# Lecture 8:
# ML for image classification, then Face and pedestrian detection

## COS 429: Computer Vision

**PRINCETON UNIVERSITY**

# Building on last lecture:
# ML for image classification

# Classification

Start with simplest example: binary classification



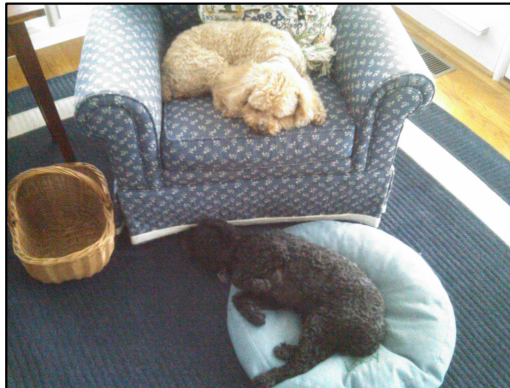| $x_1$ |
| $x_2$ |
| . . . |
| $x_N$ |

→ Cat or not cat?

Actually: a feature vector
representing the image

# Easiest Form of Classification

Just **memorize** (as in a Python dictionary)
Consider cat/dog/hippo classification.



If this:
cat.

If this:
dog.

If this:
hippo.

# Easiest Form of Classification

## Where does this go wrong?



Rule: if this, then cat



Hmmm. Not quite the same.
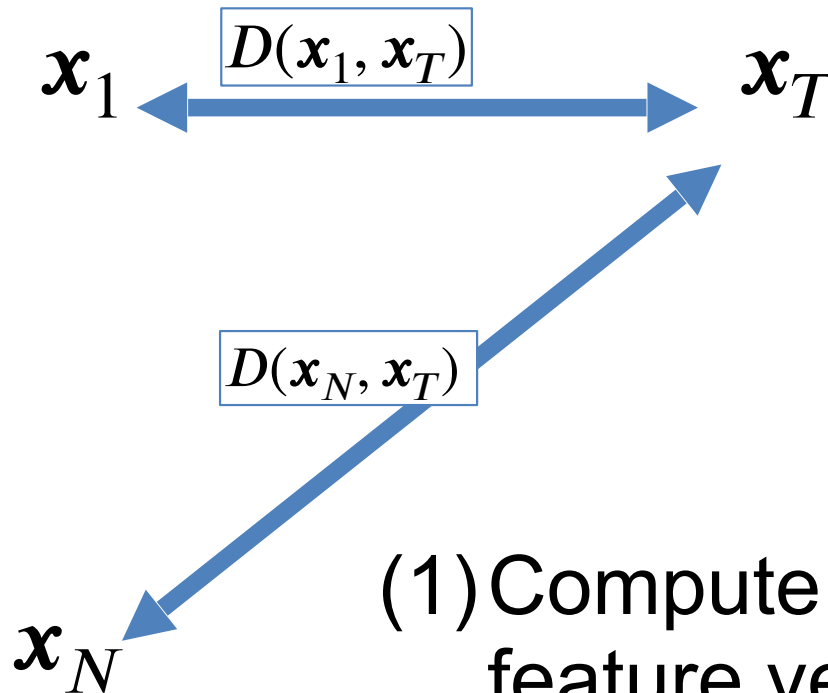
# Easiest Form of Classification

Known Images Labels

Test Image



Cat

$\boldsymbol{x}_1$ ←$D(\boldsymbol{x}_1, \boldsymbol{x}_T)$→ $\boldsymbol{x}_T$ Cat!

. . .

$D(\boldsymbol{x}_N, \boldsymbol{x}_T)$

$\boldsymbol{x}_N$

Dog

(1) Compute distance between feature vectors
(2) Find nearest
(3) Use label.

# Nearest Neighbor
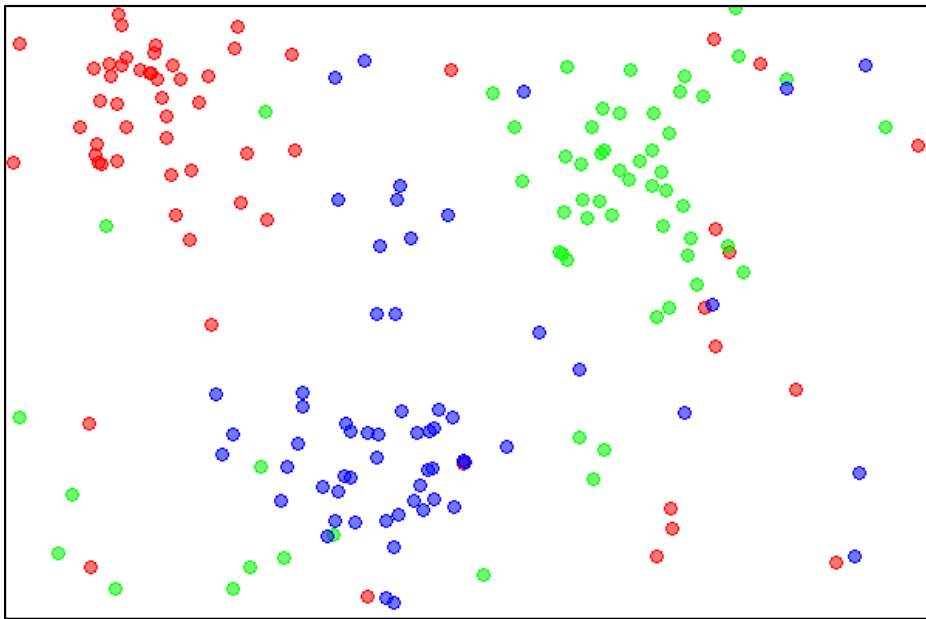
## "Algorithm"

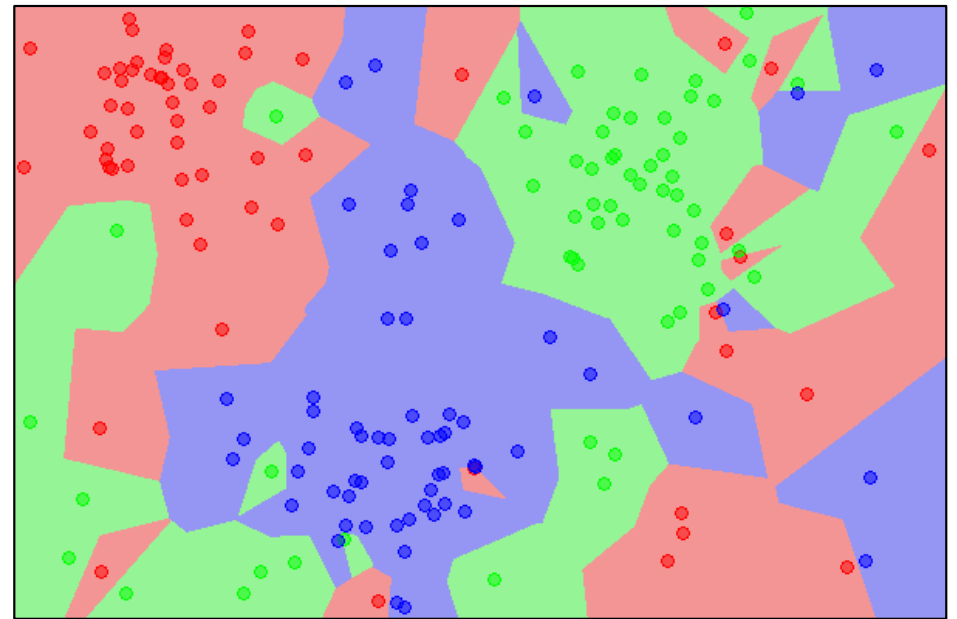Training ($\mathbf{x}_i, y_i$):      Memorize training set

Inference (x):

```
bestDist, prediction = Inf, None
for i in range(N):
        if dist(x_i,x) < bestDist:
                bestDist = dist(x_i,x)
                prediction = y_i
```

# Nearest Neighbor

## 2D Datapoints (colors = labels)



## 2D Predictions (colors = labels)



Diagram Credit: Wikipedia

Credit: D. Fouhey
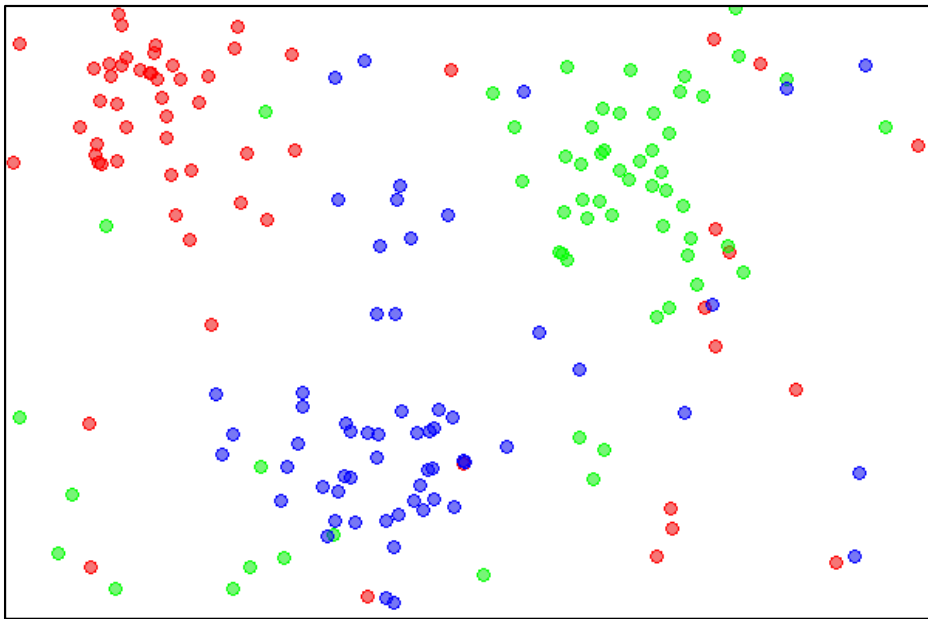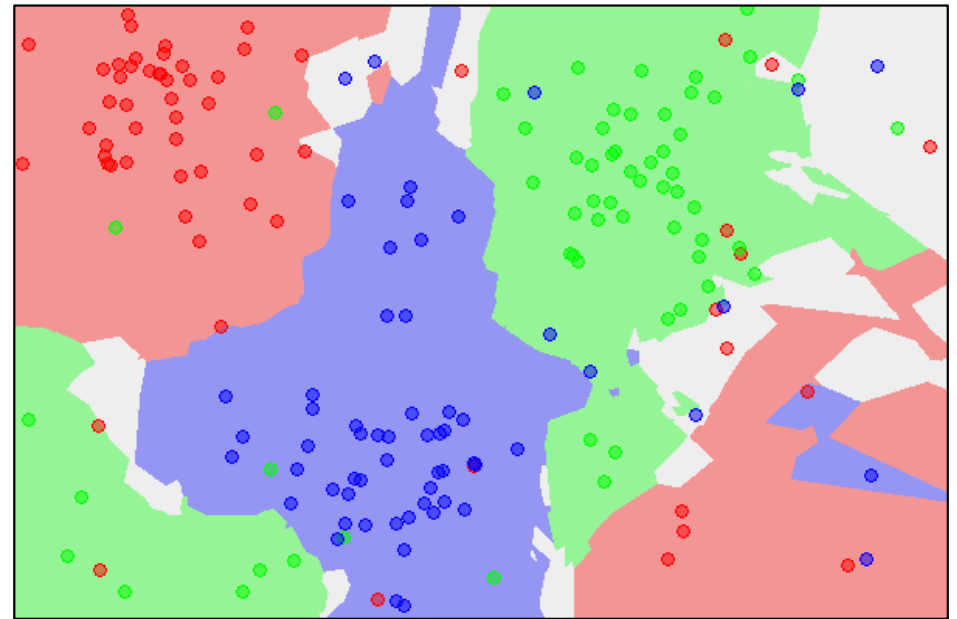
# K-Nearest Neighbors

Take top K-closest points, vote

2D Datapoints
(colors = labels)

2D Predictions
(colors = labels)



Diagram Credit: Wikipedia

Credit: D. Fouhey

# K-Nearest Neighbors

- No learning going on but usually effective

- Same algorithm for every task

- As number of datapoints → ∞, error rate is guaranteed to be at most 2x worse than optimal you could do on data

# Alternative: recall last lecture



$x_1$

$x_2$

$...$

$x_N$

Cat or not cat?

# Classification by Least-Squares

Treat as regression: $x_i$ is image feature; $y_i$ is 1 if it's a cat, 0 if it's not a cat. Minimize least-squares loss.

Training (**X**,y): $$\underset{\boldsymbol{w}}{\mathrm{argmin}} \left\| \boldsymbol{y} - \boldsymbol{X}\boldsymbol{w} \right\|_2^2 + \lambda \left\| \boldsymbol{w} \right\|_2^2$$

Loss          Trade-off          Regularization

Inference (x): $$\boldsymbol{w}^T \boldsymbol{x} > t$$

Unprincipled in theory, but often effective in practice
The reverse (regression via discrete bins) is also common

Rifkin, Yeo, Poggio. *Regularized Least Squares Classification* (http://cbcl.mit.edu/publications/ps/rlsc.pdf). 2003
Redmon, Divvala, Girshick, Farhadi. *You Only Look Once: Unified, Real-Time Object Detection.* CVPR 2016.

# Solving regularized least squares

Objective:   $\underset{w}{\mathrm{argmin}} \left\| y - Xw \right\|_2^2 + \lambda \|w\|_2^2$

Take $\dfrac{\partial}{\partial w}$, set to **0**, solve

$$w^* = \left(X^T X + \lambda I\right)^{-1} X^T y$$

$X^TX+\lambda I$ is full-rank (and thus invertible) for λ>0

Called *lots of things:* regularized least-squares, Tikhonov regularization (after Andrey Tikhonov), ridge regression, Bayesian linear regression with a multivariate normal prior.
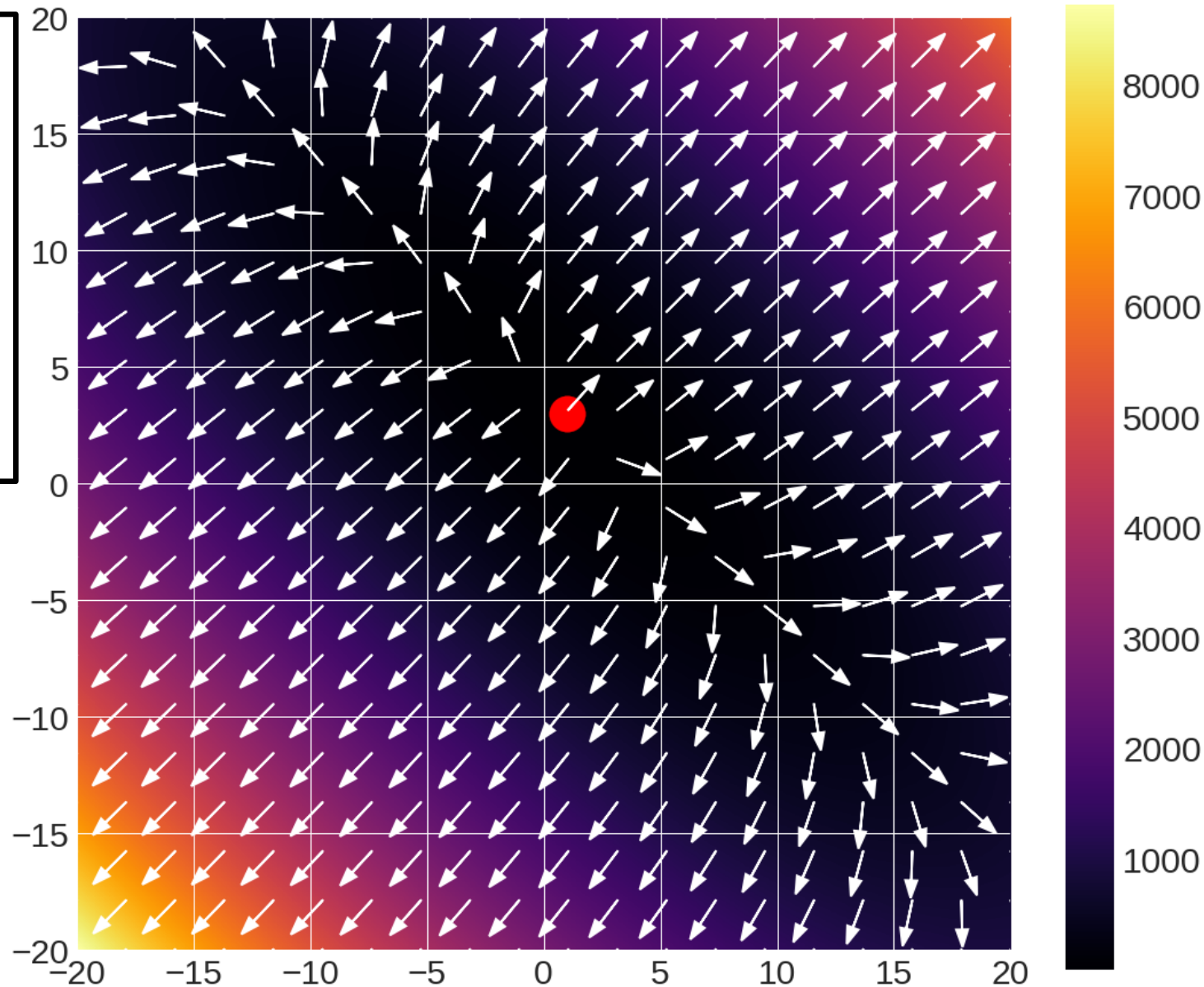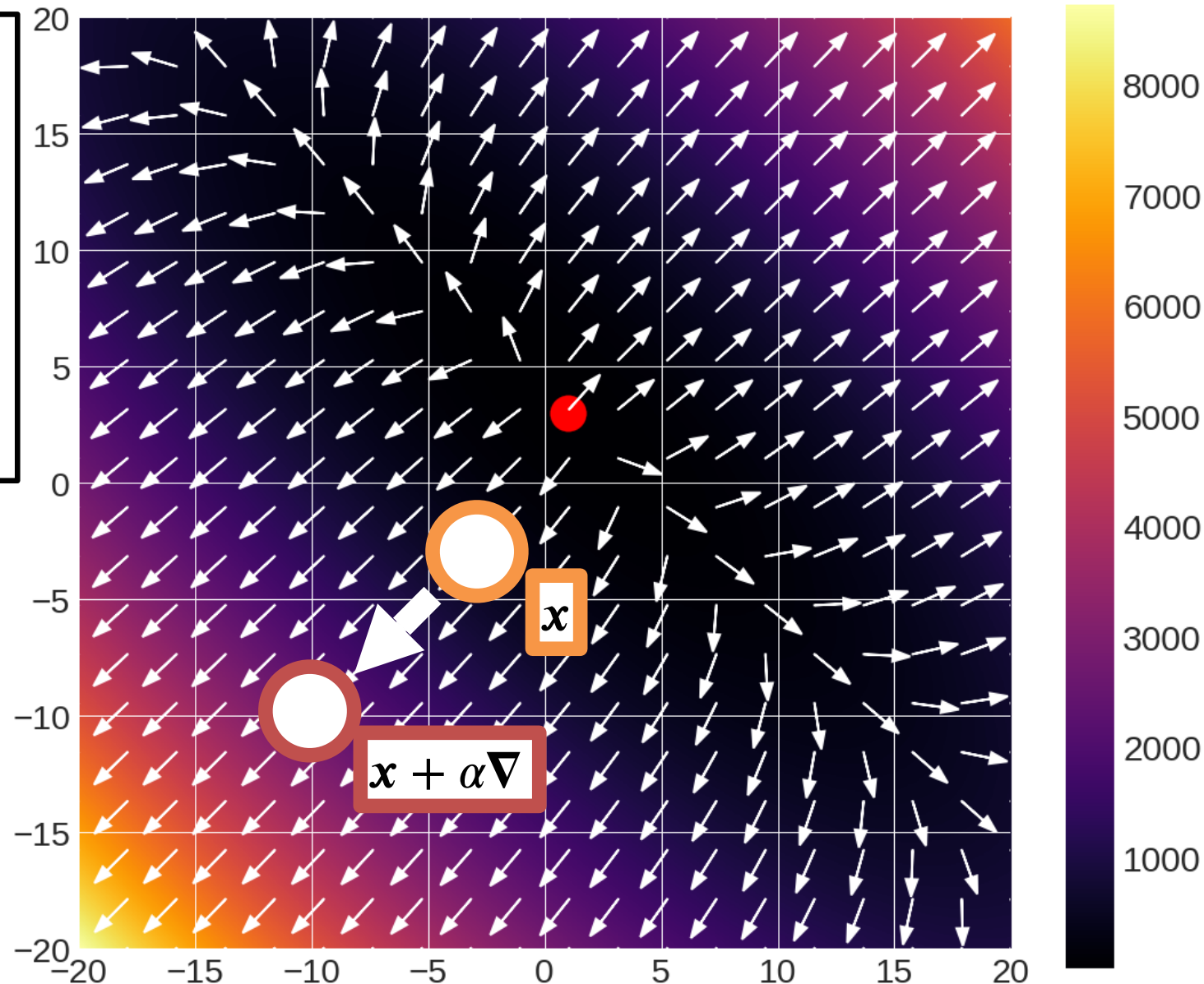
# But really, use the gradient

Arrows:
**gradient**

# Use the gradient

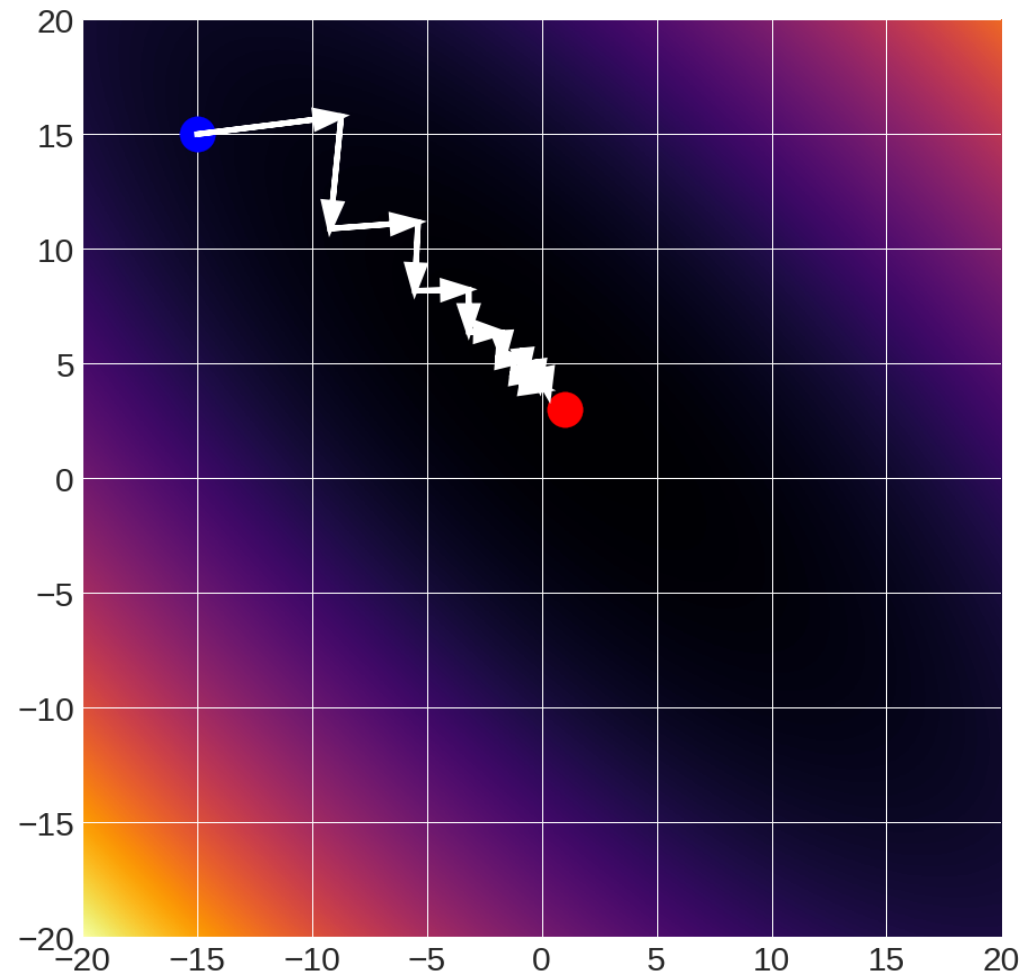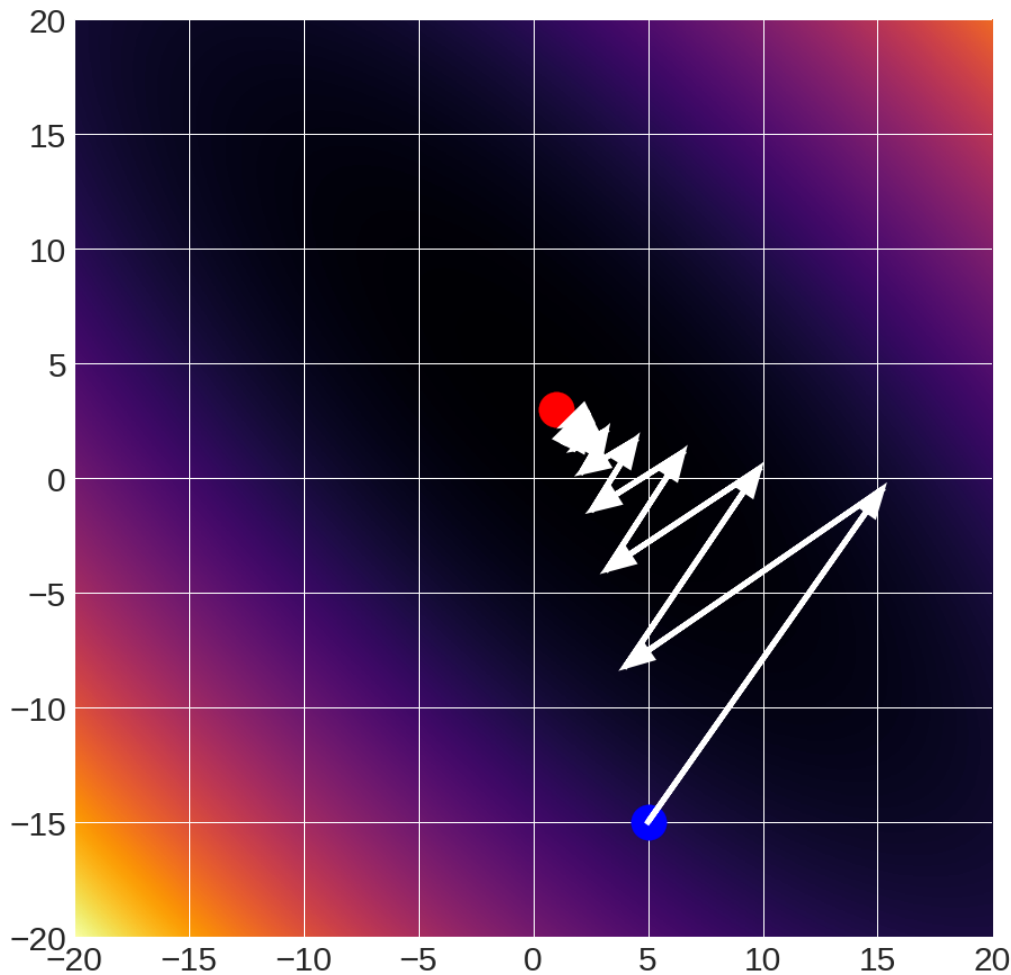Arrows: **gradient direction** (scaled to unit length)



Credit: D. Fouhey

# Use the gradient
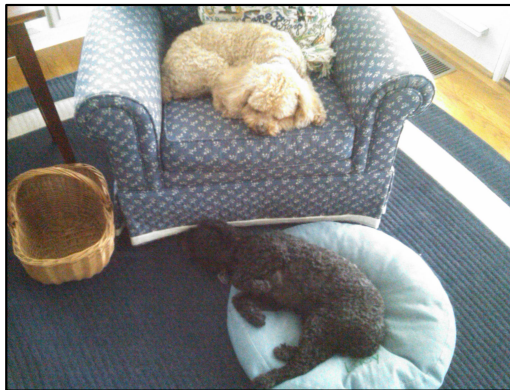
Arrows:
gradient
direction
(scaled to unit
length)



Credit: D. Fouhey

# Gradient descent

Given starting point (blue)
$$w_{i+1} = w_i + -9.8 \times 10^{-2} \times \text{gradient}$$



Credit: D. Fouhey

# Linear Models

## Example Setup: 3 classes



Model – one weight per class: $\boldsymbol{w}_0, \boldsymbol{w}_1, \boldsymbol{w}_2$

$\boldsymbol{w}_0^T \boldsymbol{x}$ big if cat

$\boldsymbol{w}_1^T \boldsymbol{x}$ big if dog

$\boldsymbol{w}_2^T \boldsymbol{x}$ big if hippo

Stack together: $\boldsymbol{W_{3xF}}$ where $\mathbf{x}$ is in $R^F$

# Linear Models



Cat weight vector | 0.2 | -0.5 | 0.1 | 2.0 | 1.1

Dog weight vector | 1.5 | 1.3 | 2.1 | 0.0 | 3.2

Hippo weight vector | 0.0 | 0.3 | 0.2 | -0.3 | -1.2

$W$

| 56 |
| 231 |
| 24 |
| 2 |
| 1 |

$x_i$

| -96.8 | Cat score |
| 437.9 | Dog score |
| 61.95 | Hippo score |

$Wx_i$

Weight matrix a collection of scoring functions, one per class

Prediction is vector where jth component is "score" for jth class.

Diagram by: Karpathy, Fei-Fei

# Geometric Intuition

## What does a linear classifier look like in 2D?



Diagram credit: Karpathy & Fei-Fei

# Visual Intuition

## CIFAR 10:
### 32x32x3 Images, 10 Classes



airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck

- Turn each image into feature by unrolling all pixels
- Fit 10 linear models

Slide credit: Karpathy & Fei-Fei

Decision rule is $\mathbf{w}^\top \mathbf{x}$. If $\mathbf{w}_i$ is big, then big values of $x_i$ are indicative of the class.

# Deer or Plane?

# Guess The Classifier

Decision rule is $\mathbf{w}^\top\mathbf{x}$. If $\mathbf{w}_i$ is big, then big values of $x_i$ are indicative of the class.
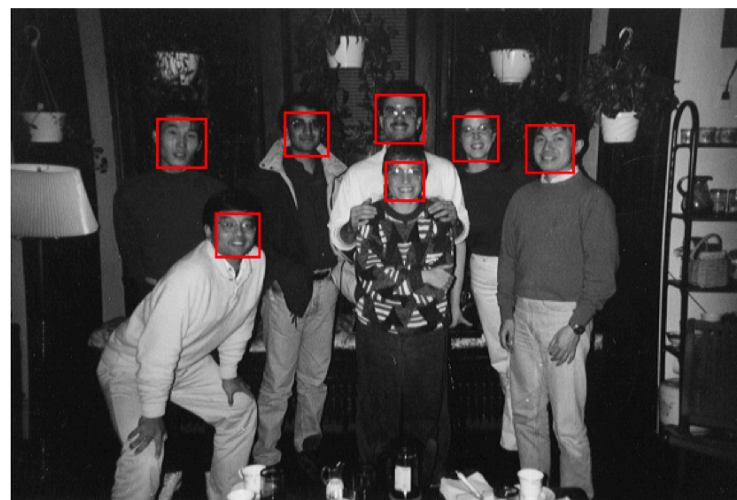
# Ship or Dog?



Diagram credit: Karpathy & Fei-Fei

# Interpreting a Linear Classifier

Decision rule is $\mathbf{w}^{\mathsf{T}}\mathbf{x}$. If $\mathbf{w}_i$ is big, then big values of $x_i$ are indicative of the class.



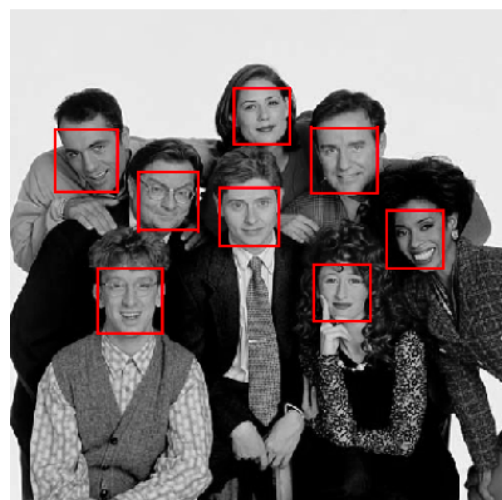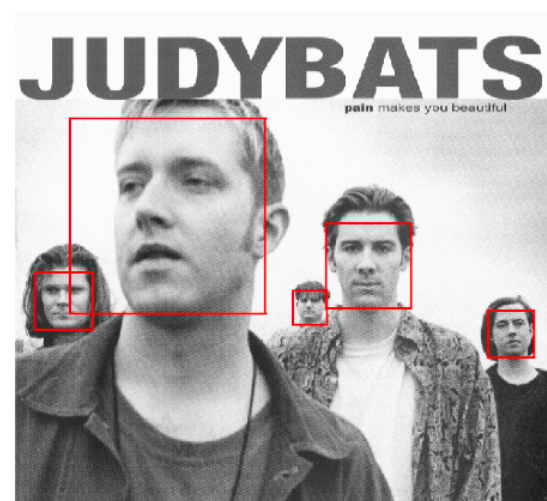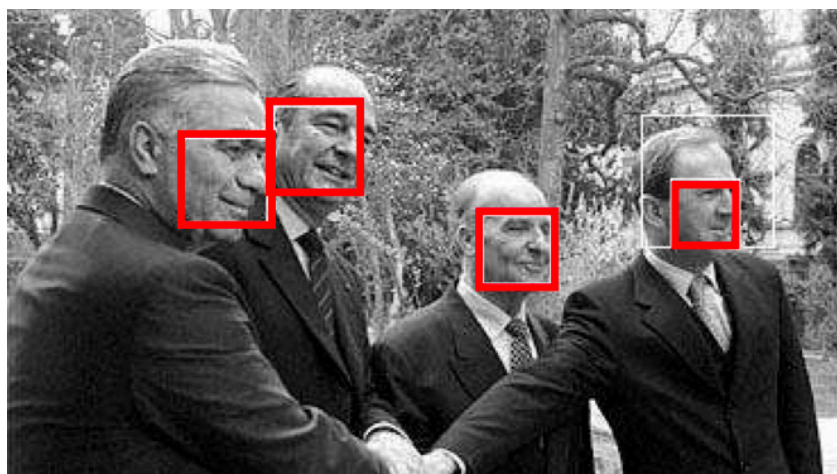Diagram credit: Karpathy & Fei-Fei

# Going forward

- Consider several computer vision tasks:
  - Today: face detection, pedestrian detection
  - Lectures 9,10: image classification, generic object detection
  - Lectures 11,12: segmentation, texture
- For each one:
  - Discuss **use cases, data, evaluation**
  - Key method: **representation, model, learning**
  - Results: **quantitative, qualitative**
- Will revisit many of these tasks in Module 4 on deep learning

# Midterm

- Everything through lecture 11 on Thursday, Oct 17th

- Will need to be able to articulate the connection between the design decisions (image representation, evaluation criteria, modeling architecture) and the outcome (expected sources of errors, comparative performance)
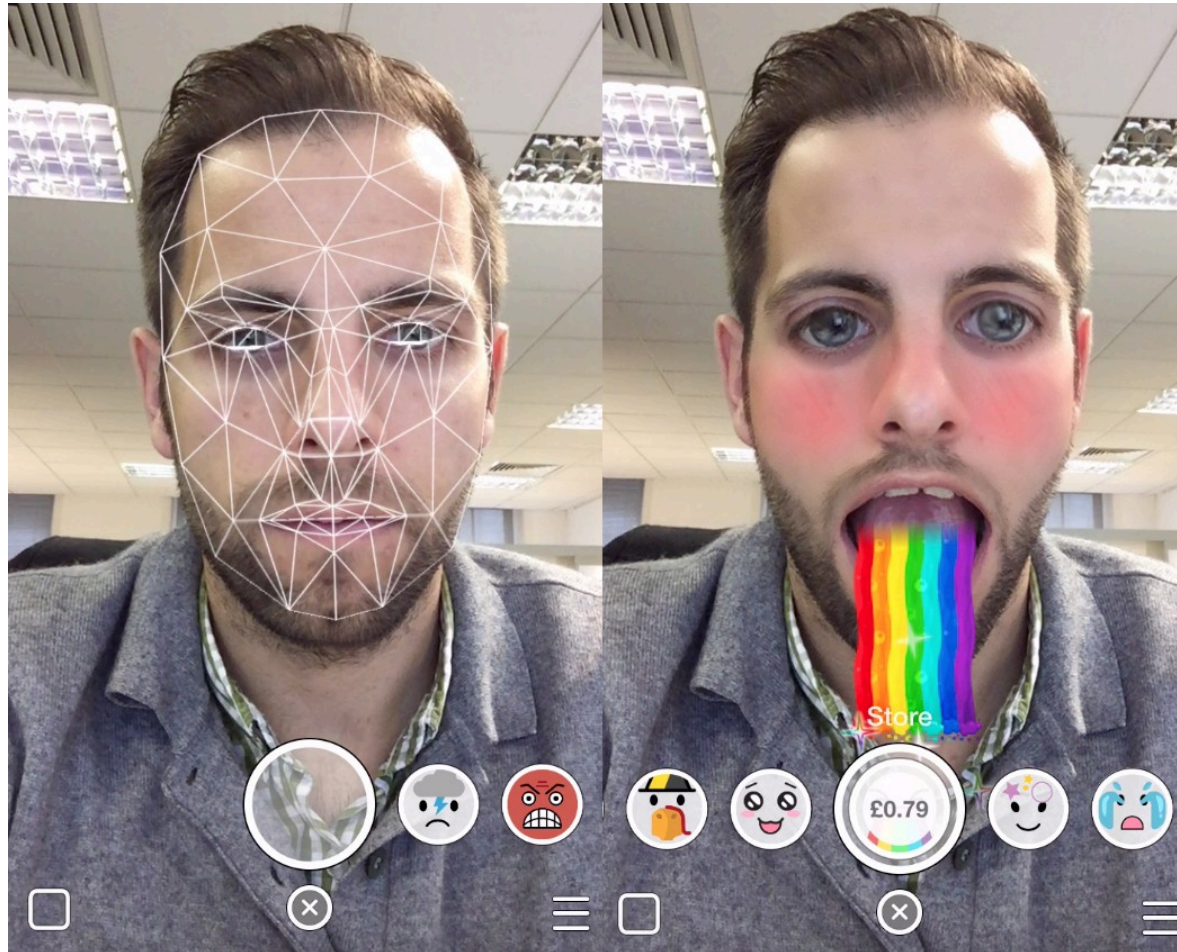
# Face detection

# Face detection in cameras
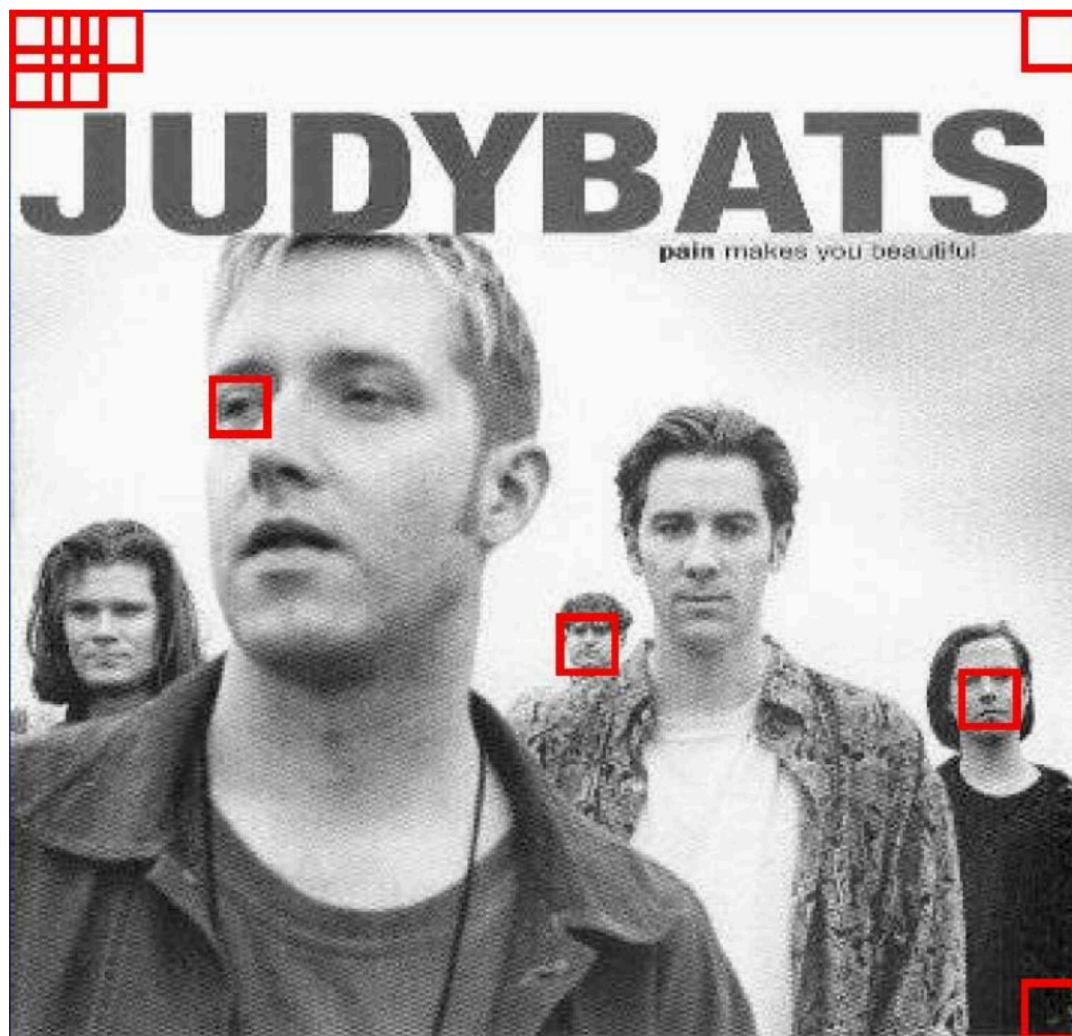


FinePix S6000fd, by Fujifilm, 2006

Viola & Jones, 2001

# E.g., Snapchat



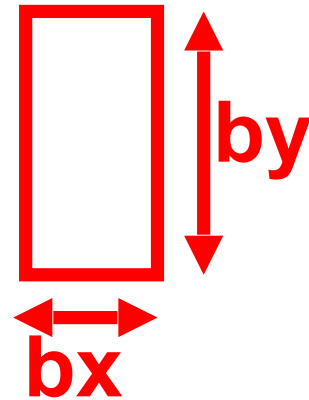https://medium.com/@anidaro/how-snapchats-filters-work-86973c3e2e9f

# Search over scale and space

# How many sliding window boxes are there?

Given a HxW image and a "template" of size by, bx.

**Q. How many sub-boxes are there of size (by,bx)?**

A. (H-by)*(W-bx)

This is before considering adding:
- *scales* (by*s,bx*s)
- *aspect ratios* (by*sy,bx*sx)
  - (although not for face detection)

# Challenges of face detection

- Sliding window detector must evaluate tens of thousands of location/scale combinations

- Faces are rare:  0–10 per image
  - A megapixel image has ~$10^6$ pixels and a comparable number of candidate face locations
  - For computational efficiency, we should try to spend as little time as possible on the non-face windows
  - To avoid having a false positive in every image, our false positive rate has to be less than $10^{-6}$

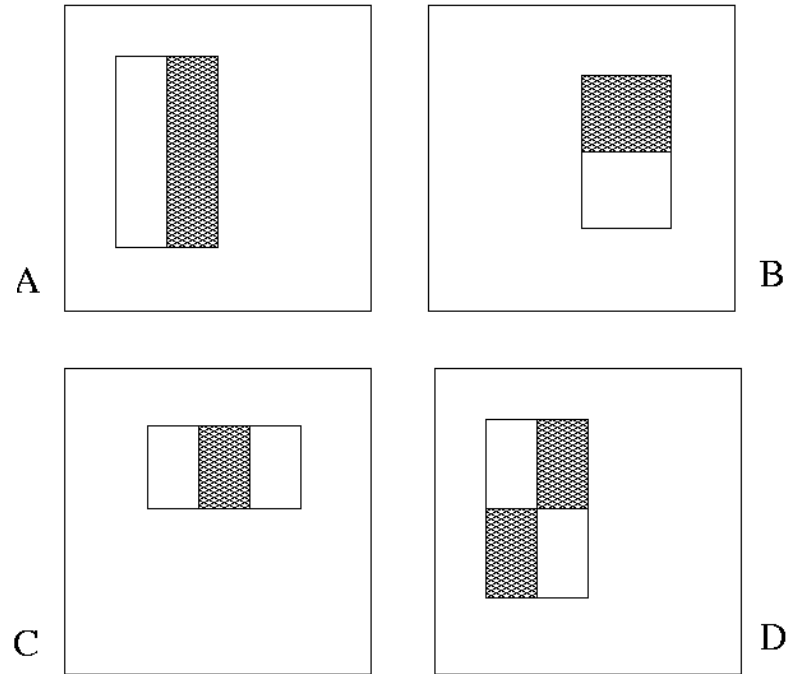# The Viola/Jones Face Detector

- A seminal approach to real-time object detection

- Training is slow, but detection is very fast

- Key ideas

  - *Integral images* for fast feature evaluation

  - *Boosting* for feature selection

  - *Attentional cascade* for fast rejection of non-face windows

P. Viola and M. Jones. *Rapid object detection using a boosted cascade of simple features.* CVPR 2001.

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

Source: S. Lazebnik
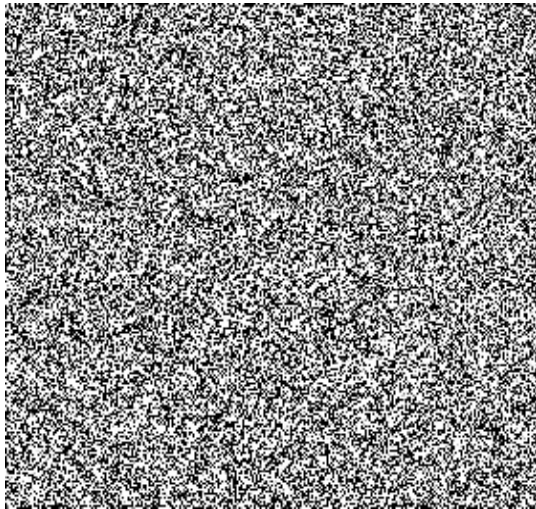
# Image Features

"Rectangle filters"



*Value =*

*∑ (pixels in white area) –*
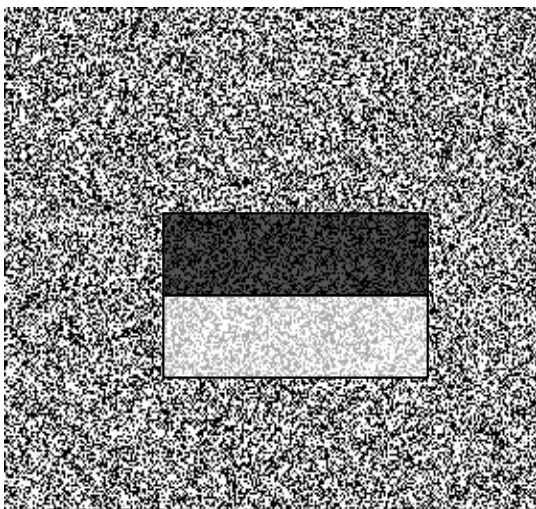*∑ (pixels in black area)*
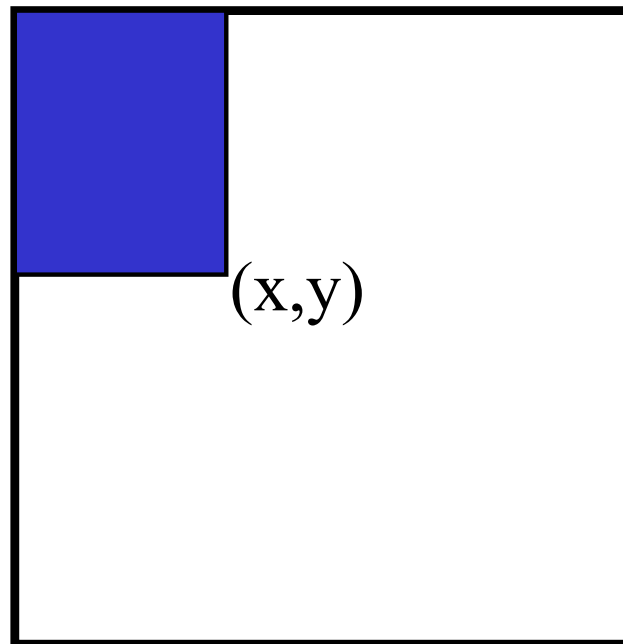
# Example

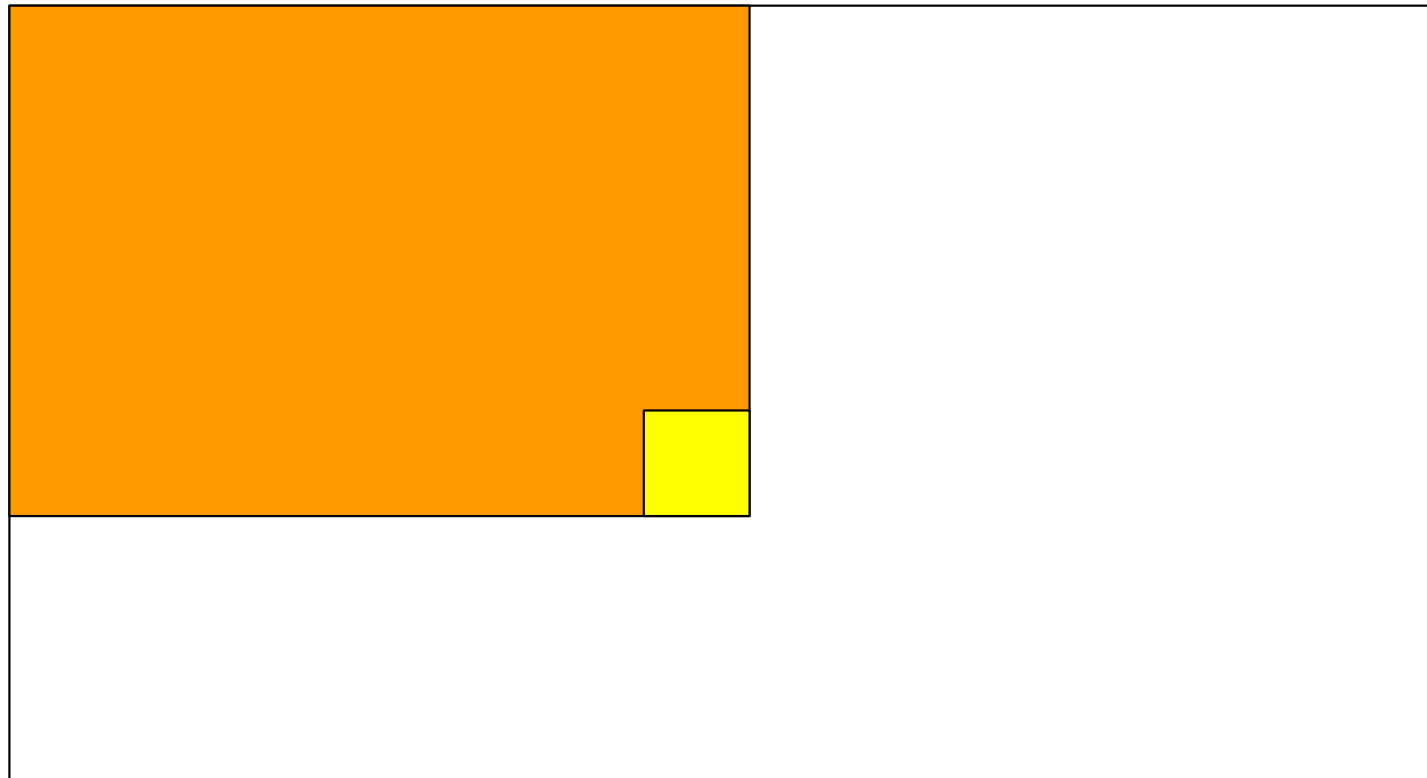

Source

Result

Source: S. Lazebnik

# Key idea 1: Fast computation with integral images

- The *integral image* computes a value at each pixel (*x,y*) that is the sum of the pixel values above and to the left of (*x,y*), inclusive

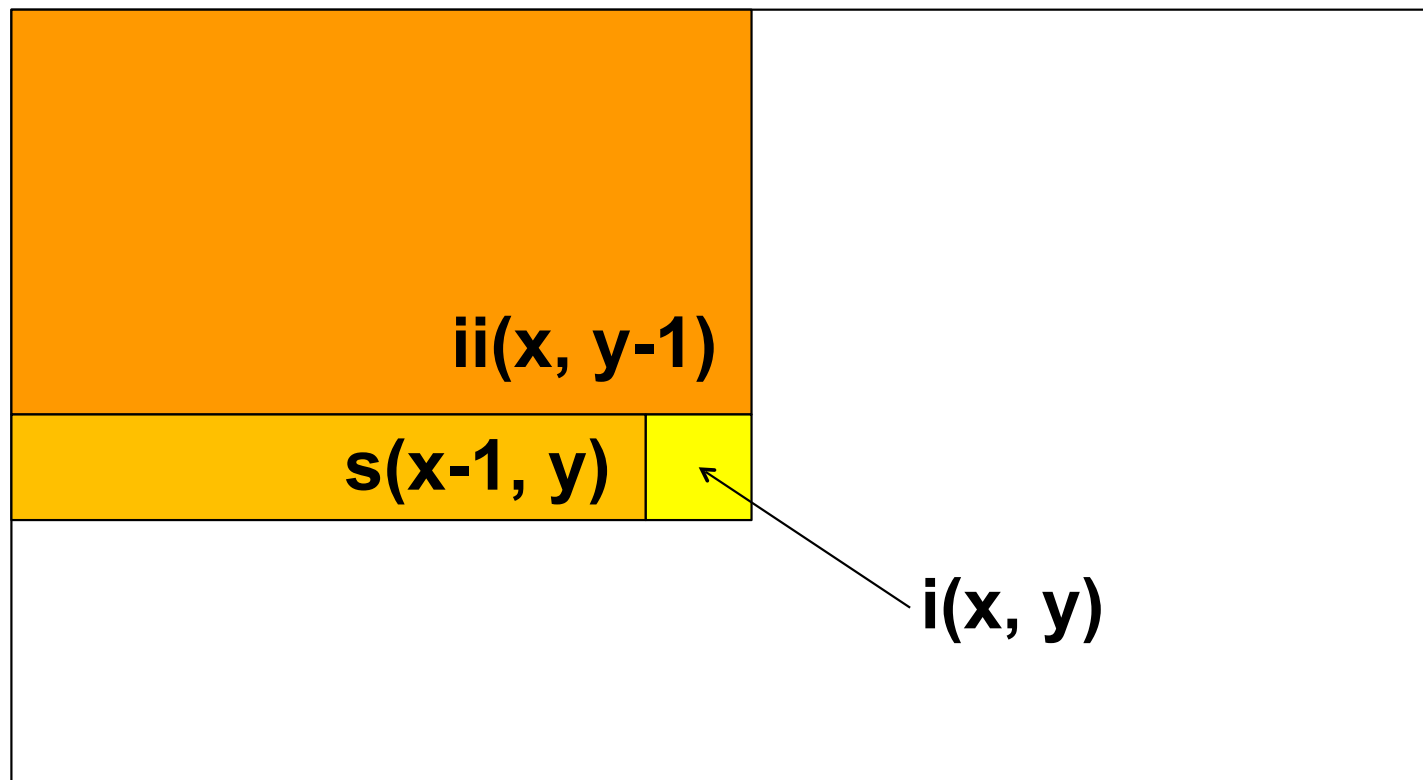- This can quickly be computed in one pass through the image

$(x,y)$

Source: S. Lazebnik

# Computing the integral image



Source: S. Lazebnik

# Computing the integral image

Cumulative row sum: $s(x, y) = s(x-1, y) + i(x, y)$

Integral image: $ii(x, y) = ii(x, y-1) + s(x, y)$



**ii(x, y-1)**

**s(x-1, y)**

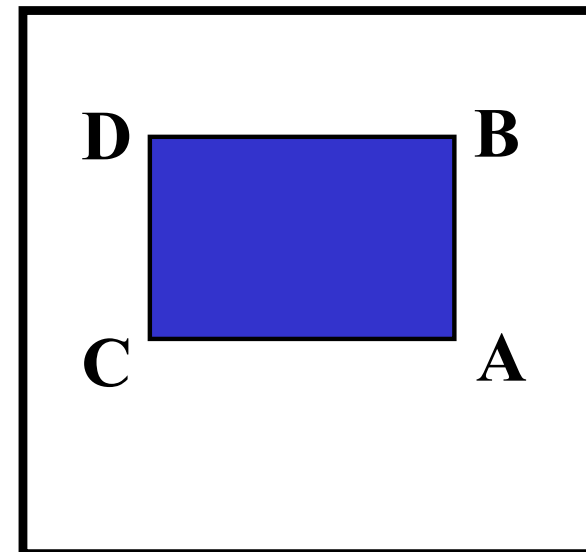**i(x, y)**

Source: S. Lazebnik
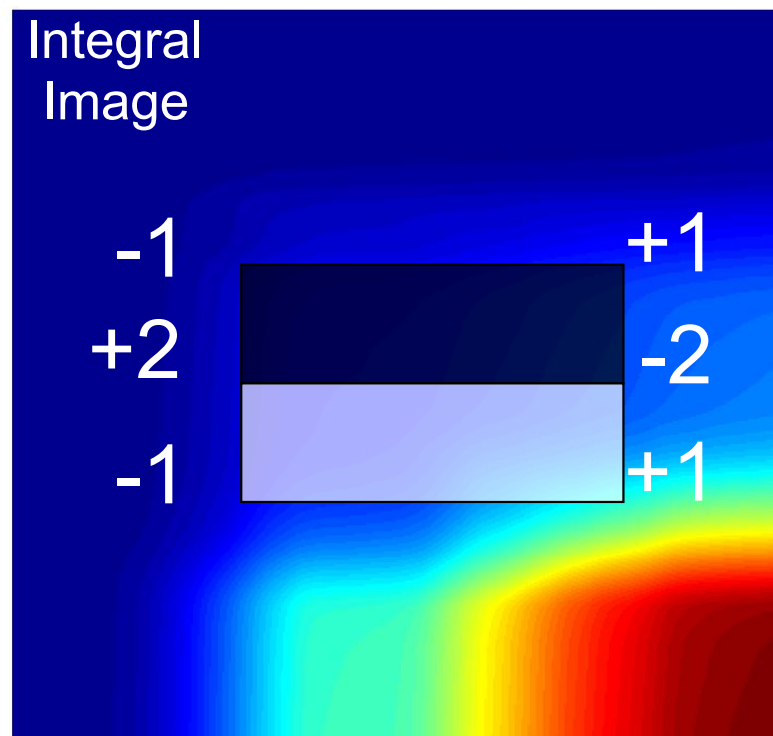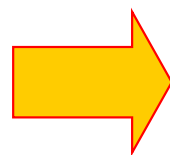
# Computing sum within a rectangle

- Let A,B,C,D be the values of the integral image at the corners of a rectangle

- Then the sum of original image values within the rectangle can be computed as:

  $$\text{sum} = A - B - C + D$$

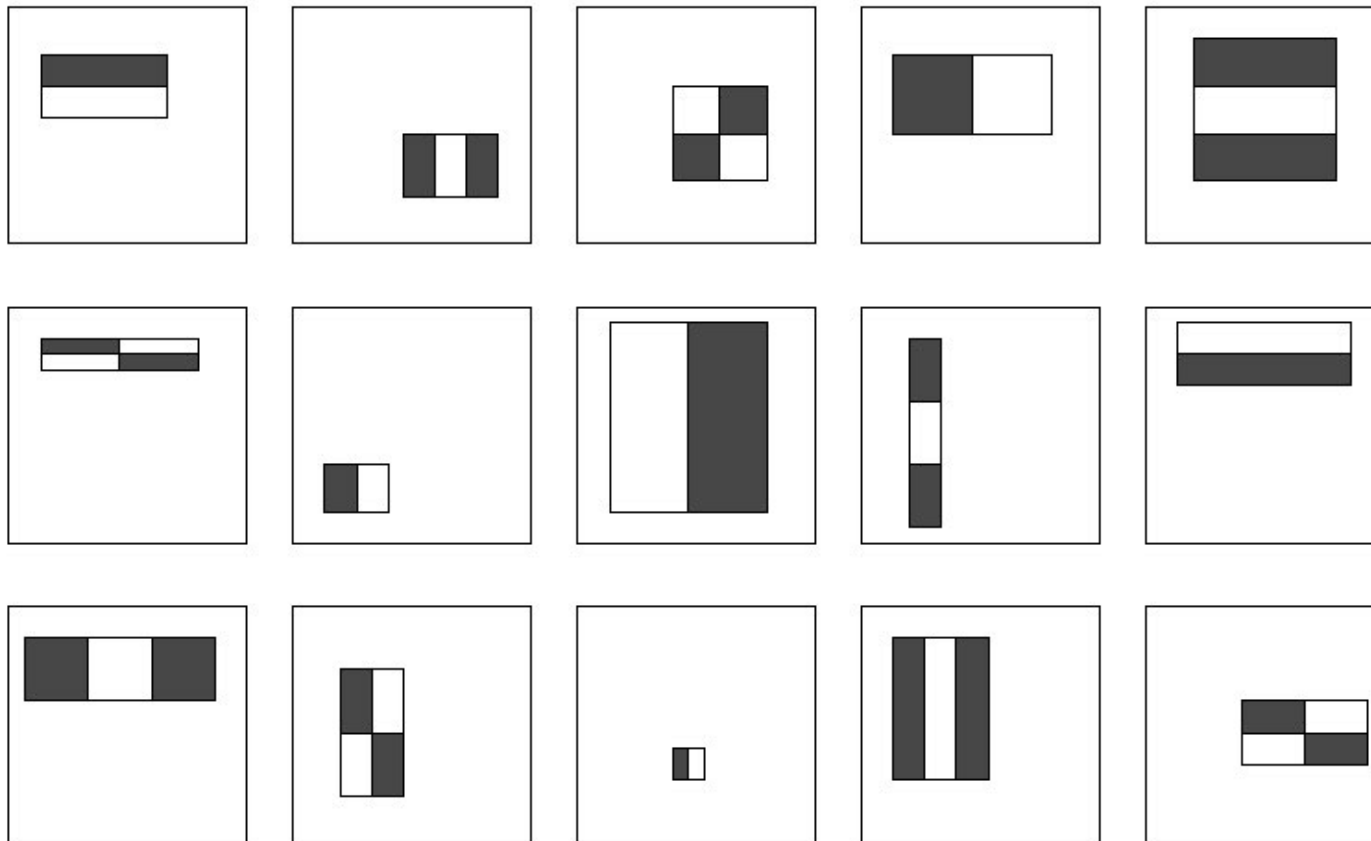- Only 3 additions are required for any size of rectangle!

# Computing a rectangle feature



Integral
Image

-1          +1

+2          -2

-1          +1

Source: S. Lazebnik

# Feature selection

- For a 24x24 detection region, the number of possible rectangle features is ~160,000!

# Feature selection

- For a 24x24 detection region, the number of possible rectangle features is ~160,000!
- At test time, it is impractical to evaluate the entire feature set
- Can we create a good classifier using just a small subset of all possible features?
- How to select such a subset?

# Key idea 2: Boosting

- *Boosting* is a classification scheme that combines *weak learners* into a more accurate *ensemble classifier*

- Weak learners based on rectangle filters:

value of rectangle feature

$$h_t(x) = \begin{cases} 1 & \text{if } p_t f_t(x) > p_t \theta_t \\ 0 & \text{otherwise} \end{cases}$$

window

parity     threshold

- Ensemble classification function:

$$C(x) = \begin{cases} 1 & \text{if } \sum_{t=1}^{T} \alpha_t h_t(x) > \frac{1}{2} \sum_{t=1}^{T} \alpha_t \\ 0 & \text{otherwise} \end{cases}$$

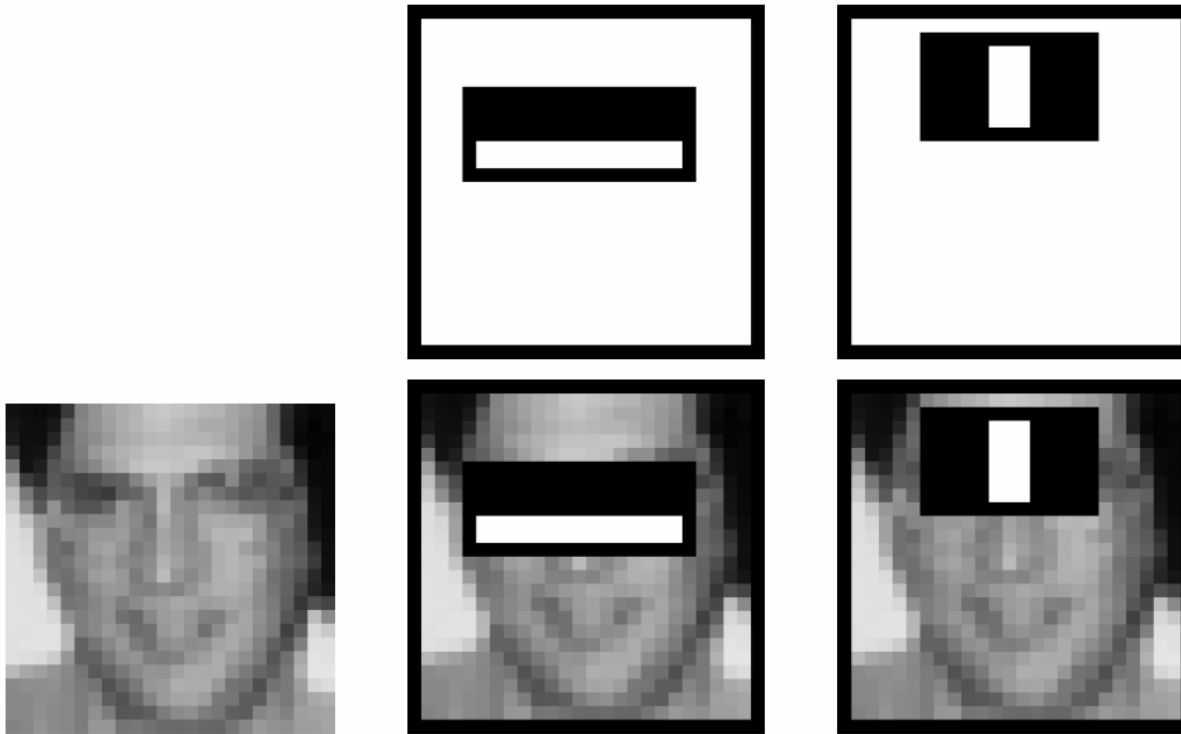learned weights

Source: S. Lazebnik

# Training procedure

- Initially, weight each training example equally

- In each boosting round:
  - Find the weak learner that achieves the lowest *weighted* training error
  - Raise the weights of training examples misclassified by current weak learner

- Compute final classifier as linear combination of all weak learners (weight of each learner is directly proportional to its accuracy)
  - Exact formulas for re-weighting and combining weak learners depend on the particular boosting scheme (e.g., AdaBoost)

Y. Freund and R. Schapire, A short introduction to boosting, *Journal of Japanese Society for Artificial Intelligence*, 14(5):771-780, September, 1999.

Source: S. Lazebnik

# Boosting for face detection

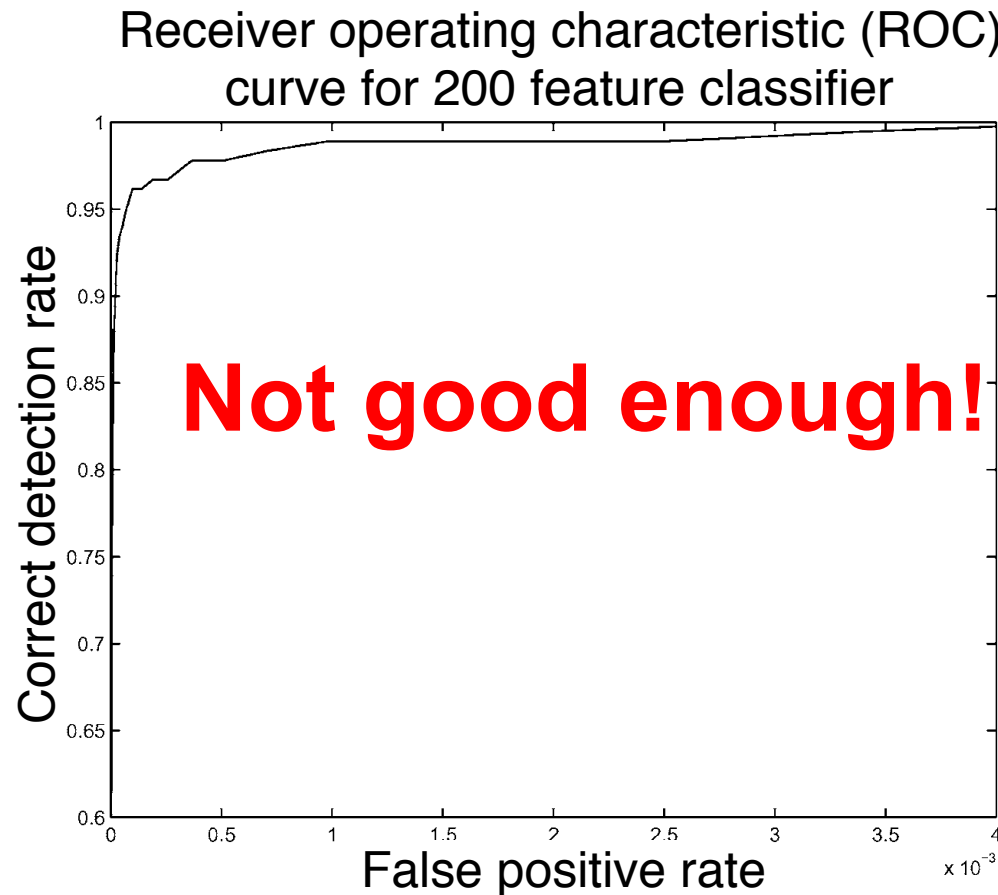- First two features selected by boosting:



This feature combination can yield 100% detection rate and 50% false positive rate

# Boosting

- ## Advantages of boosting

  - Flexibility in the choice of weak learners, boosting scheme

  - Testing is fast

  - Easy to implement

- ## Disadvantages
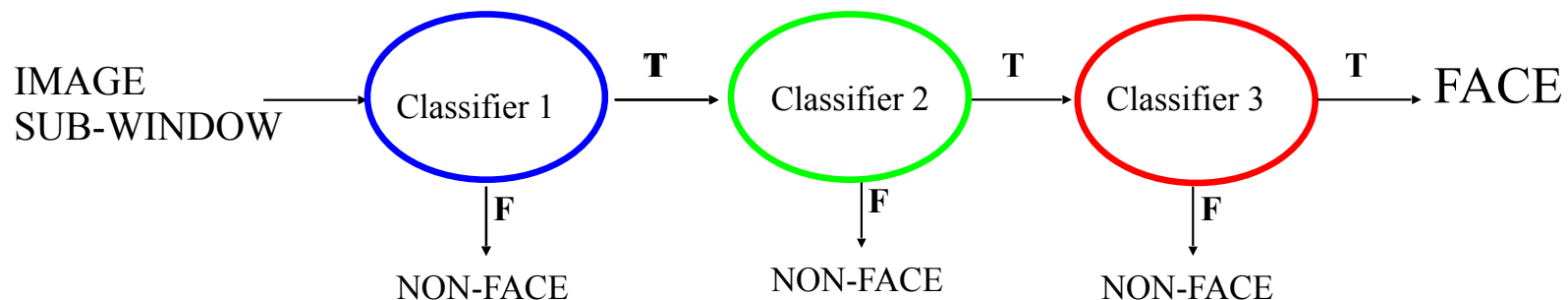
  - Needs many training examples

  - Training is slow

# Boosting for face detection

- A 200-feature classifier can yield 95% detection rate and a false positive rate of 1 in 14084



Receiver operating characteristic (ROC) curve for 200 feature classifier
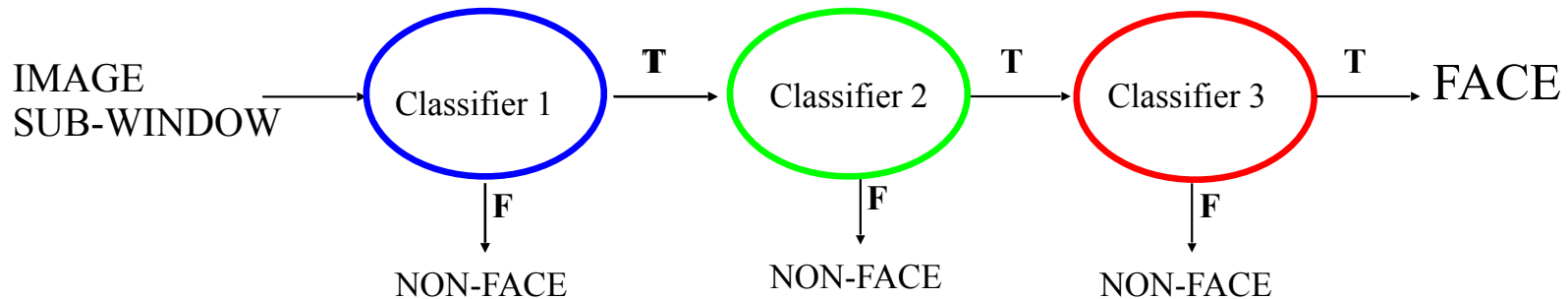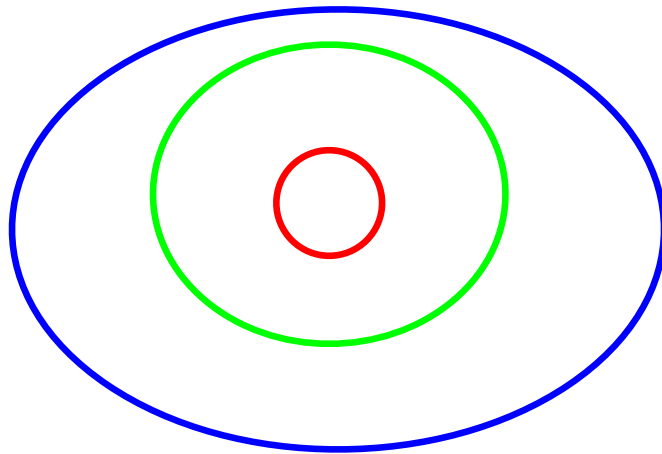
**Not good enough!**

# Key idea 3: Attentional cascade

- We start with simple classifiers which reject many of the negative sub-windows while detecting almost all positive sub-windows

- Positive response from the first classifier triggers the evaluation of a second (more complex) classifier, and so on

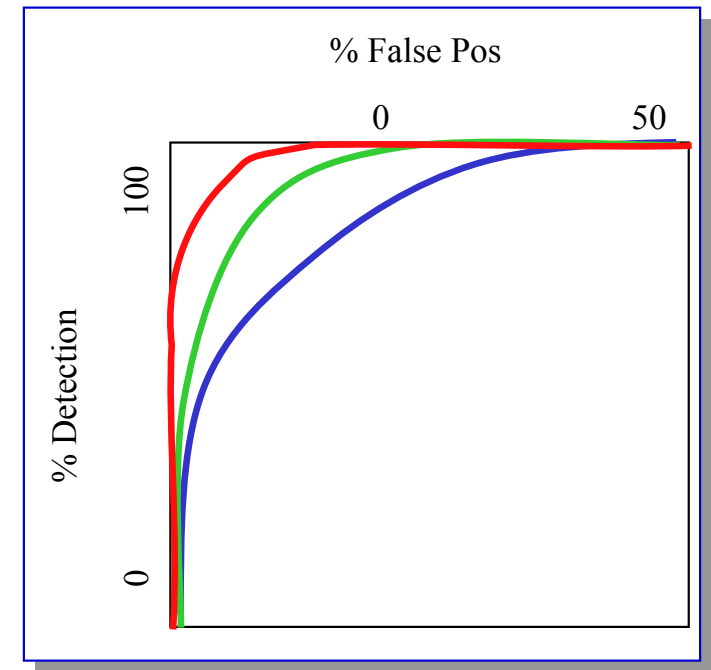- A negative outcome at any point leads to the immediate rejection of the sub-window

IMAGE SUB-WINDOW → Classifier 1 —**T**→ Classifier 2 —**T**→ Classifier 3 —**T**→ FACE

Classifier 1 —**F**→ NON-FACE

Classifier 2 —**F**→ NON-FACE

Classifier 3 —**F**→ NON-FACE

Source: S. Lazebnik

# Attentional cascade

- Chain classifiers that are progressively more complex and have lower false positive rates:

Receiver operating characteristic

% False Pos

0                    50

% Detection

100

0

IMAGE SUB-WINDOW → Classifier 1 —**T**→ Classifier 2 —**T**→ Classifier 3 —**T**→ FACE

Classifier 1 —**F**→ NON-FACE

Classifier 2 —**F**→ NON-FACE

Classifier 3 —**F**→ NON-FACE

Source: S. Lazebnik

# Training the cascade

- Set target detection and false positive rates for each stage

- Keep adding features to the current stage until its target rates have been met

    - Need to lower AdaBoost threshold to maximize detection (as opposed to minimizing total classification error)

    - Test on a *validation set*

- If the overall false positive rate is not low enough, then add another stage

- Use false positives from current stage as the negative training examples for the next stage
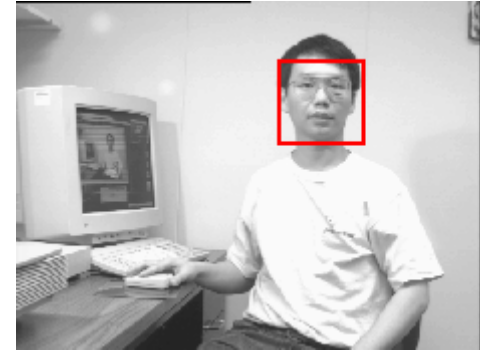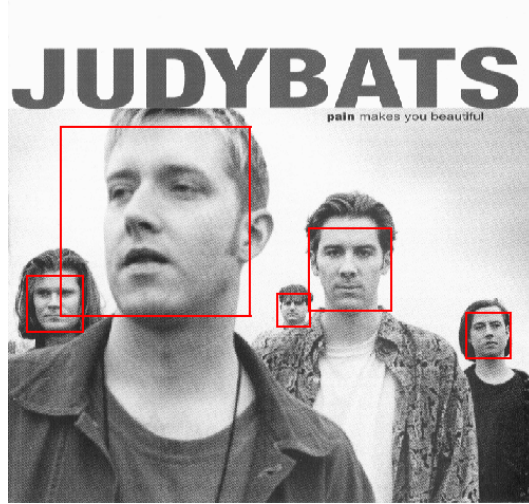
# The implemented system

- ## Training Data
  - 5000 faces
    - All frontal, rescaled to 24x24 pixels
  - 300 million non-faces
    - 9500 non-face images
  - Faces are normalized
    - Scale, translation

- ## Many variations
  - Across individuals
  - Illumination
  - Pose

# System performance

- Training time: "weeks" on 466 MHz Sun workstation

- 38 layers, total of 6061 features

- Average of 10 features evaluated per window on test set

- "On a 700 Mhz Pentium III processor, the face detector can process a 384 by 288 pixel image in about .067 seconds"

  - 15 Hz

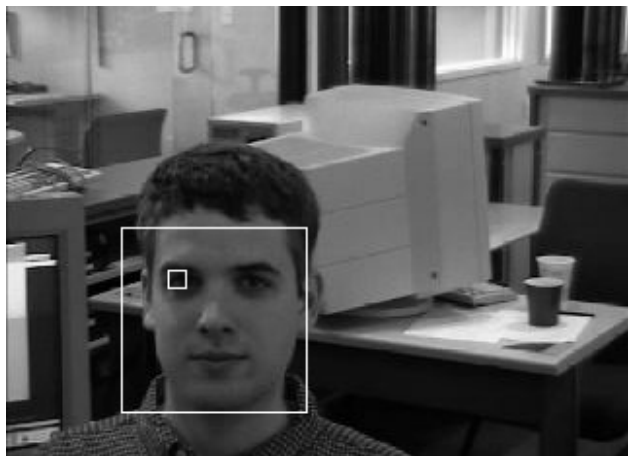  - 15 times faster than previous detector of comparable accuracy (Rowley et al., 1998)

Source: S. Lazebnik

# Output of Face Detector on Test Images
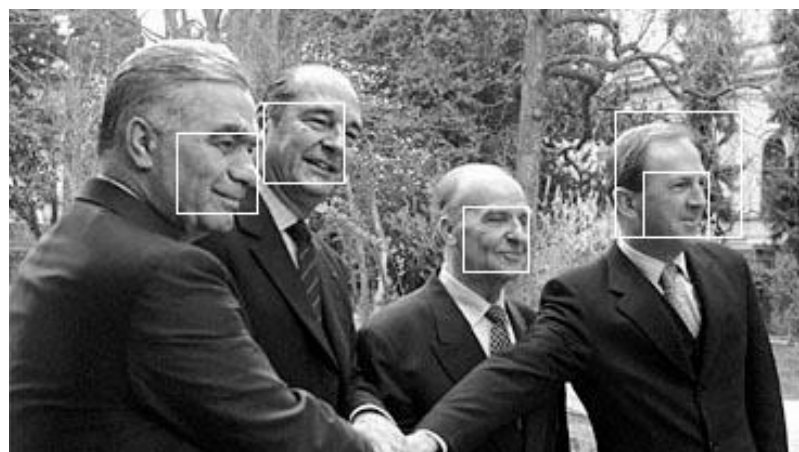
Source: S. Lazebnik

# Summary: Viola/Jones detector

- Rectangle features

- Integral images for fast computation

- Boosting for feature selection

- Attentional cascade for fast rejection of negative windows

- Caveat: Haar features work very well on structured faces, but much worse on other objects (coming up)
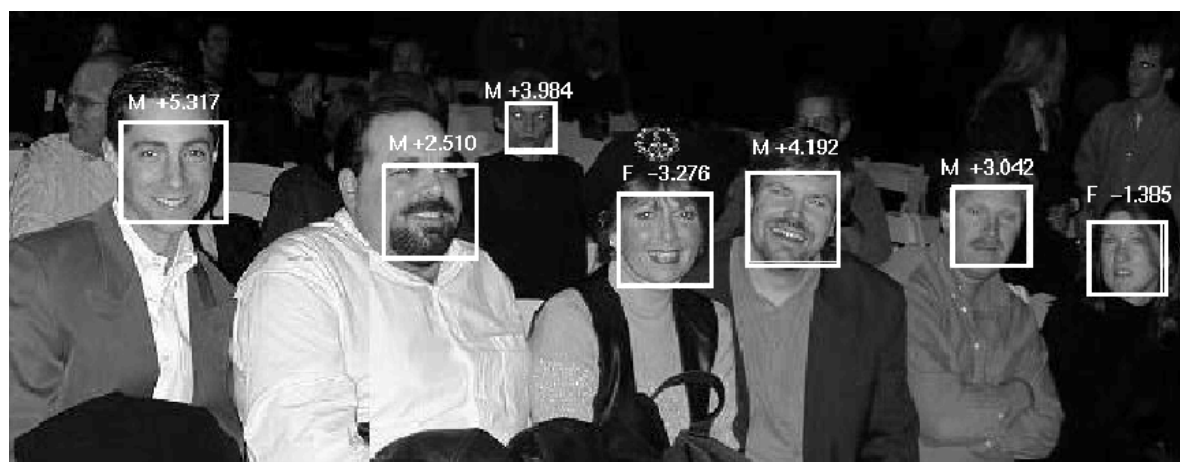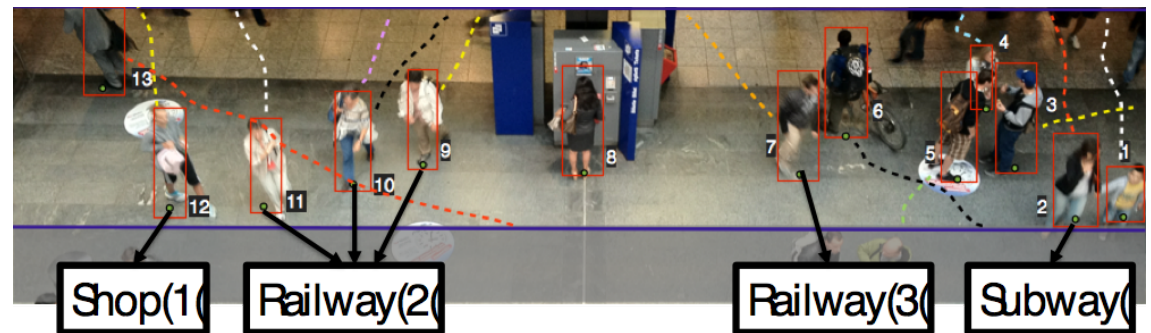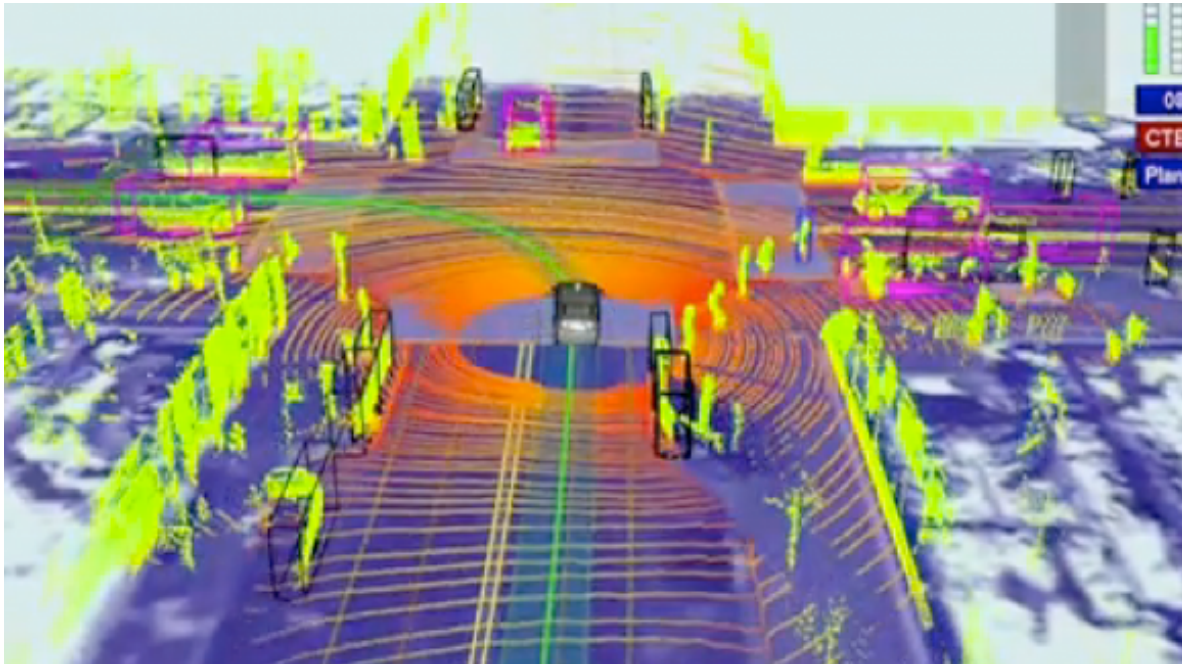
# Related detection tasks



Facial Feature Localization

Profile Detection

Male vs. female

Source: S. Lazebnik

# Pedestrian detection

# Pedestrian detection use cases



Alahi & Fei-Fei, 2014

# Faces vs pedestrians



**Faces**: the identifying features are in the structure



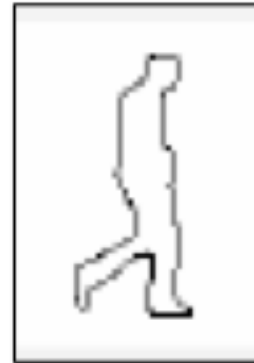**Pedestrian**: much more about shape

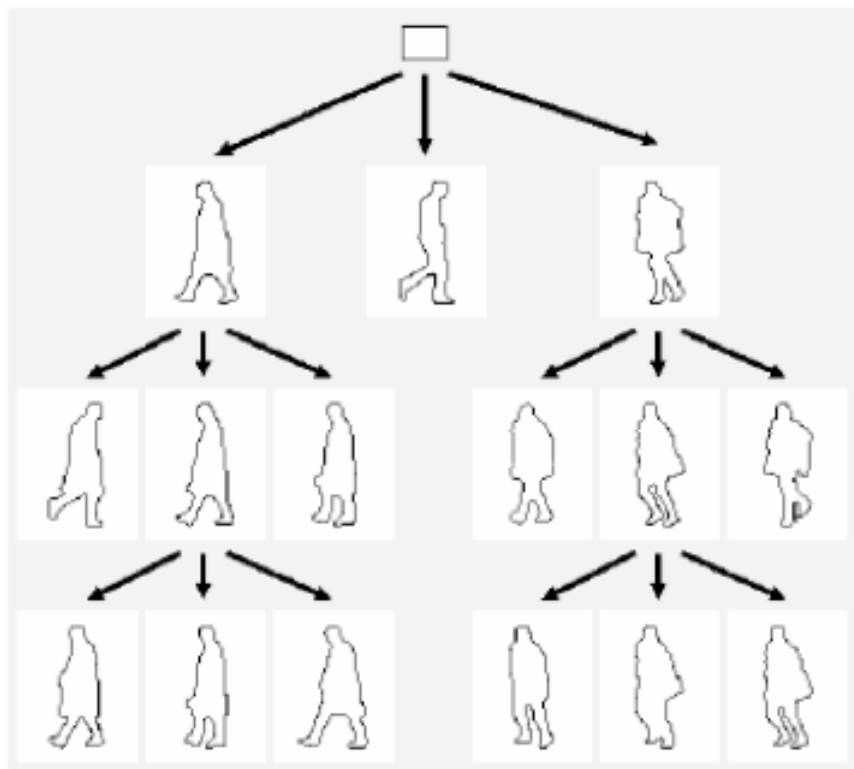# Chamfer matching



Input Image      Edge Detection      Template      Best Match

$$D_{chamfer}(T, I) \equiv \frac{1}{|T|} \sum_{t \in T} d_I(t)$$

Distance Transform

Gavrila & Philomin ICCV 1999

# Chamfer matching



Hierarchy of templates

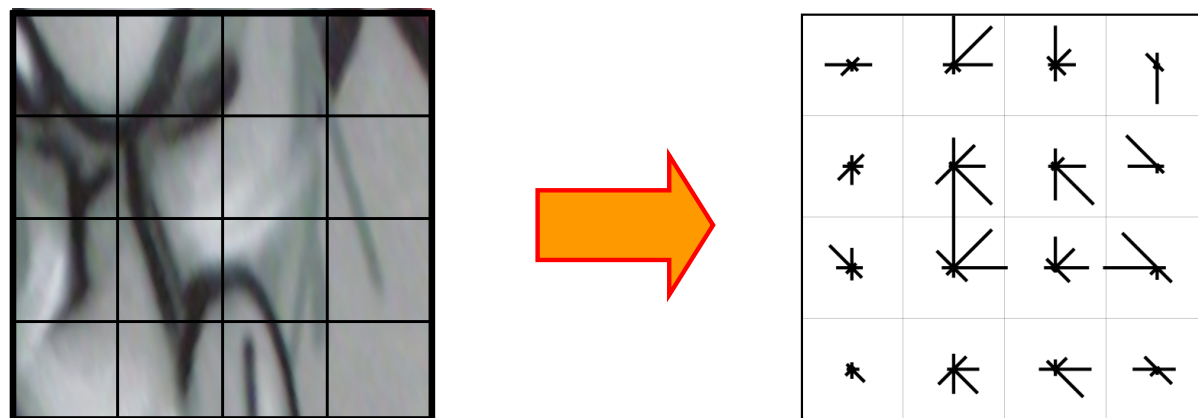Gavrila & Philomin ICCV 1999

# Positive and negative examples

+ thousands more...

+ millions more...

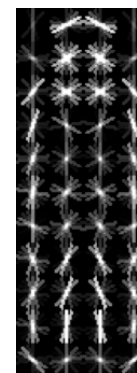N. Dalal and B. Triggs, Histograms of Oriented Gradients for Human Detection, CVPR 2005

# Window representation: edges/gradients



David G. Lowe. "Distinctive image features from scale-invariant keypoints."
*IJCV* 60 (2), pp. 91-110, 2004.

- Used to describe small patches, but now consider describing a whole detection window

# Histograms of oriented gradients (HOG)

Partition image into blocks and compute histogram of gradient orientations in each block
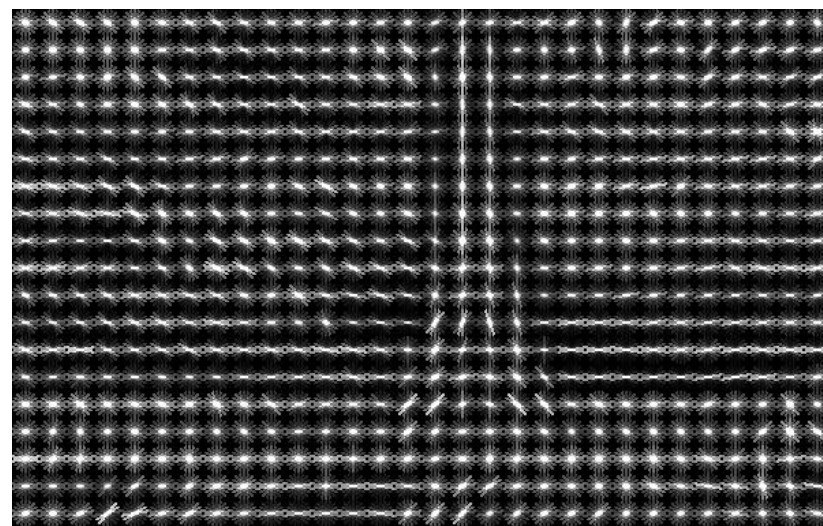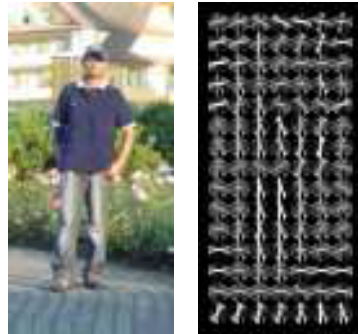
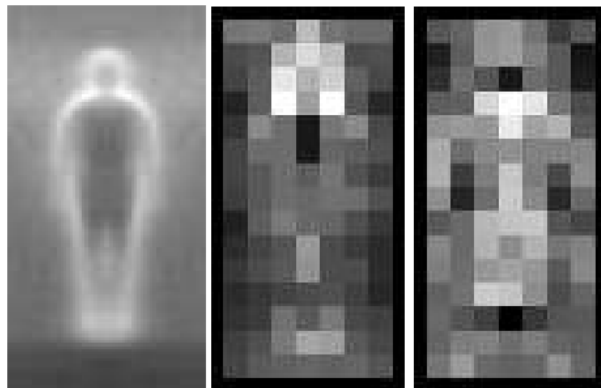## HxWx3 Image

## H'xW'xC' Image



Image credit: N. Snavely

N. Dalal and B. Triggs, Histograms of Oriented Gradients for Human Detection, CVPR 2005

# Person detection, ca. 2005 (Dalal Triggs)

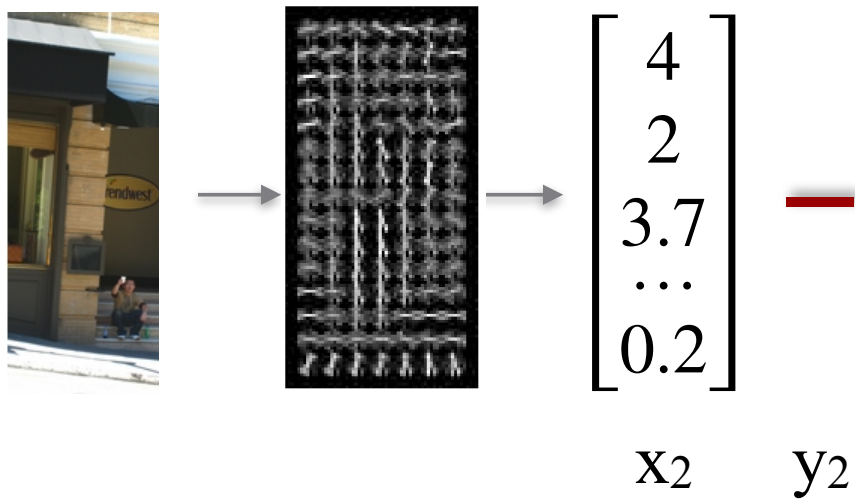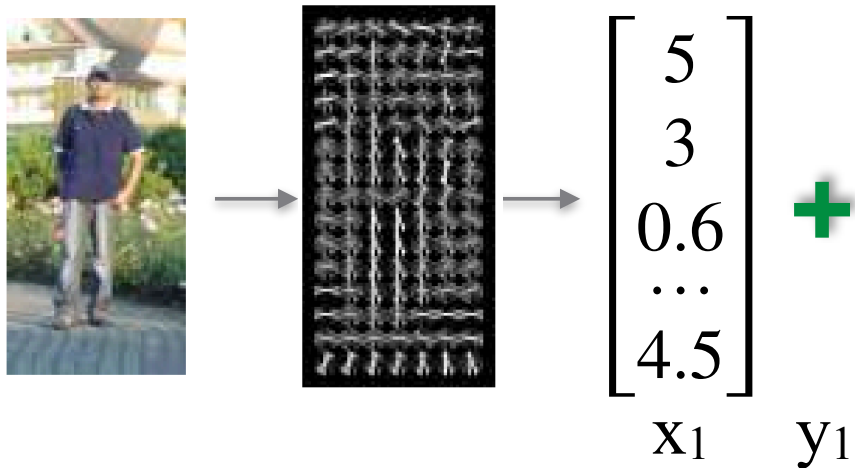1. Represent each example with a single, fixed HoG template
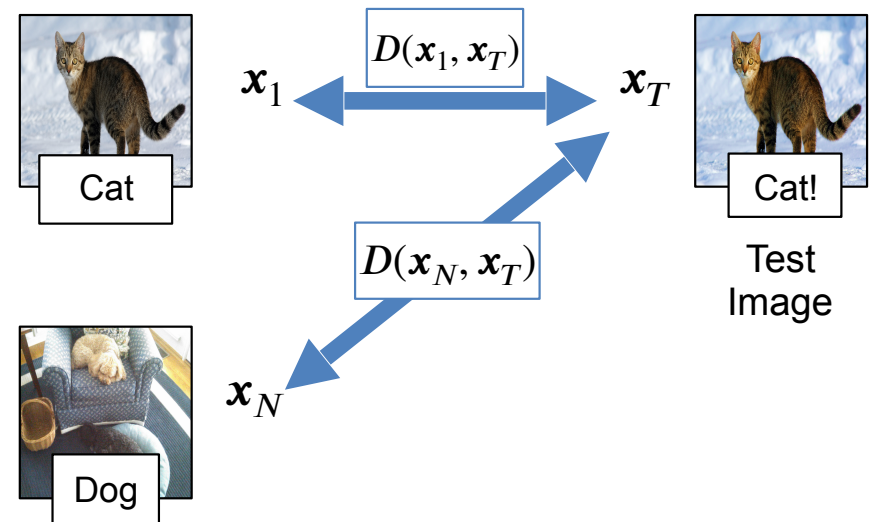


2. Learn a single linear detector



N. Dalal and B. Triggs, Histograms of Oriented Gradients for Human Detection, CVPR 2005

# Training



$$x_1 \quad y_1$$

$$x_2 \quad y_2$$

## Nearest neighbor classifier

$$x_1 \quad D(x_1, x_T) \quad x_T$$

Cat

$$D(x_N, x_T)$$

Cat!

Test Image

$$x_N$$

Dog

## Least squares classifier

Objective

$$\operatorname*{argmin}_{\boldsymbol{w}} \left\| \boldsymbol{y} - \boldsymbol{X}\boldsymbol{w} \right\|_2^2 + \lambda \left\| \boldsymbol{w} \right\|_2^2$$

Inference

$$\boldsymbol{w}^T \boldsymbol{x} > t$$

# Alternate objective: Multiclass SVM

Inference (x):

$$\underset{k}{\mathrm{argmax}}\ (\boldsymbol{Wx})_k$$

(Take the class whose weight vector gives the highest score)

# Alternate objective: Multiclass SVM

Inference (x):  $\underset{k}{\text{argmax}} \ (\boldsymbol{Wx})_k$

(Take the class whose weight vector gives the highest score)

Training ($\mathbf{x}_i, y_i$):

$$\underset{W}{\arg\min} \ \lambda\|W\|_2^2 + \sum_i^n \sum_{j \neq y_i} \max(0, \mathbf{w_j}^T \mathbf{x_i} - \mathbf{w_{y_i}}^T \mathbf{x_i} + m)$$
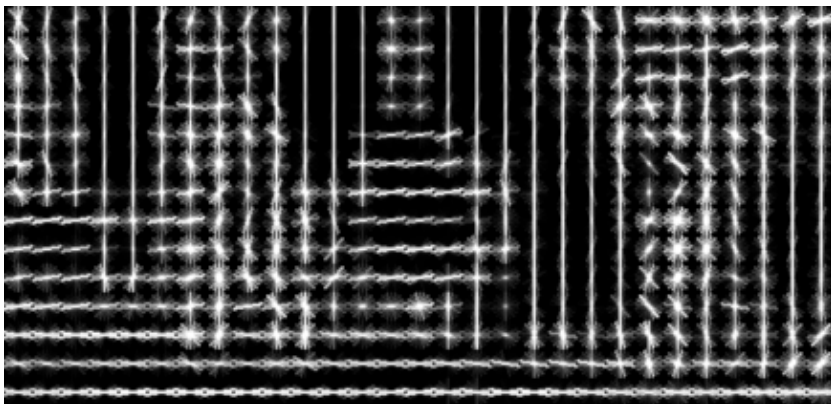
Regularization

Over all data points

For every class j that's NOT the correct one ($y_i$)

Pay no penalty if prediction for class $y_i$ is bigger than j by m ("margin"). Otherwise, pay proportional to the score of the wrong class.
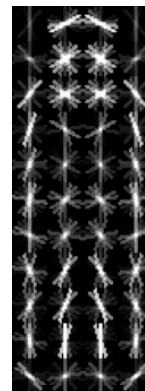
Credit: D. Fouhey

# Pedestrian detection with HOG

- Train HOG-based pedestrian "template" with a linear model

- At test time, convolve feature map with template

- Find local maxima of response

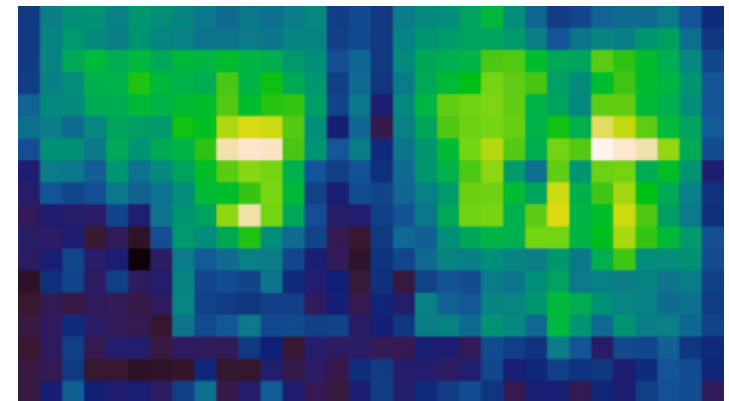- For multi-scale detection, repeat over multiple levels of a HOG *pyramid*

HOG feature map · Template · Detector response map



N. Dalal and B. Triggs, Histograms of Oriented Gradients for Human Detection, CVPR 2005
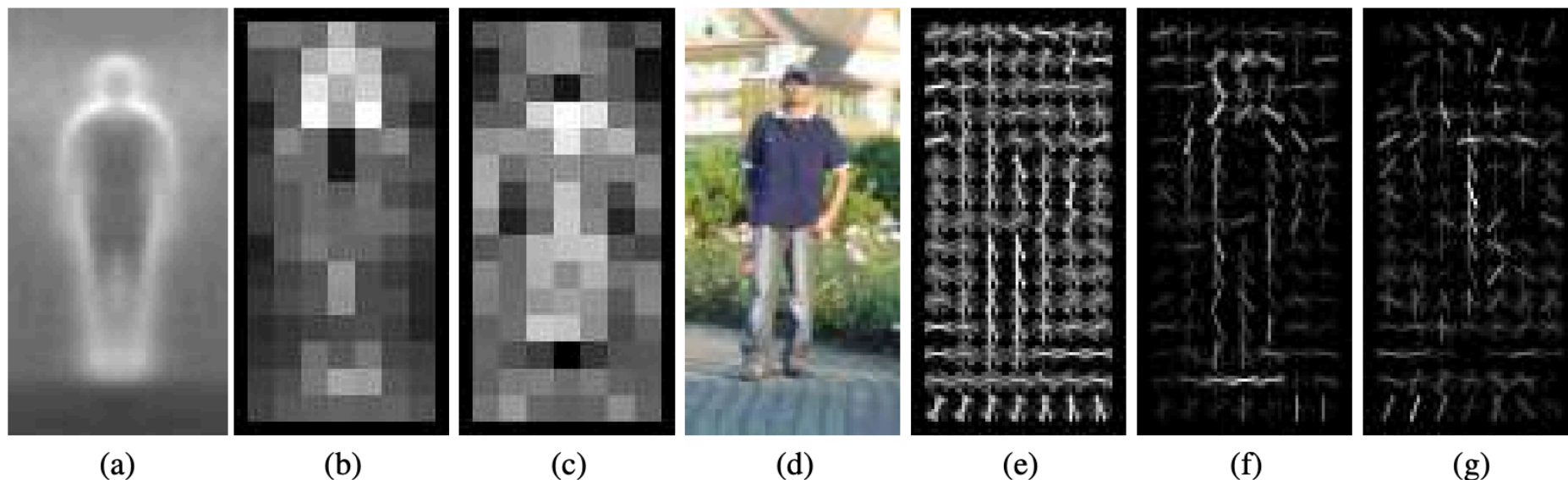
# What is learned



Figure 6. Our HOG detectors cue mainly on silhouette contours (especially the head, shoulders and feet). The most active blocks are centred on the image background just *outside* the contour. (a) The average gradient image over the training examples. (b) Each "pixel" shows the maximum positive SVM weight in the block centred on the pixel. (c) Likewise for the negative SVM weights. (d) A test image. (e) It's computed R-HOG descriptor. (f,g) The R-HOG descriptor weighted by respectively the positive and the negative SVM weights.
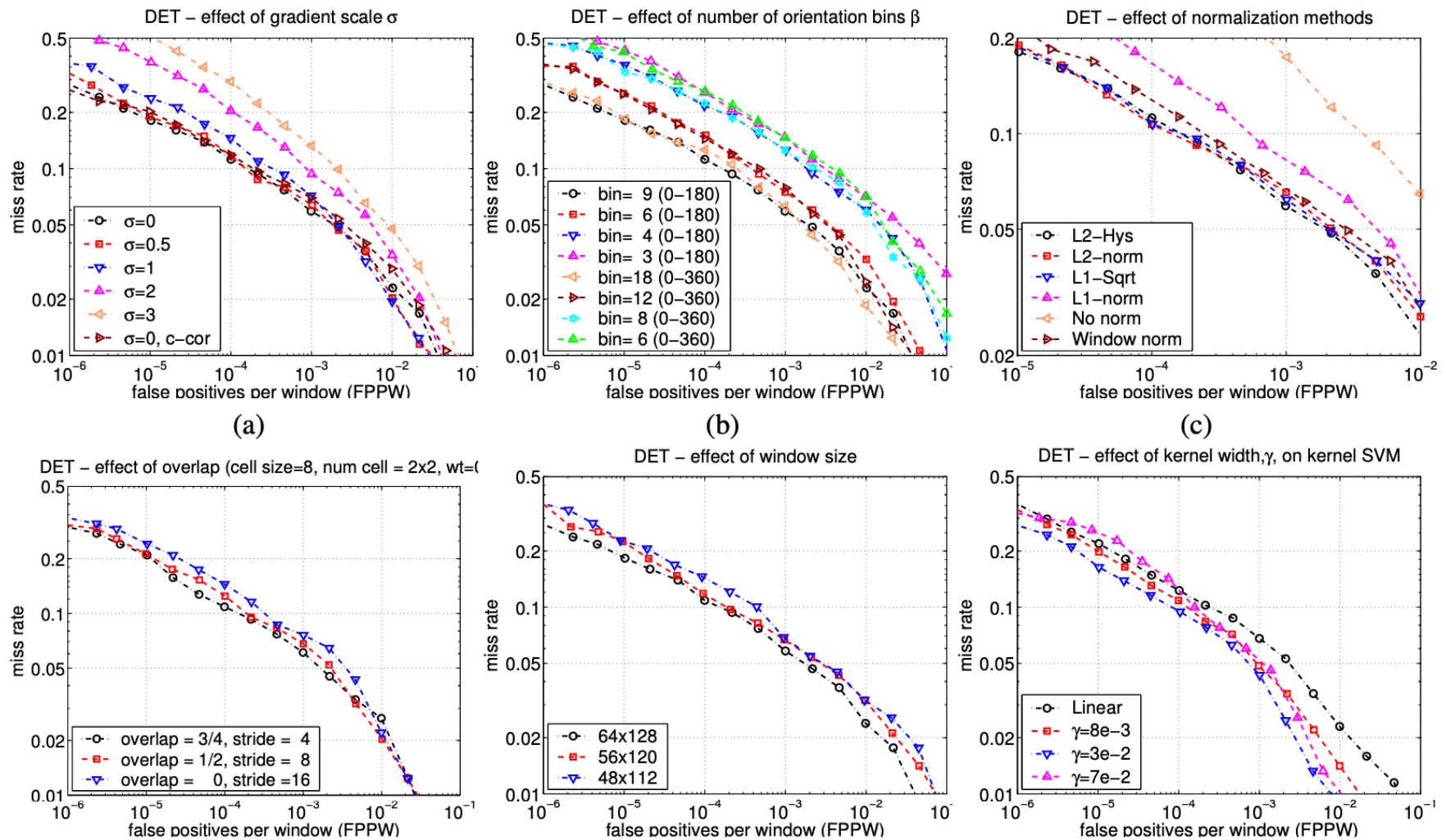
N. Dalal and B. Triggs, Histograms of Oriented Gradients for Human Detection, CVPR 2005

# Example detections



[Dalal and Triggs, CVPR 2005]

# Evaluation: effect of hyperparameters



[Dalal and Triggs, CVPR 2005]

Slide Credit: S. Lazebnik

# Summary: key concepts

1. ML for image classification (NN, kNN, linear models)

2. Viola Jones face detection

3. Dalal Triggs pedestrian detection

# Next time: object classification