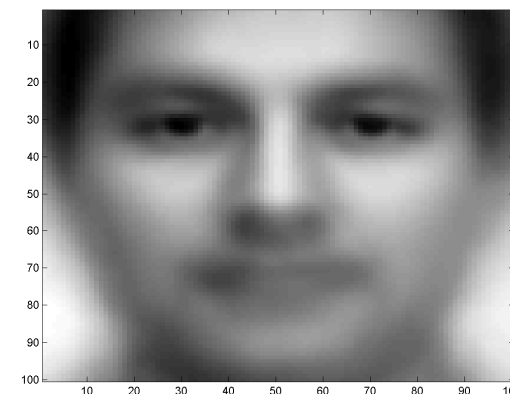
- n A collection of faces
 - q All normalized to 100x100 pixels
- n What is common among all of them?
 - q Do we have a common descriptor?



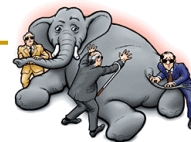
A least squares typical face



The typical face



- n Can we do better than a blank screen to find the most common portion of faces?
 - q The first checkerboard; the zeroth frequency component..
- n Assumption: There is a “typical” face that captures most of what is common to all faces
 - q Every face can be represented by a scaled version of a typical face
 - n What is this face?
- n Approximate **every** face f as $f = w_f V$
- n Estimate V to minimize the squared error
 - q How?
 - q What is V ?



A collection of least squares typical faces

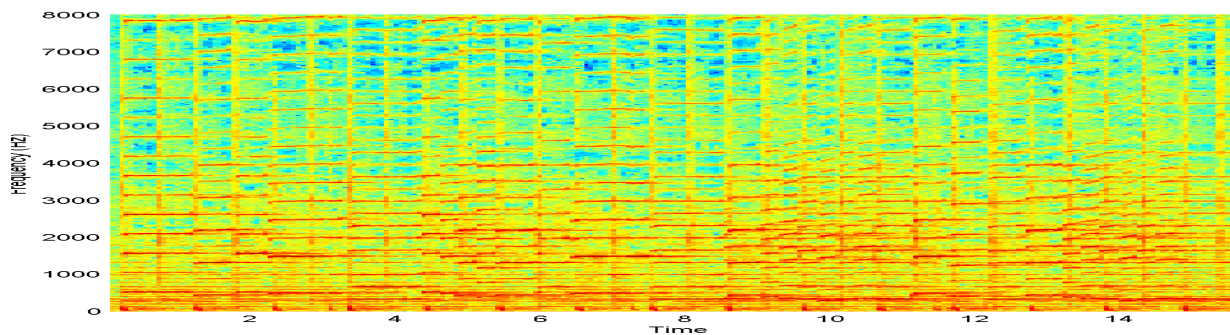


- n Assumption: There are a set of K “typical” faces that captures most of all faces
- n Approximate **every** face f as $f = w_{f,1} V_1 + w_{f,2} V_2 + w_{f,3} V_3 + \dots + w_{f,k} V_k$
 - q V_2 is used to “correct” errors resulting from using only V_1
 - n So the total energy in $w_{f,2}$ ($\sum w_{f,2}^2$) must be lesser than the total energy in $w_{f,1}$ ($\sum w_{f,1}^2$)
 - q V_3 corrects errors remaining after correction with V_2
 - n The total energy in $w_{f,3}$ must be lesser than that even in $w_{f,2}$
 - q And so on..
 - q $V = [V_1 V_2 V_3]$
- n Estimate V to minimize the squared error
 - q How?
 - q What is V ?

A recollection

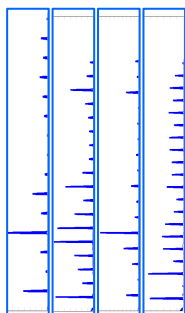


M =



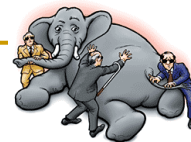
$$V = \text{PINV}(W) * M$$

W =



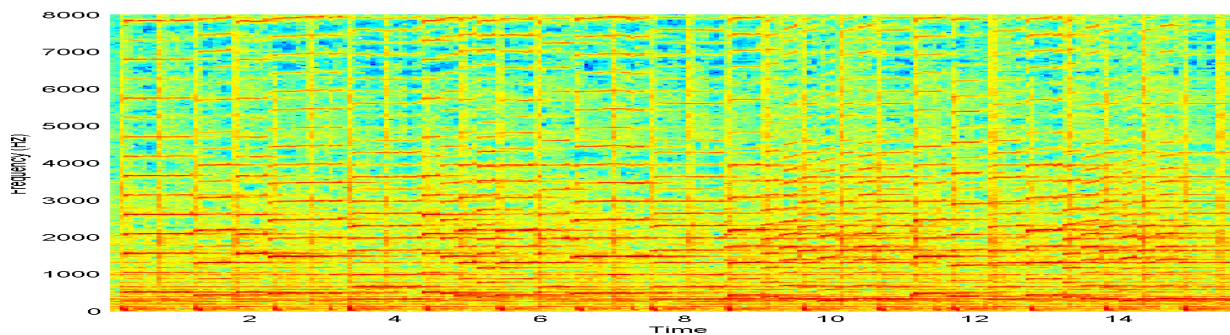
U = ?



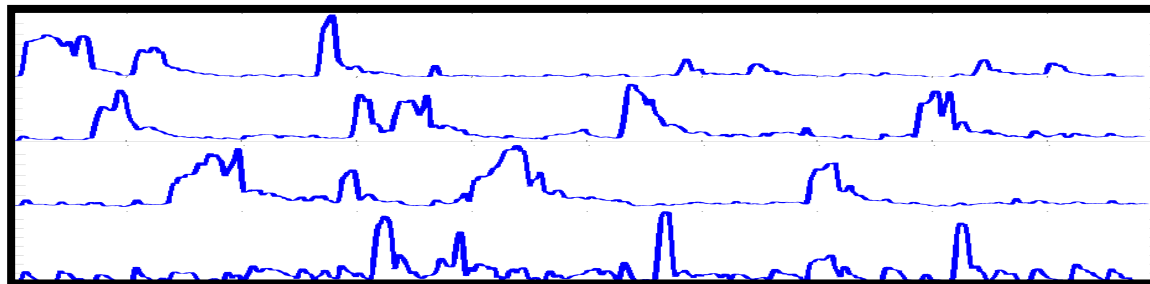


How about the other way?

M =



V =



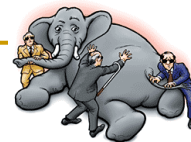
W =

?

U =

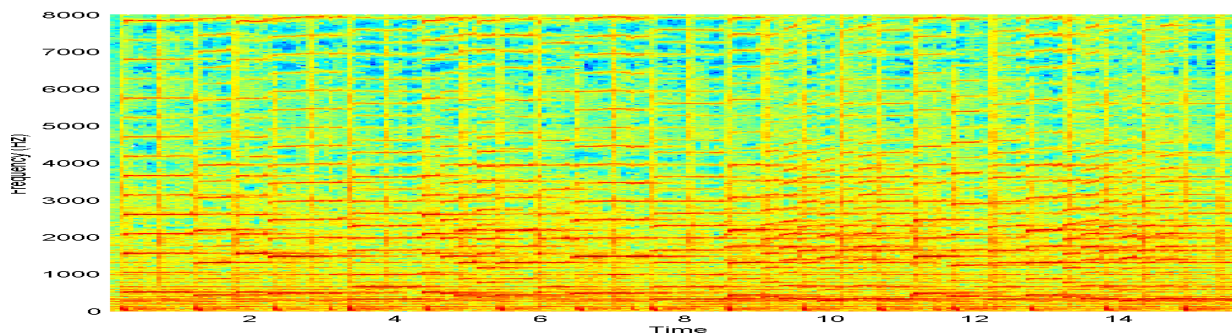
?

n $W = M * \text{Pinv}(V)$



How about the other way?

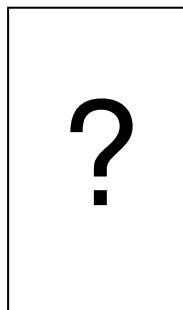
M =



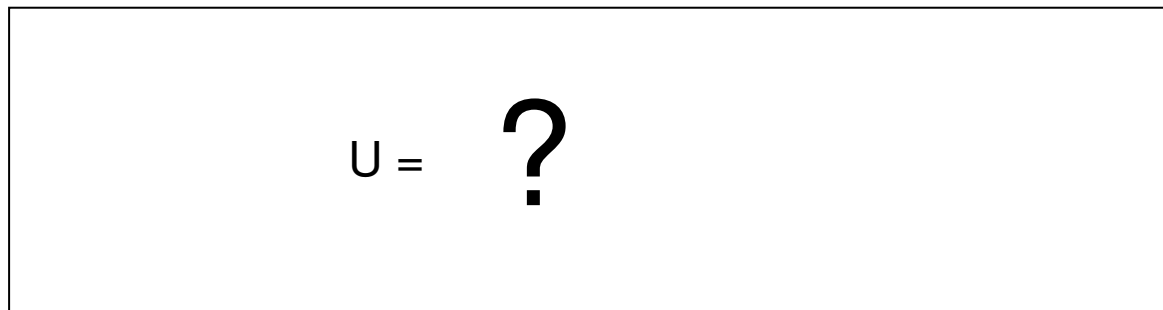
V =



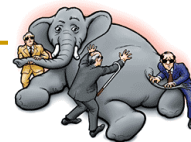
W =



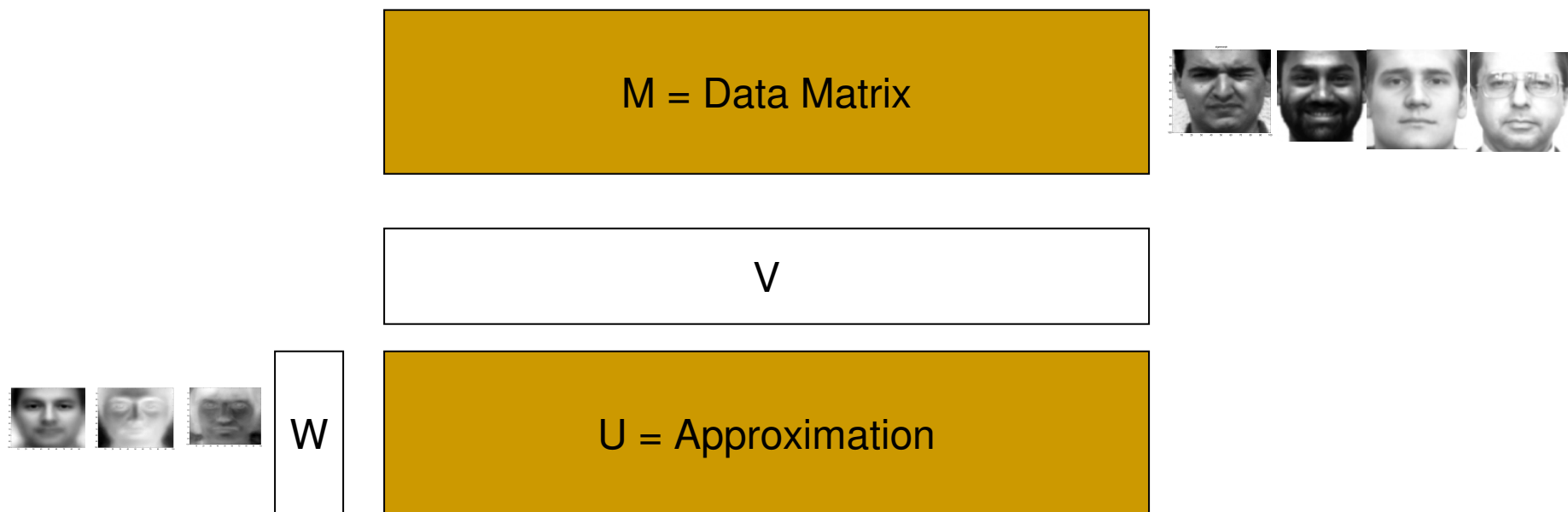
U = ?



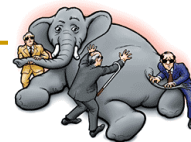
$$W V \approx M$$



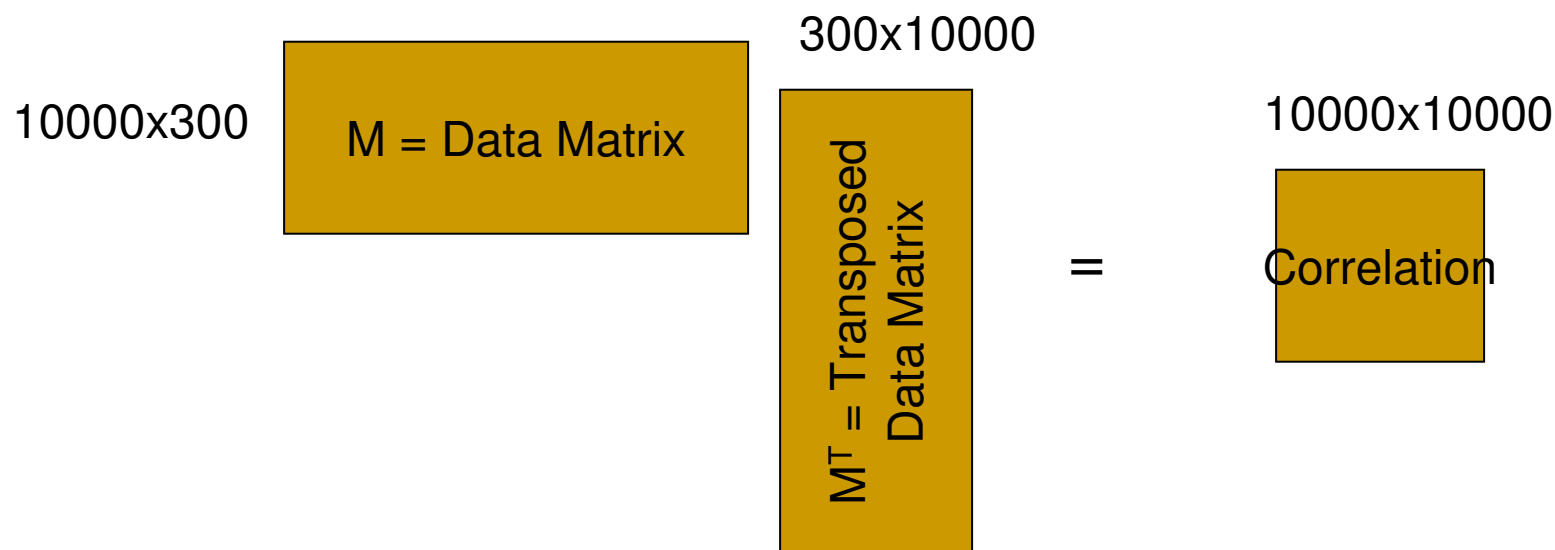
Eigen Faces!



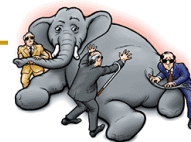
- n Here W , V and U are ALL unknown and must be determined
- q Such that the squared error between U and M is minimum
- n Eigen analysis allows you to find W and V such that $U = WV$ has the least squared error with respect to the original data M
- n If the original data are a collection of faces, the columns of W are *eigen faces*.



Eigen faces



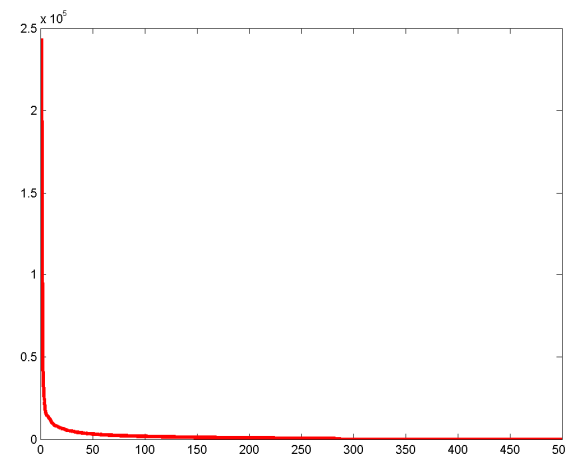
- n Lay all faces side by side in vector form to form a matrix
 - q In my example: 300 faces. So the matrix is 10000 x 300
- n Multiply the matrix by its transpose
 - q The correlation matrix is 10000x10000



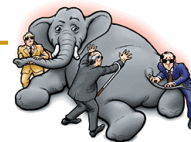
Eigen faces

$$[U, S] = \text{eig}(\text{correlation})$$

$$S = \begin{bmatrix} \lambda_1 & \cdot & 0 & \cdot & 0 \\ 0 & \lambda_2 & 0 & \cdot & 0 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & \cdot & 0 & \cdot & \lambda_{10000} \end{bmatrix} \quad U = \begin{bmatrix} \text{eigenface1} \\ \text{eigenface2} \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}$$

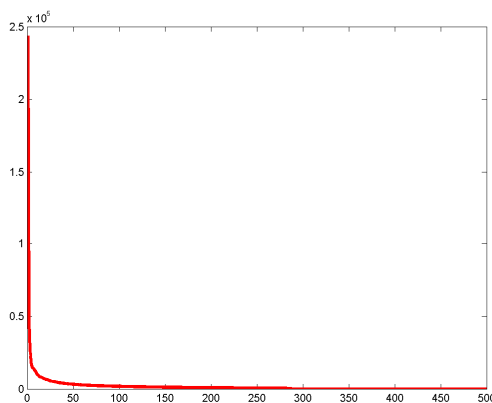


- n Compute the eigen vectors
 - q Only 300 of the 10000 eigen values are non-zero
 - n Why?
- n Retain eigen vectors with high eigen values (>0)
 - q Could use a higher threshold



Eigen Faces

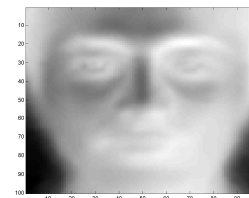
$$U = \begin{bmatrix} \text{eigenface1} \\ \text{eigenface2} \\ \bullet \\ \bullet \\ \bullet \end{bmatrix}$$



eigenface1



eigenface2

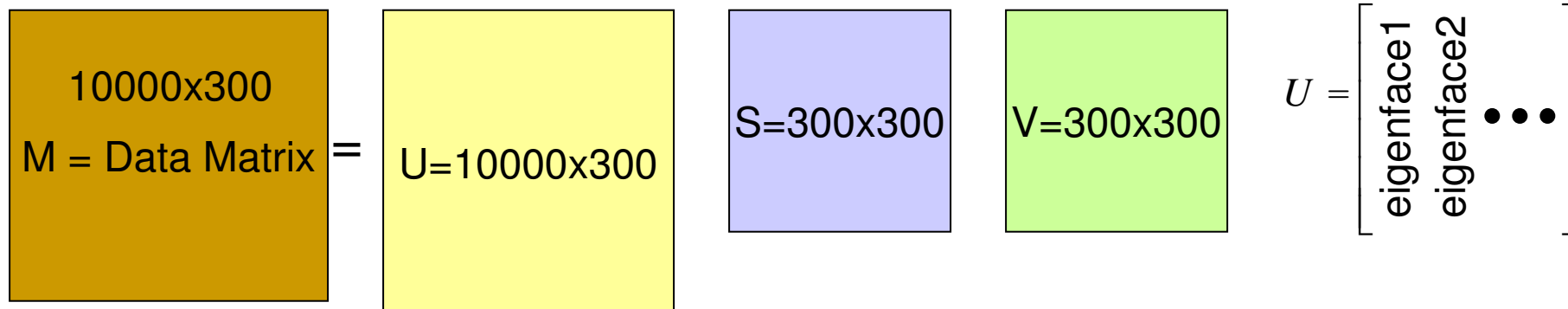


eigenface3

- n The eigen vector with the highest eigen value is the first typical face
- n The vector with the second highest eigen value is the second typical face.
- n Etc.



SVD instead of Eigen



n Do we need to compute a 10000 x 10000 correlation matrix and then perform Eigen analysis?

q Will take a very long time on your laptop

n SVD

q Only need to perform “Thin” SVD. Very fast

n $U = 10000 \times 300$

q The columns of U are the eigen faces!

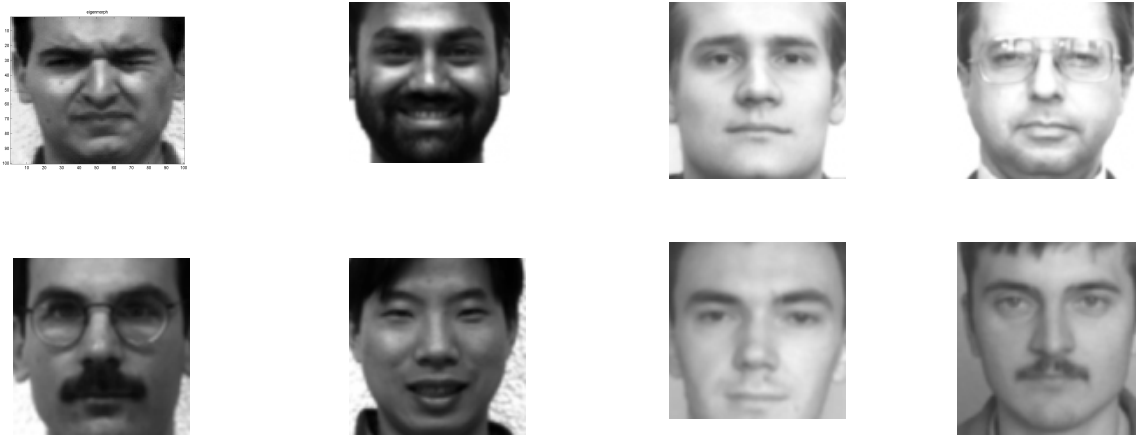
q The U s corresponding to the “zero” eigen values are not computed

n $S = 300 \times 300$

n $V = 300 \times 300$

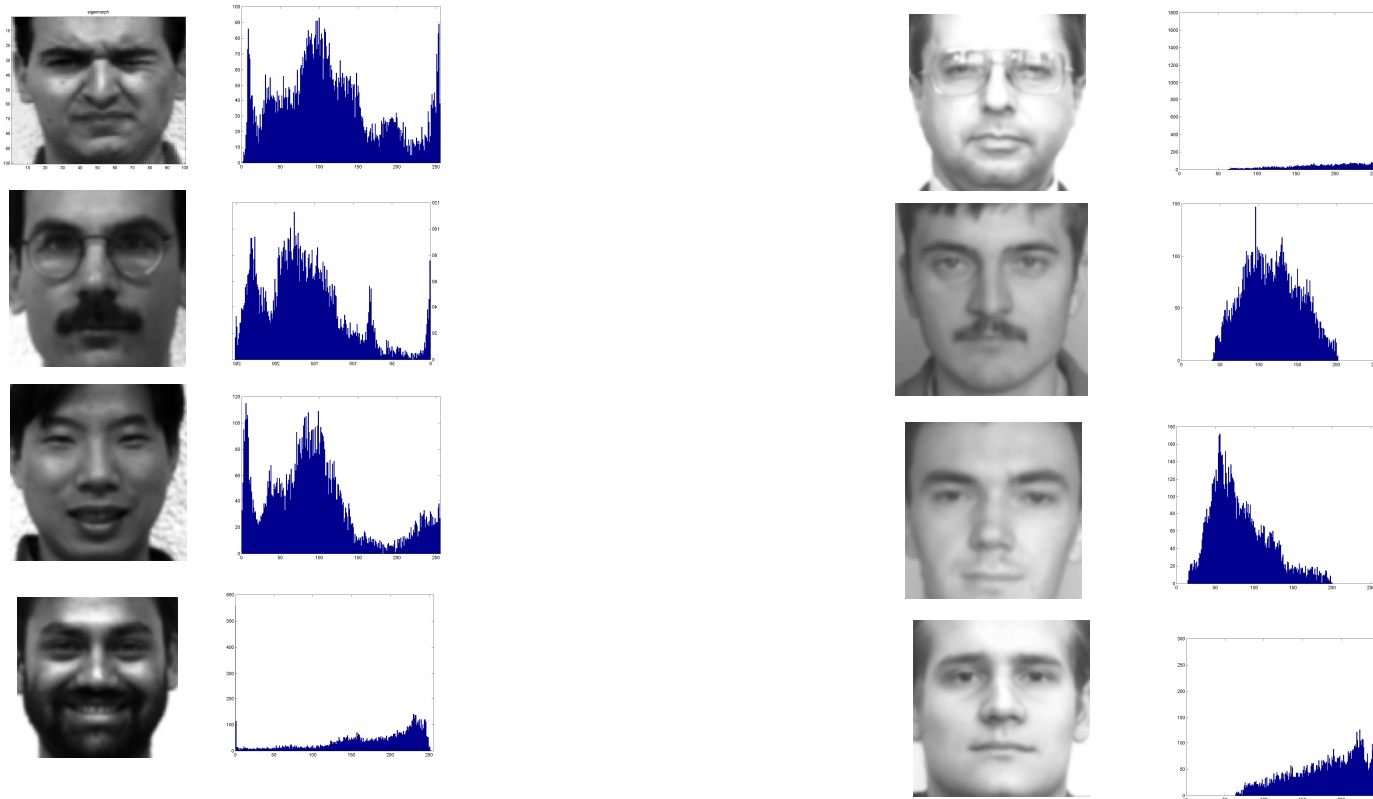
NORMALIZING OUT VARIATIONS

Images: Accounting for variations



- n What are the obvious differences in the above images
- n How can we capture these differences
 - q Hint – image histograms..

Images -- Variations



n Pixel histograms: what are the differences

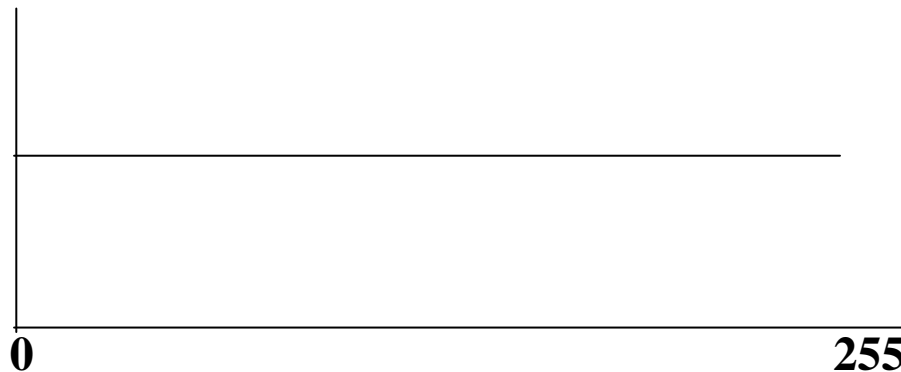
Normalizing Image Characteristics

- n Normalize the pictures
 - q Eliminate lighting/contrast variations
 - q All pictures must have “similar” lighting
 - n How?

- n Lighting and contrast are represented in the image histograms:

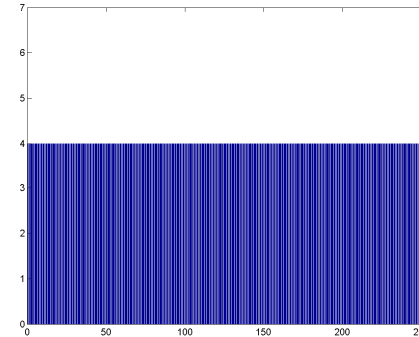
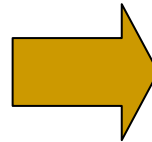
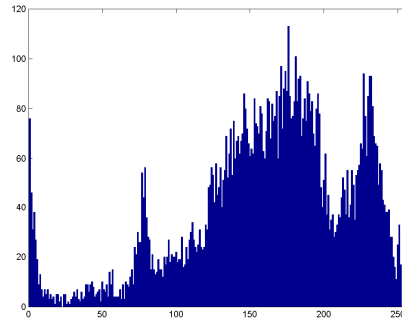
Histogram Equalization

- n Normalize histograms of images
 - q Maximize the contrast
 - n Contrast is defined as the “flatness” of the histogram
 - n For maximal contrast, every greyscale must happen as frequently as every other greyscale



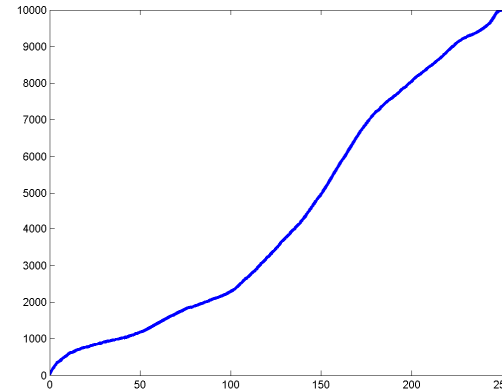
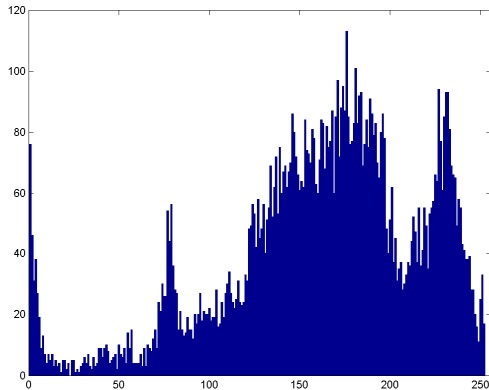
- n Maximizing the contrast: Flattening the histogram
 - q Doing it for every image ensures that every image has the same contrast
 - n I.e. exactly the same histogram of pixel values
 - q Which should be flat

Histogram Equalization



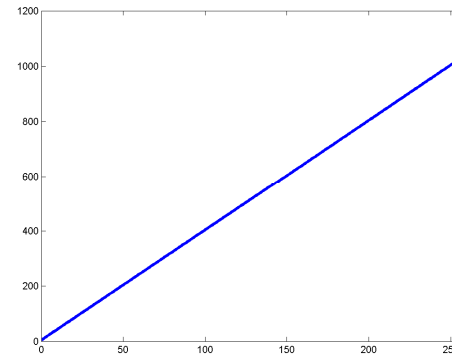
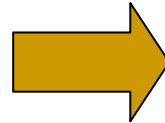
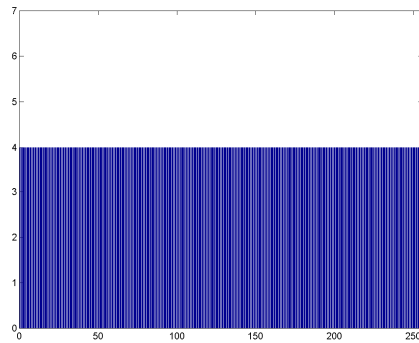
- n Modify pixel values such that histogram becomes “flat”.
- n For each pixel
 - q New pixel value = $f(\text{old pixel value})$
 - q What is $f()$?
- n Easy way to compute this function: map cumulative counts

Cumulative Count Function

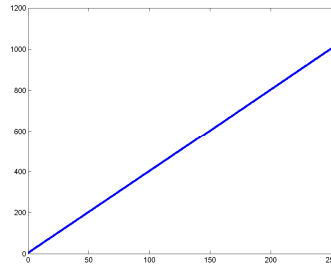
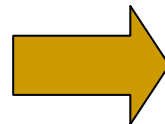
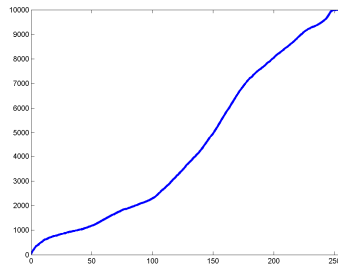


- n The *histogram (count)* of a pixel value X is the number of pixels in the image that have value X
 - q E.g. in the above image, the count of pixel value 180 is about 110
- n The *cumulative count* at pixel value X is the total number of pixels that have values in the range $0 \leq x \leq X$
 - q $CCF(X) = H(1) + H(2) + \dots + H(X)$

Cumulative Count Function

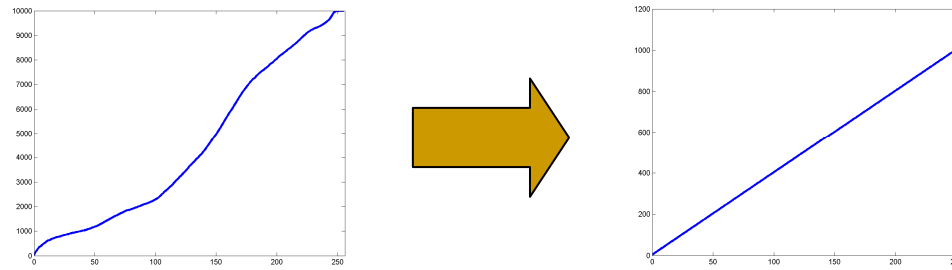


- n The cumulative count function of a uniform histogram is a line



- n We must modify the pixel values of the image so that its cumulative count is a line

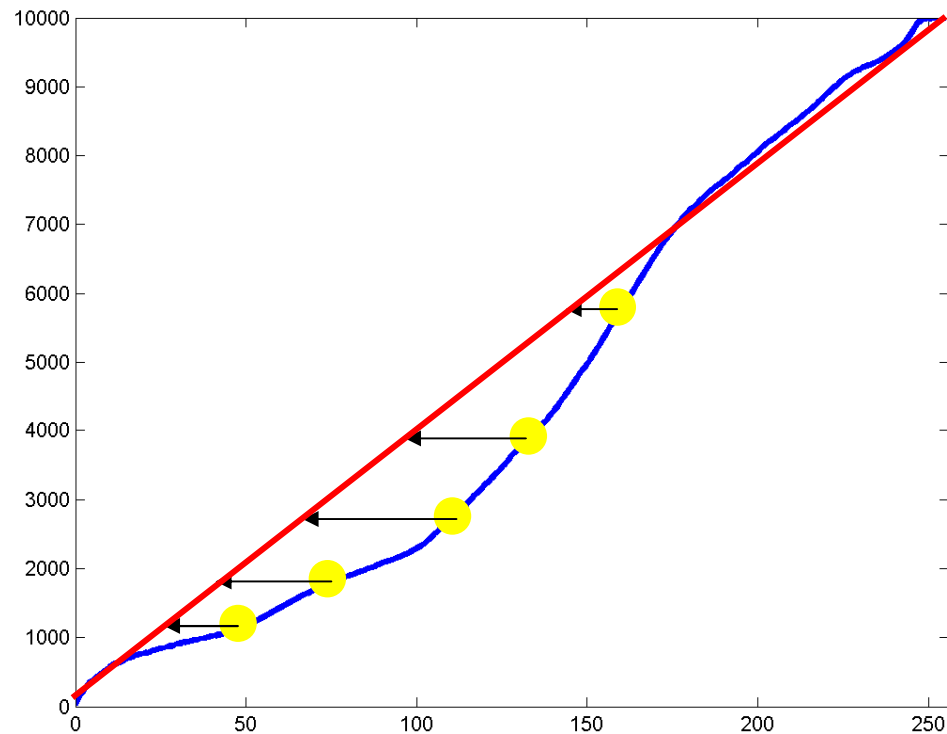
Mapping CCFs



Move x axis levels around until the plot to the left looks like the plot to the right

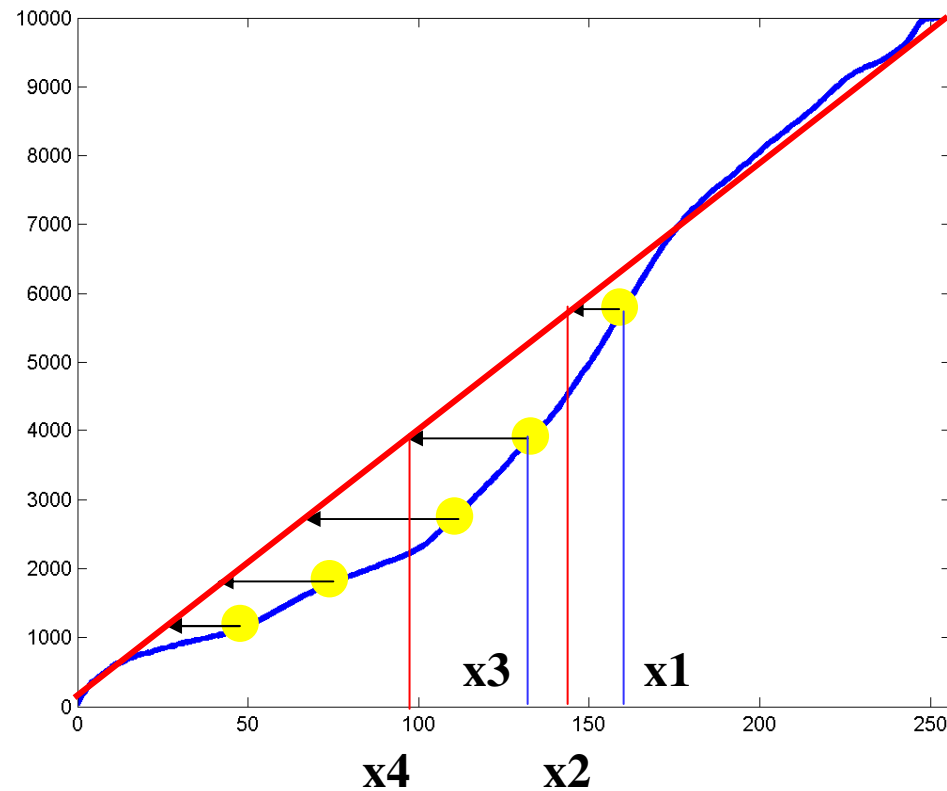
- n CCF($f(x)$) $\rightarrow a * f(x)$ [of $a * (f(x)+1)$ if pixels can take value 0]
 - q x = pixel value
 - q $f()$ is the function that converts the old pixel value to a new (normalized) pixel value
 - q $a = (\text{total no. of pixels in image}) / (\text{total no. of pixel levels})$
 - n The no. of pixel levels is 256 in our examples
 - n Total no. of pixels is 10000 in a 100x100 image

Mapping CCFs



- n For each pixel value x :
 - q Find the location on the red line that has the closest Y value to the observed CCF at x

Mapping CCFs



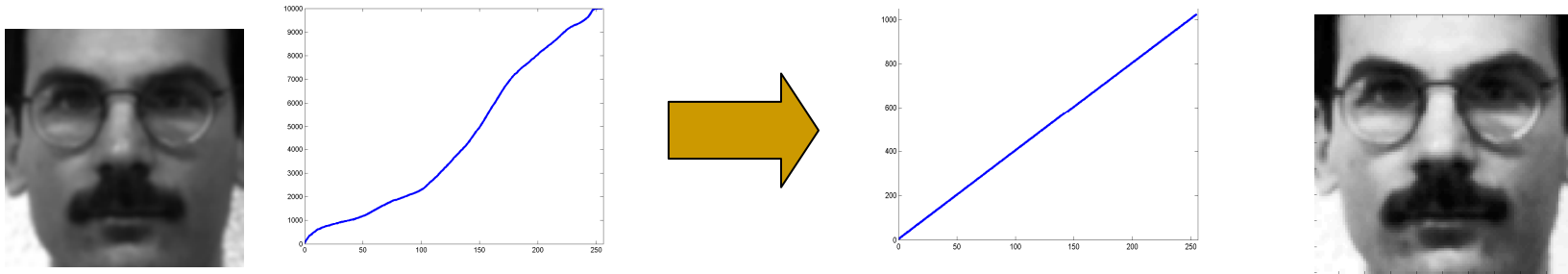
$$f(x_1) = x_2$$

$$f(x_3) = x_4$$

Etc.

- n For each pixel value x :
 - q Find the location on the red line that has the closet Y value to the observed CCF at x

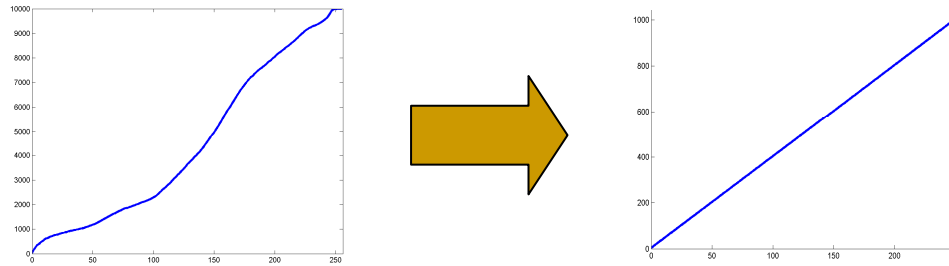
Mapping CCFs



Move x axis levels around until the plot to the left looks like the plot to the right

- n For each pixel in the image to the left
 - q The pixel has a value x
 - q Find the CCF at that pixel value $CCF(x)$
 - q Find x' such that $CCF(x')$ in the function to the right equals $CCF(x)$
 - n x' such that $CCF_flat(x') = CCF(x)$
 - q Modify the pixel value to x'

Doing it Formulaically



$$f(x) = \text{round} \left(\frac{CCF(x) - CCF_{\min}}{N_{\text{pixels}} - CCF_{\min}} \text{Max.pixel.value} \right)$$

- n CCF_{\min} is the smallest non-zero value of $CCF(x)$
 - q The value of the CCF at the smallest observed pixel value
- n N_{pixels} is the total no. of pixels in the image
 - q 10000 for a 100x100 image
- n Max.pixel.value is the highest pixel value
 - q 255 for 8-bit pixel representations

Or even simpler

n Matlab:

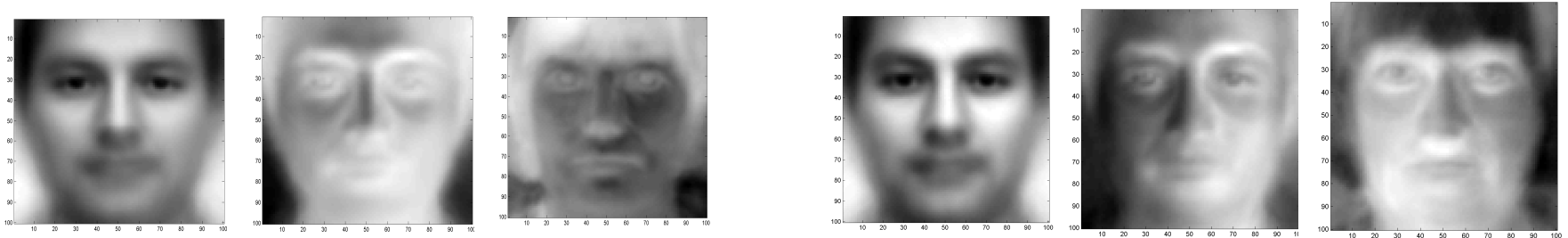
q `Newimage = histeq(oldimage)`

Histogram Equalization



- n Left column: Original image
- n Right column: Equalized image
- n All images now have similar contrast levels

Eigenfaces after Equalization



n Left panel : Without HEQ

n Right panel: With HEQ

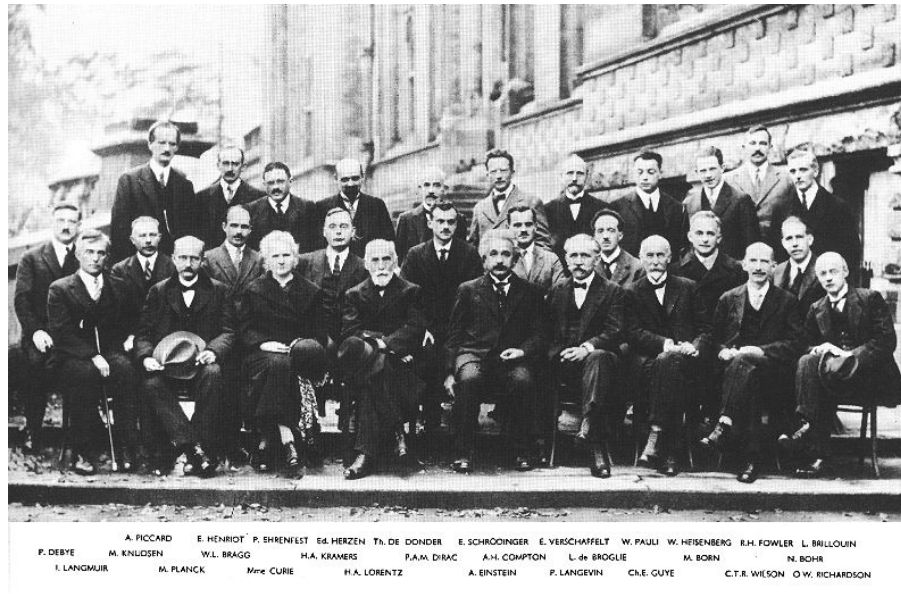
q Eigen faces are more face like..

n Need not always be the case



Detecting Faces in Images

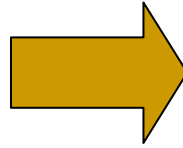
Detecting Faces in Images



- n Finding face like patterns
 - q How do we find if a picture has faces in it
 - q Where are the faces?

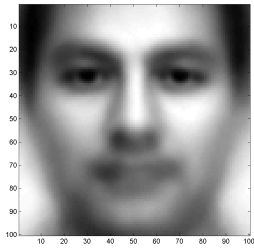
- n A simple solution:
 - q Define a “typical face”
 - q Find the “typical face” in the image

Finding faces in an image



- n Picture is larger than the “typical face”
 - q E.g. typical face is 100x100, picture is 600x800
- n First convert to greyscale
 - q $R + G + B$
 - q Not very useful to work in color

Finding faces in an image



- n Goal .. To find out if and where images that look like the “typical” face occur in the picture

Finding faces in an image



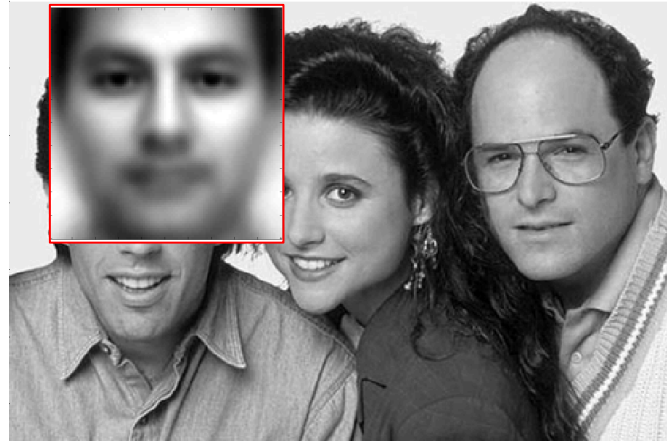
- n Try to “match” the typical face to each location in the picture

Finding faces in an image



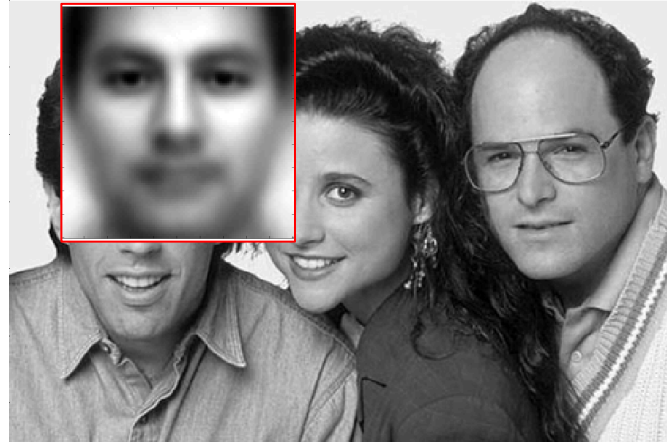
- n Try to “match” the typical face to each location in the picture

Finding faces in an image



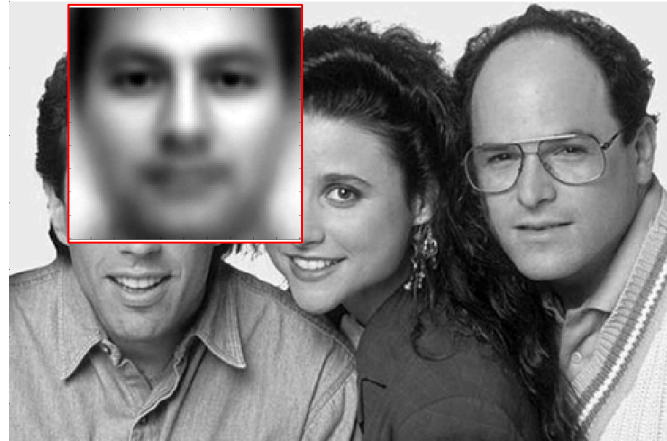
- n Try to “match” the typical face to each location in the picture

Finding faces in an image



- n Try to “match” the typical face to each location in the picture

Finding faces in an image



- n Try to “match” the typical face to each location in the picture

Finding faces in an image



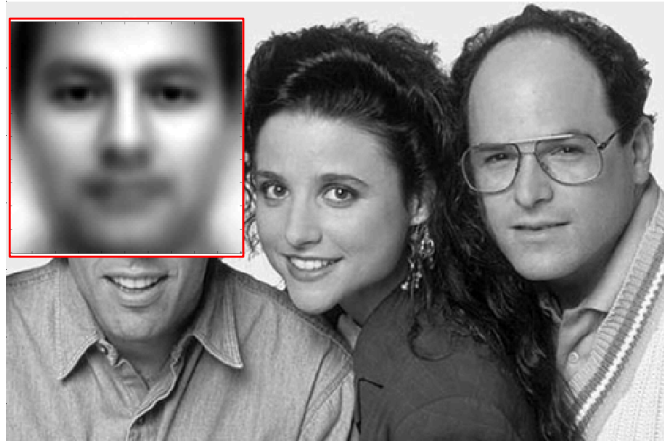
- n Try to “match” the typical face to each location in the picture

Finding faces in an image



- n Try to “match” the typical face to each location in the picture

Finding faces in an image



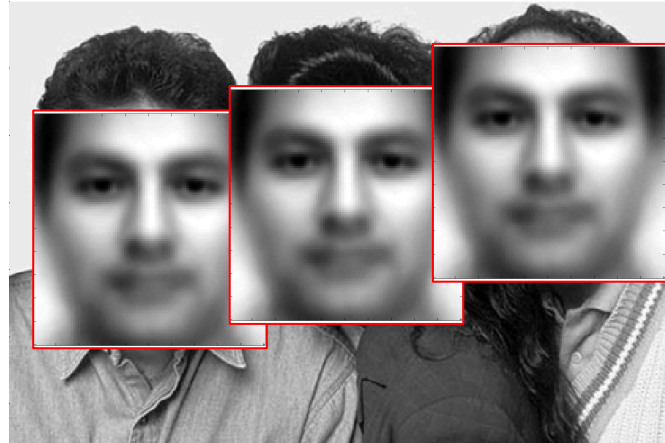
- n Try to “match” the typical face to each location in the picture

Finding faces in an image



- n Try to “match” the typical face to each location in the picture

Finding faces in an image



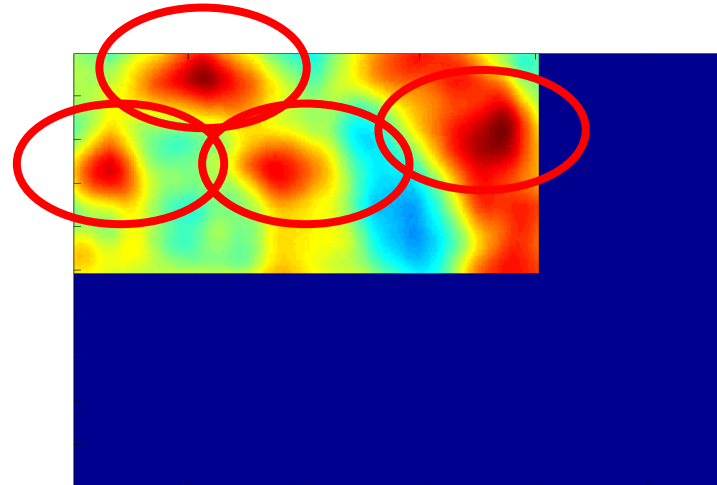
- n Try to “match” the typical face to each location in the picture
- n The “typical face” will explain some spots on the image much better than others
 - q These are the spots at which we probably have a face!

How to “match”



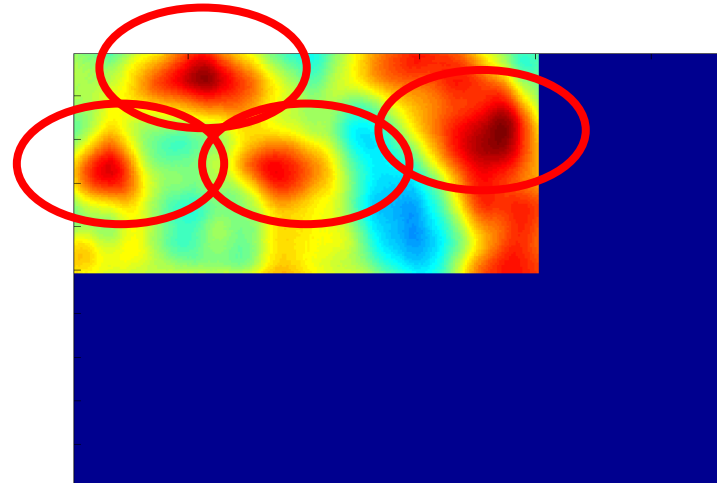
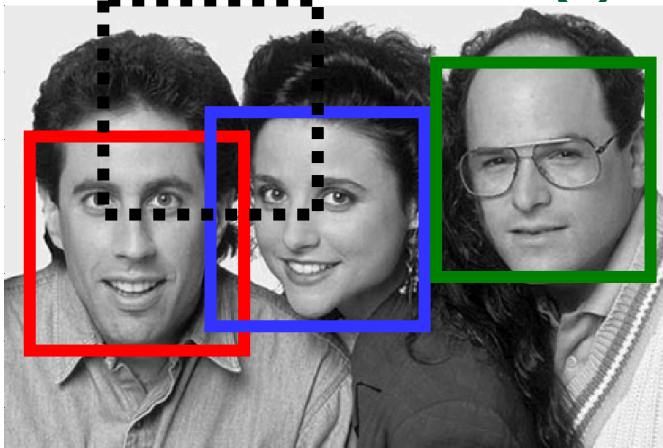
- n What exactly is the “match”
 - q What is the match “score”
- n The DOT Product
 - q Express the typical face as a vector
 - q Express the region of the image being evaluated as a vector
 - n But first histogram equalize the region
 - q Just the section being evaluated, without considering the rest of the image
 - q Compute the dot product of the typical face vector and the “region” vector

What do we get



- n The right panel shows the dot product at various locations
- q Redder is higher
 - n The locations of peaks indicate locations of faces!

What do we get

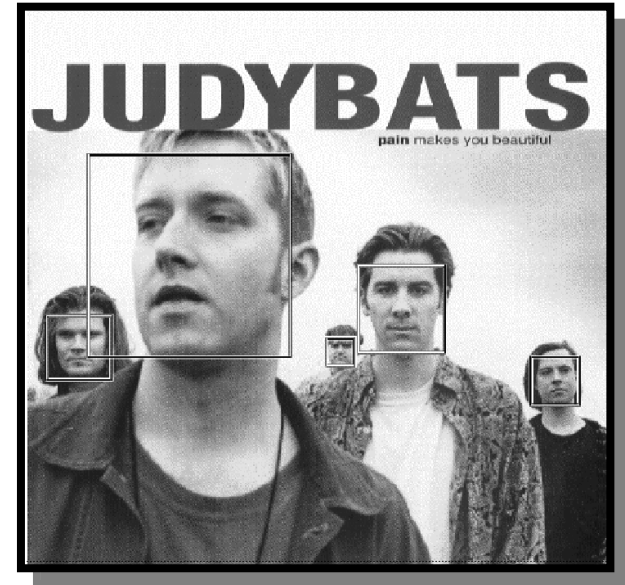


- n The right panel shows the dot product a various loctions
 - q Redder is higher
 - n The locations of peaks indicate locations of faces!
- n Correctly detects all three faces
 - q Likes George's face most
 - n He looks most like the typical face
- n Also finds a face where there is none!
 - q A false alarm

Scaling and Rotation Problems

n Scaling

- q Not all faces are the same size
- q Some people have bigger faces
- q The size of the face on the image changes with perspective
- q Our “typical face” only represents one of these sizes

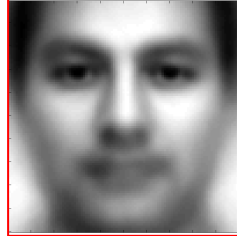
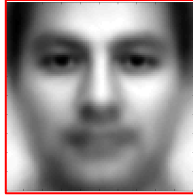


n Rotation

- q The head need not always be upright!
 - n Our typical face image was upright



Solution



- n Create many “typical faces”
 - q One for each scaling factor
 - q One for each rotation
 - n How will we do this?
- n Match them all
- n Does this work
 - q Kind of .. Not well enough at all
 - q We need more sophisticated models



Face Detection: A Quick Historical Perspective

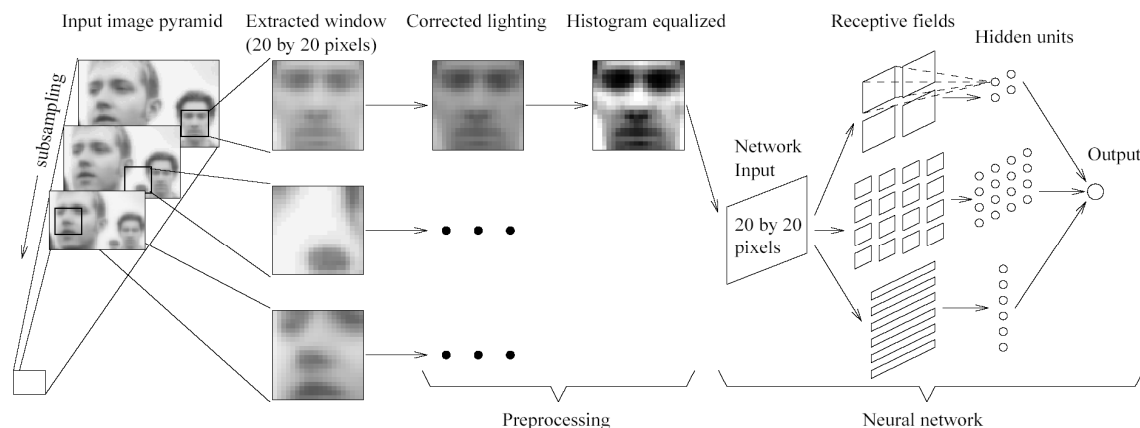


Figure 1: The basic algorithm used for face detection.

- n Many more complex methods
 - q Use edge detectors and search for face like patterns
 - q Find “feature” detectors (noses, ears..) and employ them in complex neural networks..

- n The Viola Jones method
 - q Boosted cascaded classifiers

- n But first, what is boosting

And even before that – what is classification?

- n Given “features” describing an entity, determine the category it belongs to
 - q Walks on two legs, has no hair. Is this
 - n A Chimpanzee
 - n A Human
 - q Has long hair, is 5’4” tall, is this
 - n A man
 - n A woman
 - q Matches “eye” pattern with score 0.5, “mouth pattern” with score 0.25, “nose” pattern with score 0.1. Are we looking at
 - n A face
 - n Not a face?

Classification

- n Multi-class classification

- q Many possible categories

- n E.g. Sounds “AH, IY, UW, EY..”

- n E.g. Images “Tree, dog, house, person..”

- n Binary classification

- q Only two categories

- n Man vs. Woman

- n Face vs. not a face..

- n Face detection: Recast as binary face classification

- q For each little square of the image, determine if the square represents a face or not

Face Detection as Classification



For each square, run a classifier to find out if it is a face or not

- n Faces can be many sizes
- n They can happen anywhere in the image
- n For each face size
 - q For each location
 - n Classify a rectangular region of the face size, at that location, as a face or not a face
- n This is a series of **binary** classification problems

Introduction to Boosting

- n An *ensemble* method that sequentially combines many simple **BINARY** classifiers to construct a final complex classifier
 - q Simple classifiers are often called “weak” learners
 - q The complex classifiers are called “strong” learners
- n Each weak learner focuses on instances where the previous classifier failed
 - q Give greater weight to instances that have been incorrectly classified by previous learners
- n Restrictions for weak learners
 - q Better than 50% correct
- n Final classifier is *weighted* sum of weak classifiers

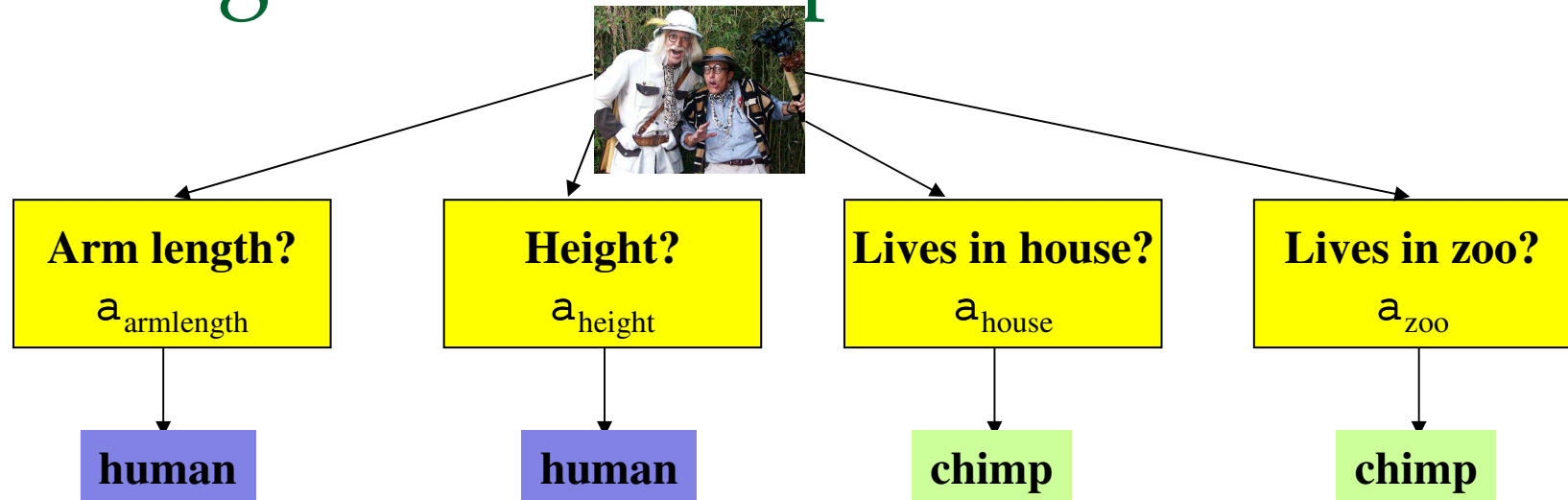
Boosting: A very simple idea

- n One can come up with many rules to classify
 - q E.g. Chimpanzee vs. Human classifier:
 - q If arms == long, entity is chimpanzee
 - q If height > 5'6" entity is human
 - q If lives in house == entity is human
 - q If lives in zoo == entity is chimpanzee

- n Each of them is a reasonable rule, but makes many mistakes
 - q Each rule has an intrinsic error rate

- n *Combine* the predictions of these rules
 - q But not equally
 - q Rules that are less accurate should be given lesser weight

Boosting and the Chimpanzee Problem



- n The total confidence in all classifiers that classify the entity as a chimpanzee is

$$Score_{chimp} = \sum_{\text{classifier favors chimpanzee}} a_{\text{classifier}}$$

- n The total confidence in all classifiers that classify it as a human is

$$Score_{human} = \sum_{\text{classifier favors human}} a_{\text{classifier}}$$

- n If $Score_{chimpanzee} > Score_{human}$ then the our belief that we have a chimpanzee is greater than the belief that we have a human

Boosting as defined by Freund

- n A gambler wants to write a program to predict winning horses. His program must encode the expertise of his brilliant winner friend
- n The friend has no single, encodable algorithm. Instead he has many rules of thumb
 - q He uses a different rule of thumb for each set of races
 - n E.g. “in this set, go with races that have black horses with stars on their foreheads”
 - q But cannot really enumerate what rules of thumbs go with what sets of races: he simply “knows” when he encounters a set
 - n A common problem that faces us in many situations
- n Problem:
 - q How best to combine all of the friend’s rules of thumb
 - q What is the best set of races to present to the friend, to extract the various rules of thumb

Boosting

- n The basic idea: Can a “weak” learning algorithm that performs just slightly better than random guessing be *boosted* into an arbitrarily accurate “strong” learner
 - q Each of the gambler’s rules may be just better than random guessing
- n This is a “meta” algorithm, that poses no constraints on the form of the weak learners themselves
 - q The gambler’s rules of thumb can be anything

Boosting: A Voting Perspective

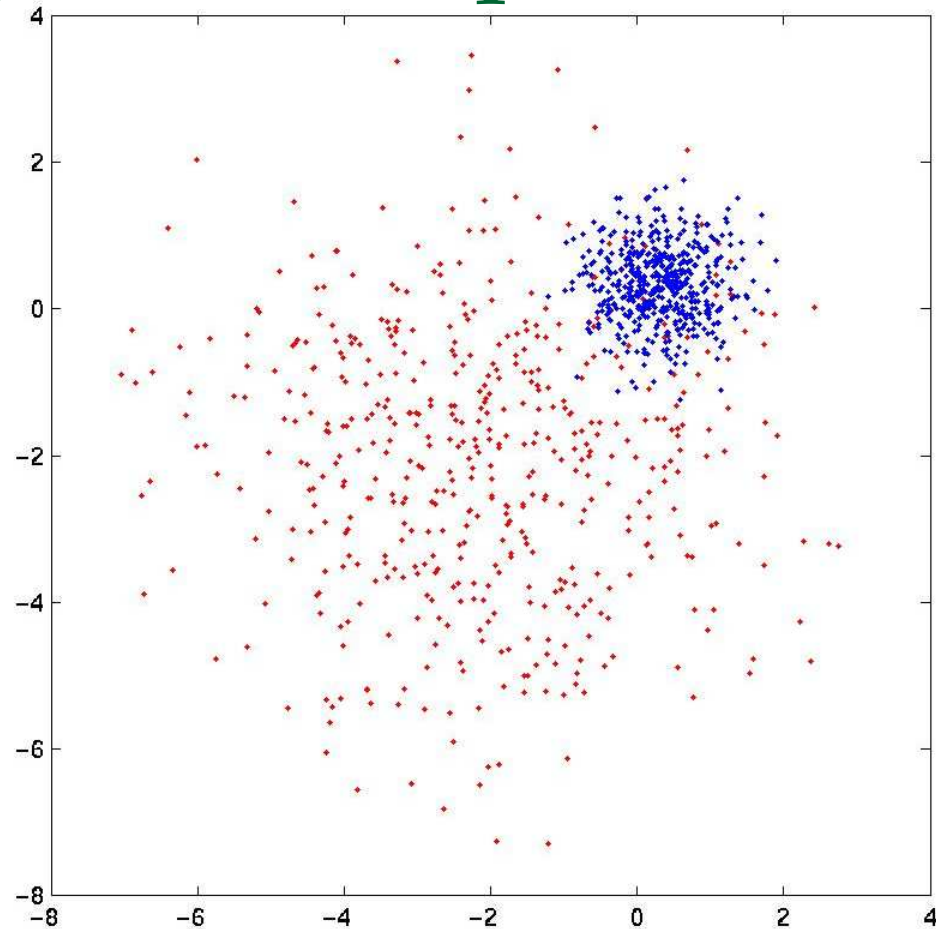
- n Boosting can be considered a form of voting
 - q Let a number of different classifiers classify the data
 - q Go with the majority
 - q Intuition says that as the number of classifiers increases, the dependability of the majority vote increases

- n The corresponding algorithms were called Boosting by majority
 - q A (weighted) majority vote taken over all the classifiers
 - q How do we compute weights for the classifiers?
 - q How do we actually train the classifiers

ADA Boost: Adaptive algorithm for learning the weights

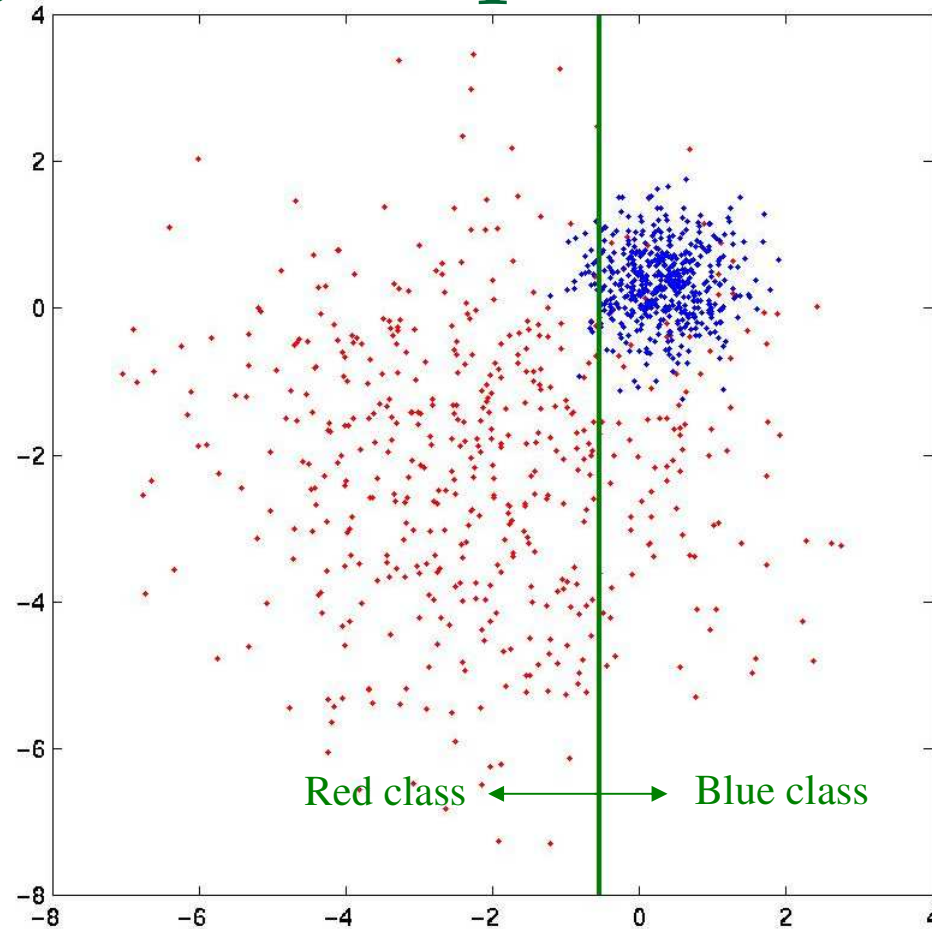
- n ADA Boost: Not named of ADA Lovelace
- n An *adaptive* algorithm that learns the weights of each classifier sequentially
 - q Learning adapts to the current accuracy
- n Iteratively:
 - q Train a simple classifier from training data
 - n It will make errors even on training data
 - n Train a new classifier that focuses on the training data points that have been misclassified

Boosting: An Example



- n Red dots represent training data from Red class
- n Blue dots represent training data from Blue class

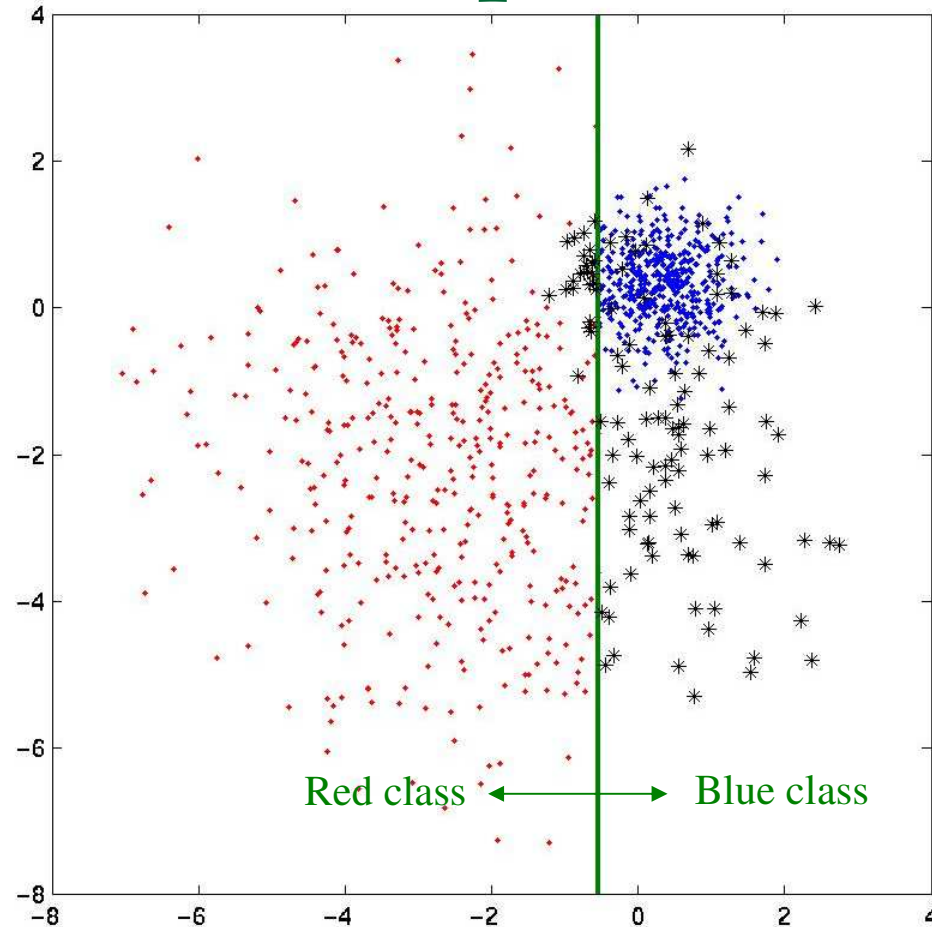
Boosting: An Example



n Very simple weak learner

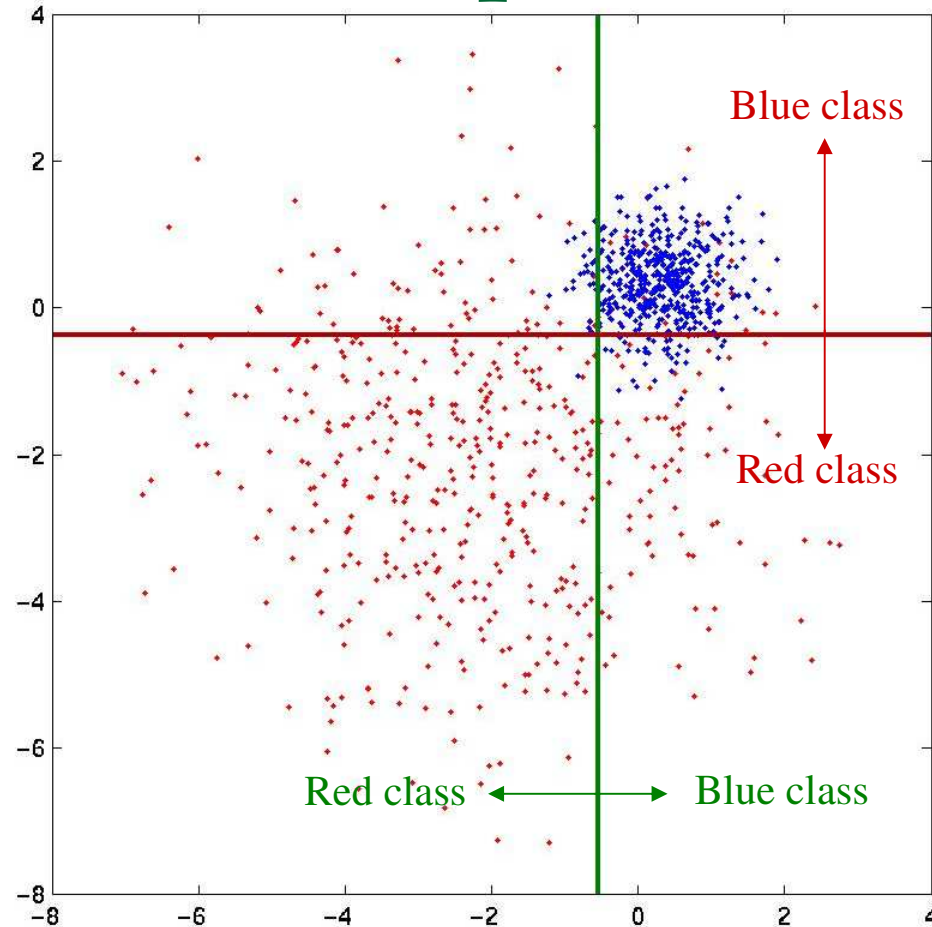
q A line that is parallel to one of the two axes

Boosting: An Example



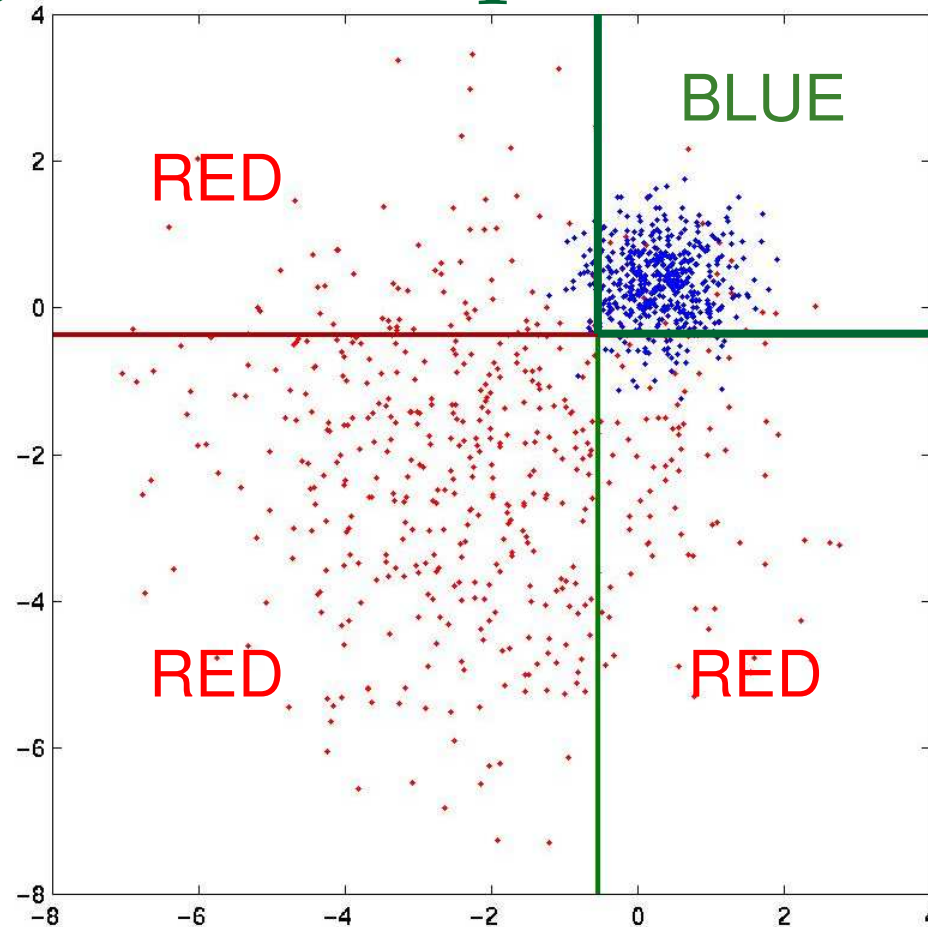
- n First weak learner makes many mistakes
 - q Errors coloured black

Boosting: An Example



- n Second weak learner focuses on errors made by first learner

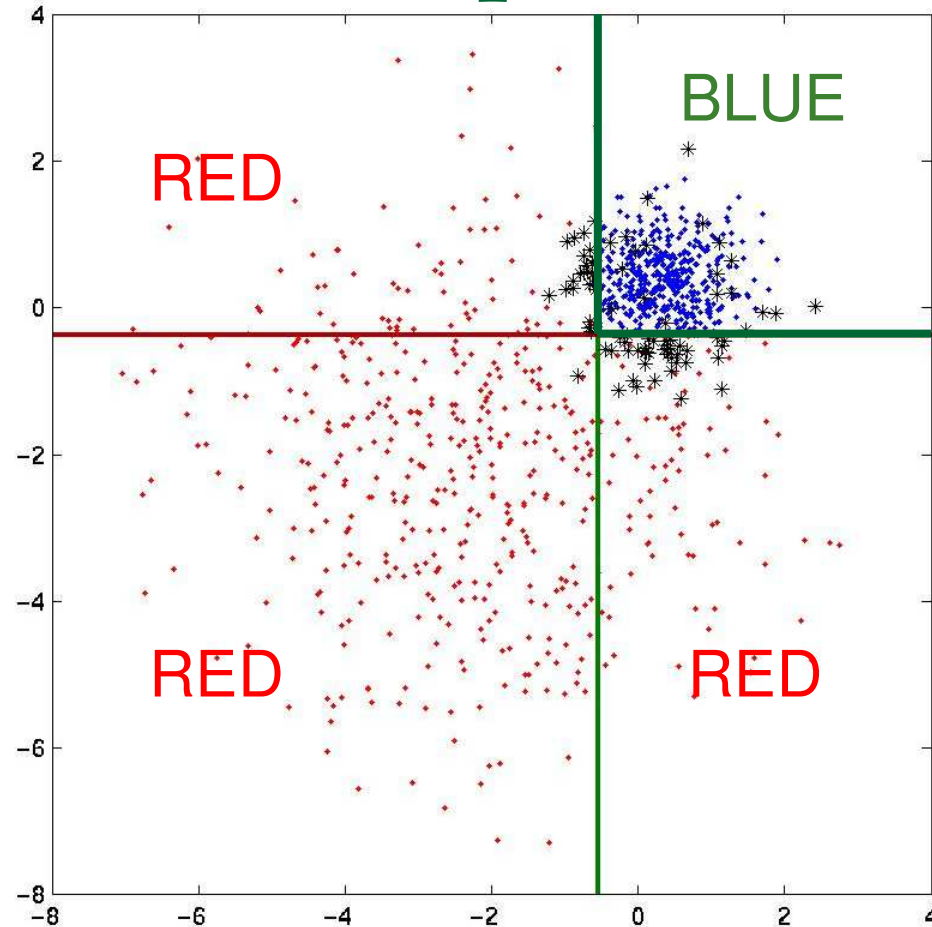
Boosting: An Example



§ **Second strong learner:** weighted combination of first and second weak learners

- Decision boundary shown by black lines

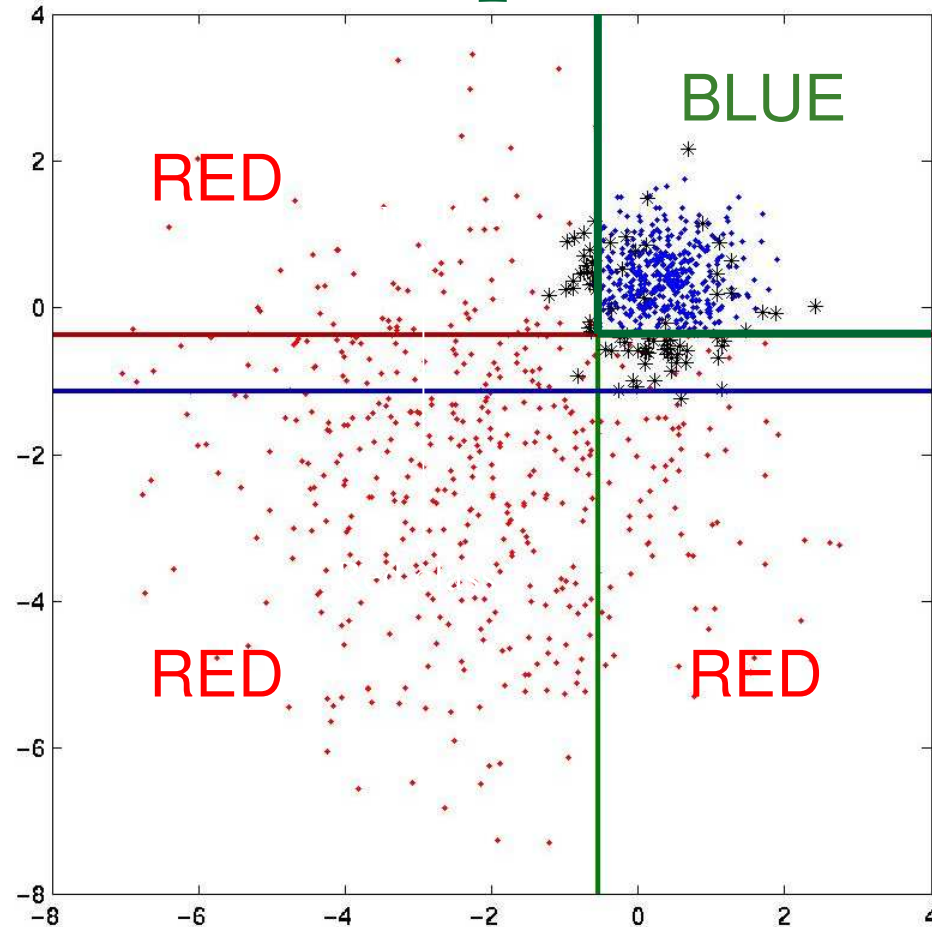
Boosting: An Example



n The second strong learner also makes mistakes

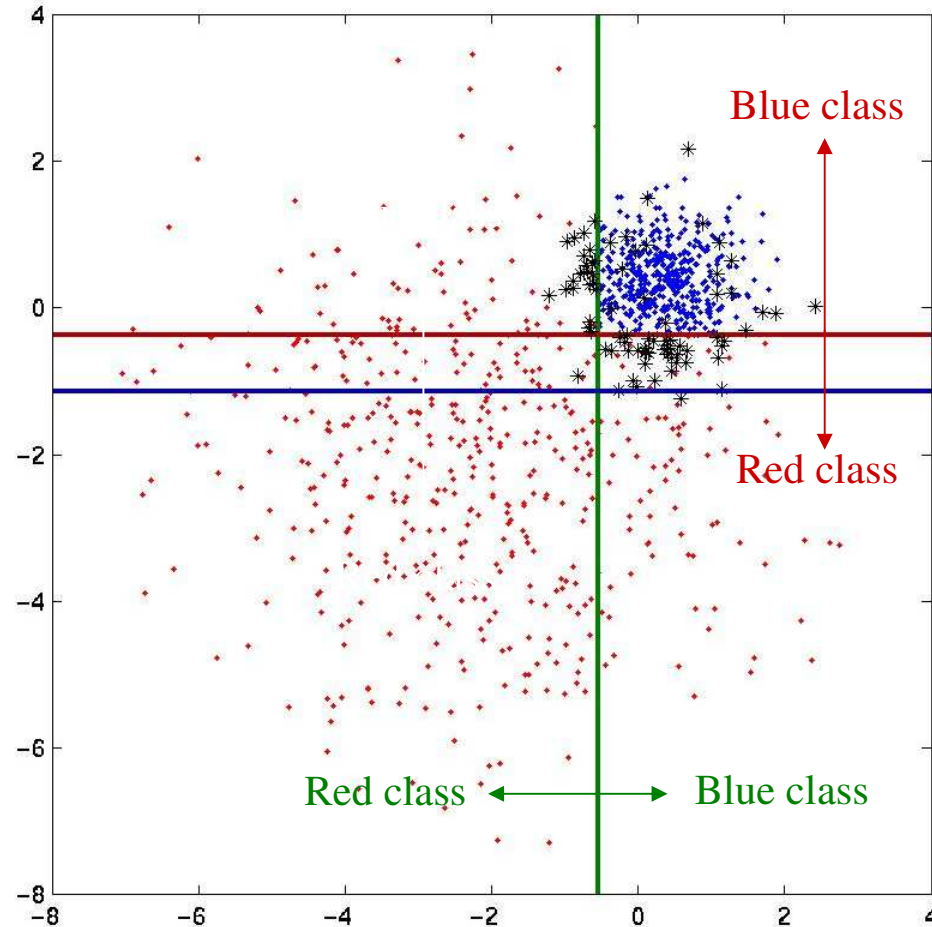
q Errors colored black

Boosting: An Example



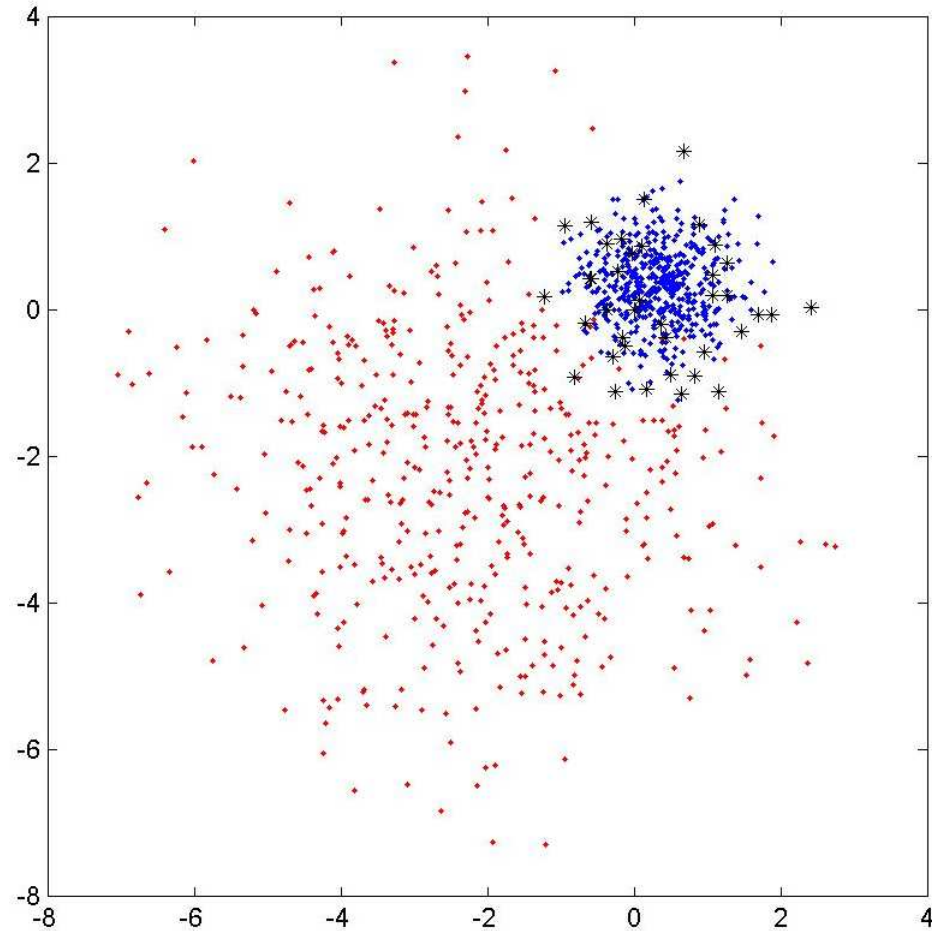
- n Third weak learner concentrates on errors made by second strong learner

Boosting: An Example



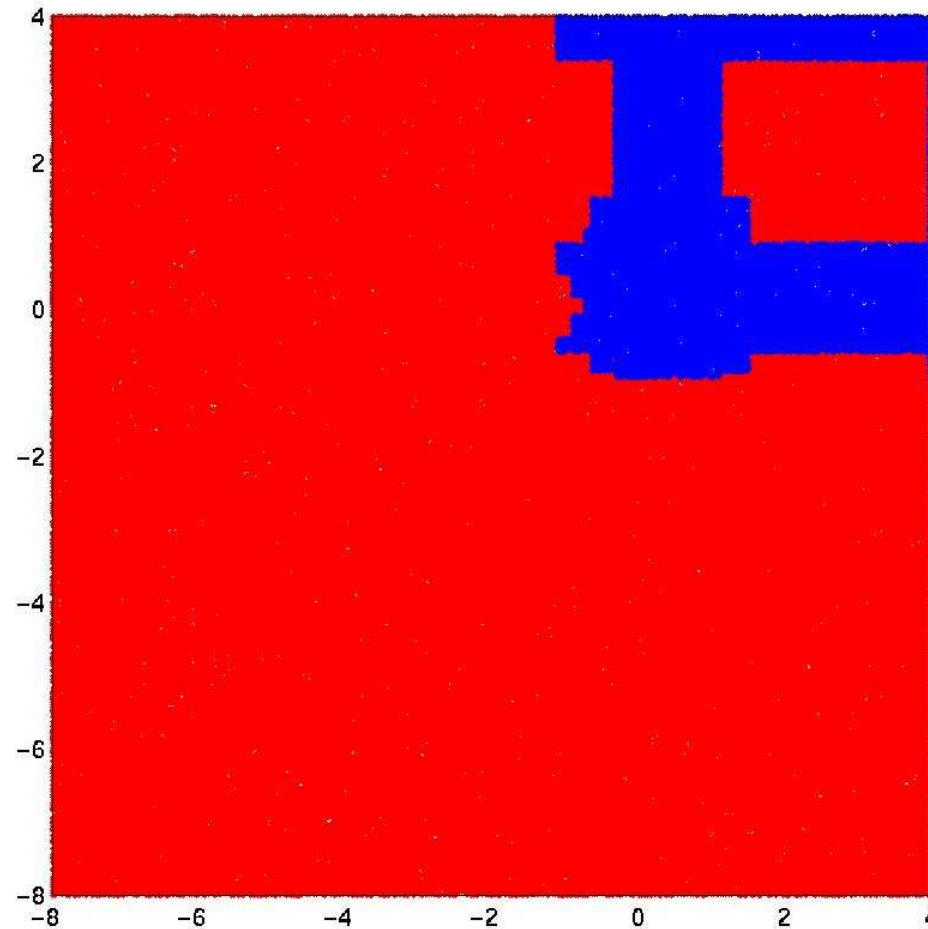
- § Third weak learner concentrates on errors made by combination of previous weak learners
- § Continue adding weak learners until....

Boosting: An Example



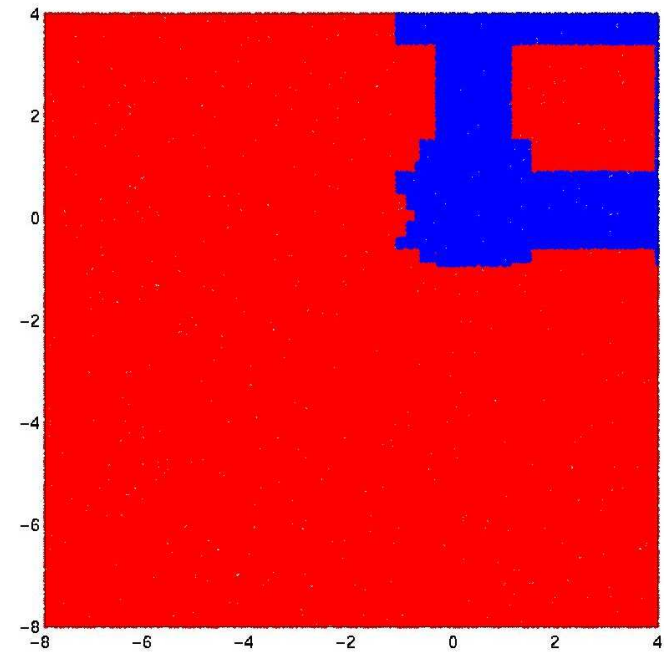
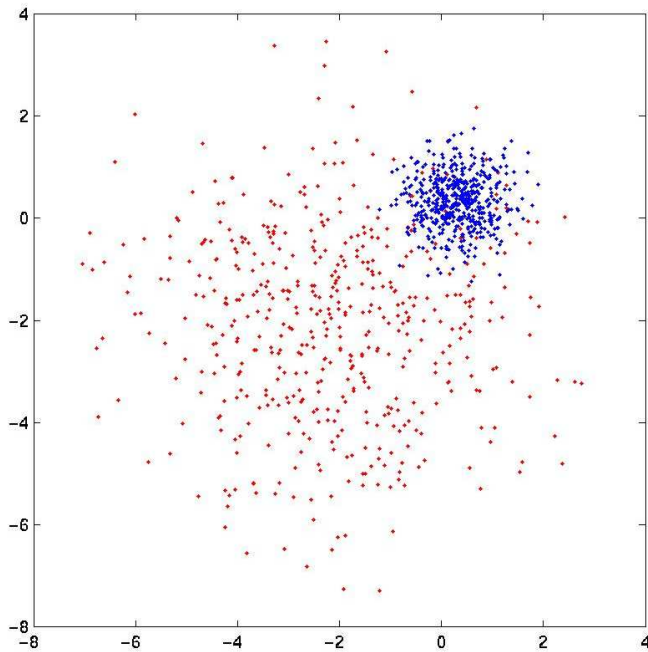
n Voila! Final strong learner: very few errors on the training data

Boosting: An Example



- n The final **strong** learner has learnt a complicated decision boundary

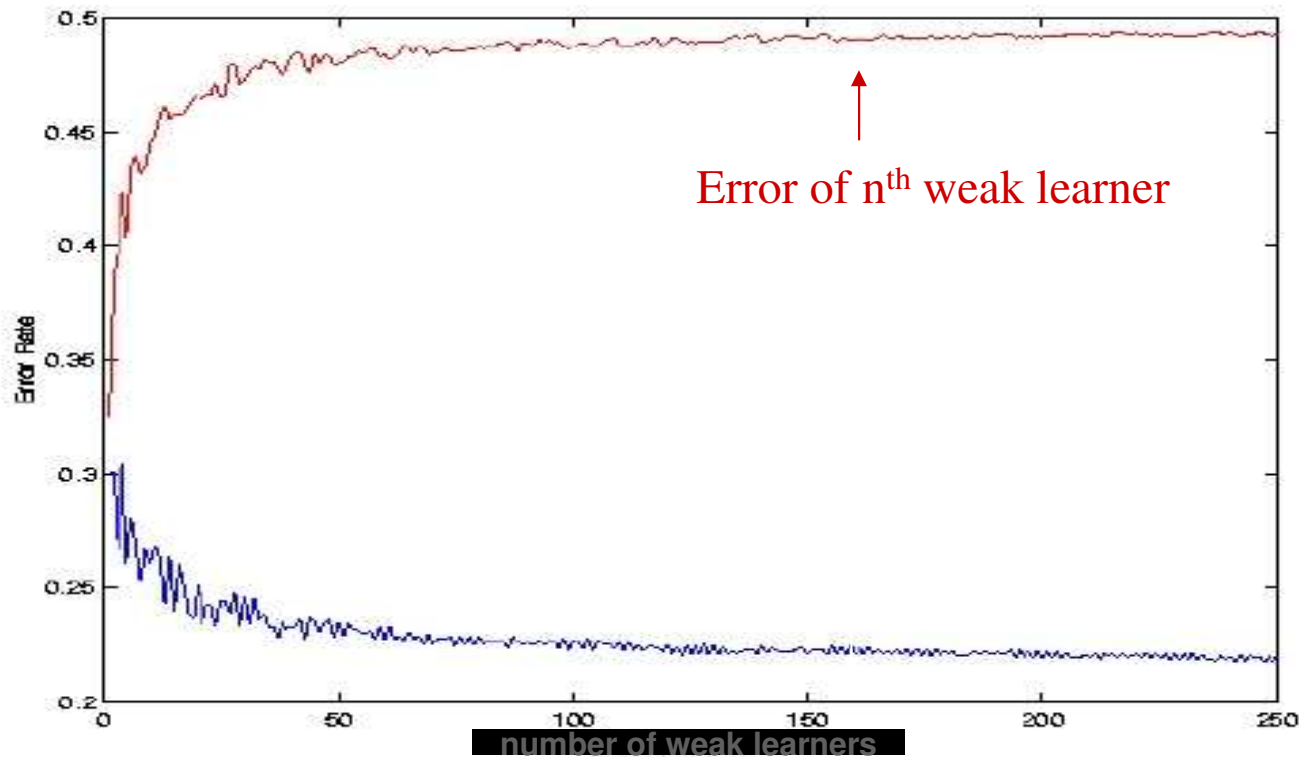
Boosting: An Example



- n The final **strong** learner has learnt a complicated decision boundary
- n Decision boundaries in areas with low density of training points assumed inconsequential

Overall Learning Pattern

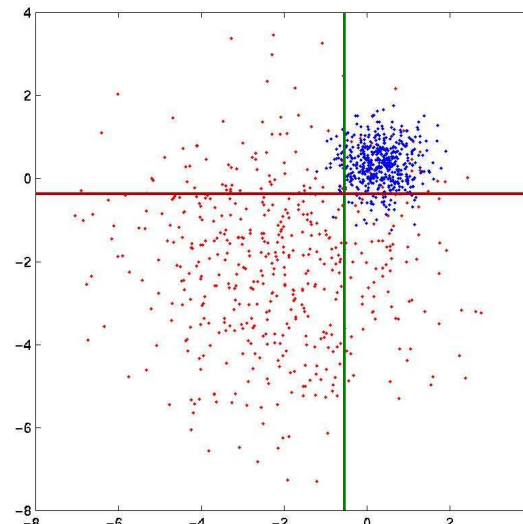
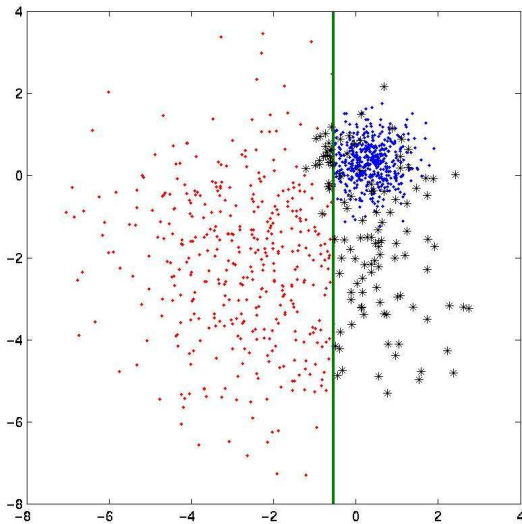
- § Strong learner increasingly accurate with increasing number of weak learners
- § Residual errors increasingly difficult to correct
 - Additional weak learners less and less effective



ADABOOST

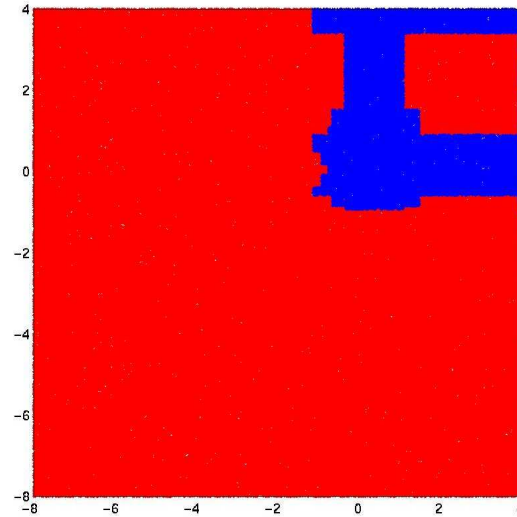
- n Cannot just add new classifiers that work well only on the previously misclassified data
- n Problem: The new classifier will make errors on the points that the **earlier** classifiers got right
 - q Not good
 - q On test data we have no way of knowing which points were correctly classified by the first classifier
- n Solution: Weight the data when training the second classifier
 - q Use all the data but assign them weights
 - n Data that are already correctly classified have less weight
 - n Data that are currently incorrectly classified have more weight

ADA Boost



- n The red and blue points (correctly classified) will have a weight $a < 1$
- n Black points (incorrectly classified) will have a weight $b (= 1/a) > 1$
- n To compute the optimal second classifier, we minimize the total weighted error
 - q Each data point contributes a or b to the total count of correctly and incorrectly classified points
 - n E.g. if one of the red points is misclassified by the new classifier, the total error of the new classifier goes up by a

ADA Boost



- n Each new classifier modifies the weights of the data points based on the accuracy of the *current* classifier
- n ***The final classifier too is a weighted combination of all component classifiers***

Formalizing the Boosting Concept

- n Given a set of instances $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$
 - q x_i is the set of attributes of the i^{th} instance
 - q y_i is the class for the i^{th} instance
 - n y_i can be 1 or -1 (binary classification only)

- n Given a set of classifiers h_1, h_2, \dots, h_T
 - q h_i classifies an instance with attributes x as $h_i(x)$
 - q $h_i(x)$ is either -1 or +1 (for a binary classifier)

 - q $y^*h(x)$ is 1 for all correctly classified points and -1 for incorrectly classified points

- n Devise a function $f(h_1(x), h_2(x), \dots, h_T(x))$ such that classification based on $f()$ is superior to classification by any $h_i(x)$
 - q The function is succinctly represented as $f(x)$

The Boosting Concept

n A simple combiner function: Voting

q $f(x) = \sum_j s_j h_j(x)$

q Classifier $H(x) = \text{sign}(f(x)) = \text{sign}(\sum_j s_j h_j(x))$

q Simple majority classifier

n A simple voting scheme

n A better combiner function: Boosting

q $f(x) = \sum_j s_j a_j h_j(x)$

n Can be any real number

q Classifier $H(x) = \text{sign}(f(x)) = \text{sign}(\sum_j s_j a_j h_j(x))$

q A weighted majority classifier

n The weight a_j for any $h_j(x)$ is a measure of our trust in $h_j(x)$

Adaptive Boosting

n As before:

q y is either -1 or +1

q $H(x)$ is +1 or -1

q If the instance is correctly classified, both y and $H(x)$ will have the same sign

n The product $y.H(x)$ is 1

n For incorrectly classified instances the product is -1

n Define the error for x : $\frac{1}{2}(1 - yH(x))$

q For a correctly classified instance, this is 0

q For an incorrectly classified instance, this is 1

The ADABOOST Algorithm

- n Given: a set $(x_1, y_1), \dots, (x_N, y_N)$ of training instances
 - q x_i is the set of attributes for the i^{th} instance
 - q y_i is the class for the i^{th} instance and can be either +1 or -1

The ADABOOST Algorithm

- n Initialize $D_1(x_j) = 1/N$
- n For $t = 1, \dots, T$
 - q Train a weak classifier h_t using distribution D_t
 - q Compute total error on training data
 - n $e_t = \text{Sum} \{1/2 (1 - y_i h_t(x_i))\}$
 - q Set $a_t = 1/2 \ln ((1 - e_t) / e_t)$
 - q For $i = 1 \dots N$
 - n set $D_{t+1}(x_i) = D_t(x_i) \exp(- a_t y_i h_t(x_i))$
 - q Normalize D_{t+1} to make it a distribution
- n The final classifier is
 - q $H(x) = \text{sign}(\sum_t a_t h_t(x))$

The ADABOOST Algorithm

Initialize $D_1(x_i) = 1/N$

n For $t = 1, \dots, T$

q Train a weak classifier h_t using distribution D_t

q Compute total error on training data

n $e_t = \text{Sum} \{1/2 (1 - y_i h_t(x_i))\}$

q Set $a_t = 1/2 \ln ((1 - e_t) / e_t)$

q For $i = 1 \dots N$

n set $D_{t+1}(x_i) = D_t(x_i) \exp(- a_t y_i h_t(x_i))$

q Normalize D_{t+1} to make it a distribution

n The final classifier is

q $H(x) = \text{sign}(\sum_t a_t h_t(x))$

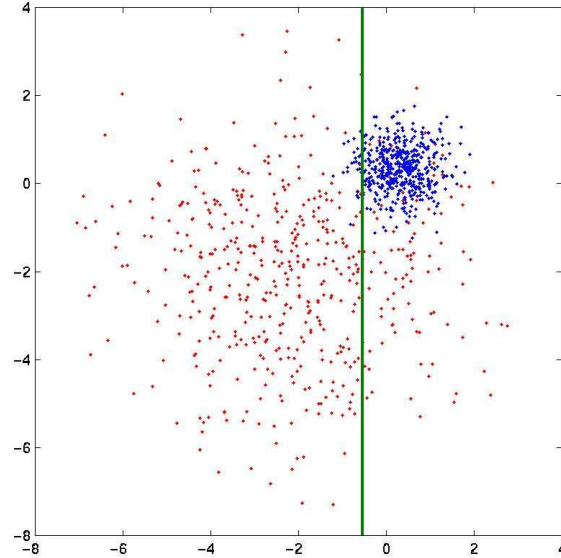
ADA Boost

- n Initialize $D_1(x_i) = 1/N$
- n Just a normalization: total weight of all instances is 1
 - q Makes the algorithm invariant to training data set size

The ADABOOST Algorithm

- n Initialize $D_1(x_j) = 1/N$
- n For $t = 1, \dots, T$
 - q Train a weak classifier h_t using distribution D_t
 - q Compute total error on training data
 - n $e_t = \text{Sum} \{1/2 (1 - y_i h_t(x_i))\}$
 - q Set $a_t = 1/2 \ln ((1 - e_t) / e_t)$
 - q For $i = 1 \dots N$
 - n set $D_{t+1}(x_i) = D_t(x_i) \exp(- a_t y_i h_t(x_i))$
 - q Normalize D_{t+1} to make it a distribution
- n The final classifier is
 - q $H(x) = \text{sign}(\sum_t a_t h_t(x))$

ADA Boost



- n Train a weak classifier h_t using distribution D_t
- n Simply train the simple classifier that that classifies data with error 50%
 - q Where each data x point contributes $D(x)$ towards the count of errors or correct classification
 - q Initially $D(x) = 1/N$ for all data
- n Better to actually train a *good* classifier

The ADABOOST Algorithm

n Initialize $D_1(x_j) = 1/N$

n For $t = 1, \dots, T$

q Train a weak classifier h_t using distribution D_t

q Compute total error on training data

$$n e_t = \text{Sum} \{1/2 (1 - y_i h_t(x_i))\}$$

q Set $a_t = 1/2 \ln ((1 - e_t) / e_t)$

q For $i = 1 \dots N$

$$n \text{ set } D_{t+1}(x_i) = D_t(x_i) \exp(- a_t y_i h_t(x_i))$$

q Normalize D_{t+1} to make it a distribution

n The final classifier is

$$q H(x) = \text{sign}(\sum_t a_t h_t(x))$$

ADA Boost

n Compute total error on training data

q $e_t = \text{Sum} \{ \frac{1}{2} (1 - y_i h_t(x_i)) \}$

n For each data point x , $\frac{1}{2}(1-y \cdot h(x)) = 0$ for correct classification, 1 for error

n e_t is simply the sum of the weights $D(x)$ for all points that are misclassified by the latest classifier $h_t(x)$

q Will lie between 0 and 1

$$e_t = \sum_{x \text{ such that } x \text{ is misclassified by } h_t(x)} D(x)$$

The ADABOOST Algorithm

n Initialize $D_1(x_j) = 1/N$

n For $t = 1, \dots, T$

q Train a weak classifier h_t using distribution D_t

q Compute total error on training data

$$n e_t = \text{Sum} \{1/2 (1 - y_i h_t(x_i))\}$$

q Set $a_t = 1/2 \ln ((1 - e_t) / e_t)$

q For $i = 1 \dots N$

$$n \text{ set } D_{t+1}(x_i) = D_t(x_i) \exp(- a_t y_i h_t(x_i))$$

q Normalize D_{t+1} to make it a distribution

n The final classifier is

$$q H(x) = \text{sign}(\sum_t a_t h_t(x))$$

Classifier Weight

n Set $a_t = \frac{1}{2} \ln \left(\frac{1-e_t}{e_t} \right)$

n The a_t for any classifier is always positive

n The weight for the t^{th} classifier is a function of its error

q The poorer the classifier is, the closer a_t is to 0

q If the error of the classifier is exactly 0.5, a_t is 0.

n We don't trust such classifiers at all \mathcal{J}

q If the error approaches 0, a_t becomes high

n We trust these classifiers

The ADABOOST Algorithm

n Initialize $D_1(x_i) = 1/N$

n For $t = 1, \dots, T$

q Train a weak classifier h_t using distribution D_t

q Compute total error on training data

n $e_t = \text{Average} \{1/2 (1 - y_i h_t(x_i))\}$

q Set $a_t = 1/2 \ln ((1 - e_t) / e_t)$

q For $i = 1 \dots N$

n set $D_{t+1}(x_i) = D_t(x_i) \exp(- a_t y_i h_t(x_i))$

q Normalize D_{t+1} to make it a distribution

n The final classifier is

q $H(x) = \text{sign}(\sum_t a_t h_t(x))$

ADA Boost

- n For $i = 1 \dots N$
 - q set $D_{t+1}(x_i) = D_t(x_i) \exp(-a_t y_i h_t(x_i))$
- n Normalize D_{t+1} to make it a distribution

- n Readjusting the weights of all training instances
 - q If the instance is correctly classified, multiply its weight by $b (= \exp(-a_t)) < 1$
 - q If it is *misclassified*, multiply its weight by $b (= \exp(a_t)) > 1$

- n Renormalize, so they all sum to 1

- q
$$D_{renormalized}(x) = D(x) / \sum_{x'} D(x')$$

The ADABOOST Algorithm

- n Initialize $D_1(x_j) = 1/N$
- n For $t = 1, \dots, T$
 - q Train a weak classifier h_t using distribution D_t
 - q Compute total error on training data
 - n $e_t = \text{Average} \{1/2 (1 - y_i h_t(x_i))\}$
 - q Set $a_t = 1/2 \ln ((1 - e_t) / e_t)$
 - q For $i = 1 \dots N$
 - n set $D_{t+1}(x_i) = D_t(x_i) \exp(- a_t y_i h_t(x_i))$
 - q Normalize D_{t+1} to make it a distribution
- n The final classifier is
 - q $H(x) = \text{sign}(\sum_t a_t h_t(x))$

ADA Boost

n The final classifier is

q $H(x) = \text{sign}(S_t a_t h_t(x))$

n The output is 1 if the total weight of all weak learners that classify x as 1 is greater than the total weight of all weak learners that classify it as -1

Next Class

- n Fernando De La Torre
- n We will continue with Viola Jones after a few classes