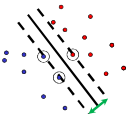


Support vector machines and kernels

Thurs April 20
Kristen Grauman
UT Austin



Last time

- Sliding window object detection wrap-up
 - Attentional cascade
 - Applications / examples
 - Pros and cons
- Supervised classification continued
 - Nearest neighbors

Today

- Supervised classification continued
 - Nearest neighbors (wrap up)
 - Support vector machines
 - HoG pedestrians example
 - Kernels
 - Multi-class from binary classifiers
 - Pyramid match kernels
 - Evaluation
 - Scoring an object detector
 - Scoring a multi-class recognition system

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

Black = negative
Red = positive

k = 5

If query lands here, the 5 NN consist of 3 negatives and 2 positives, so we classify it as negative.

Source: D. Lowe

Window-based models: Three case studies

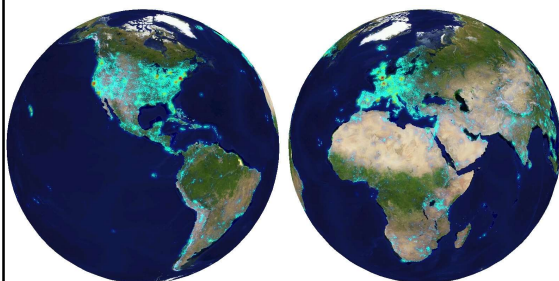
 Boosting + face detection Viola & Jones	 NN + scene Gist classification e.g., Hays & Efros	 SVM + person detection e.g., Dalal & Triggs
--	--	--

Where in the World?



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

6+ million geotagged photos
by 109,788 photographers



Annotated by Flickr users

Which scene properties are relevant?

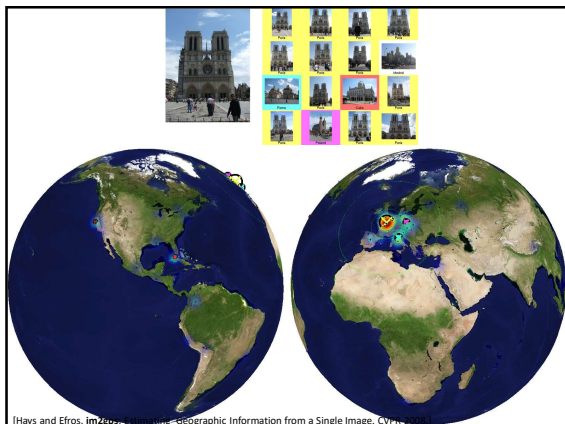
- **Gist scene descriptor**
- **Color Histograms** - L*A*B* 4x14x14 histograms
- **Texton Histograms** - 512 entry, filter bank based
- **Line Features** - Histograms of straight line stats

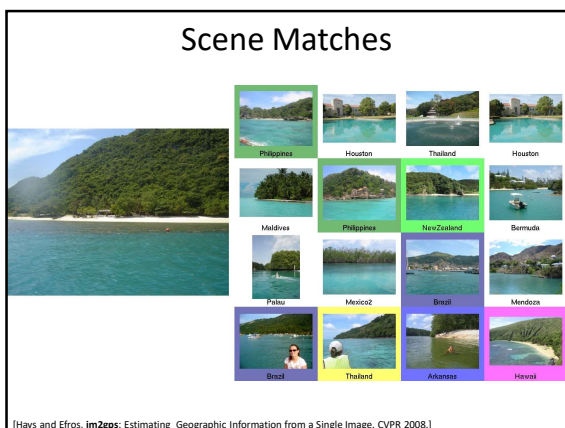
Im2gps: Scene Matches

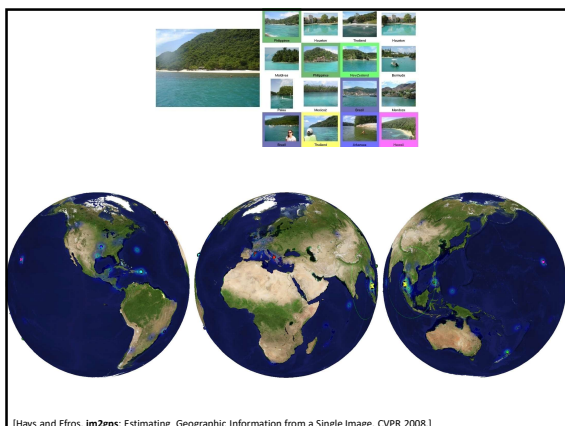
[Hays and Efros. *Im2gps: Estimating Geographic Information from a Single Image*. CVPR 2008.]

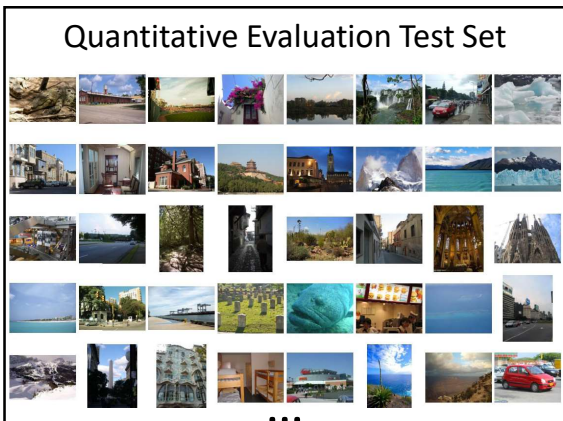
Im2gps: Scene Matches

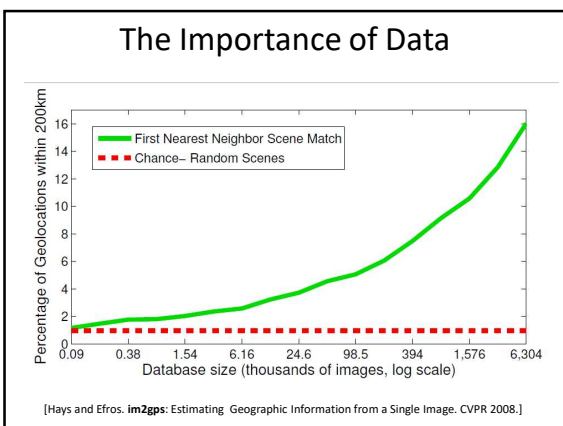
[Hays and Efros. *Im2gps: Estimating Geographic Information from a Single Image*. CVPR 2008.]





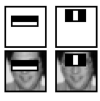

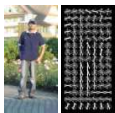


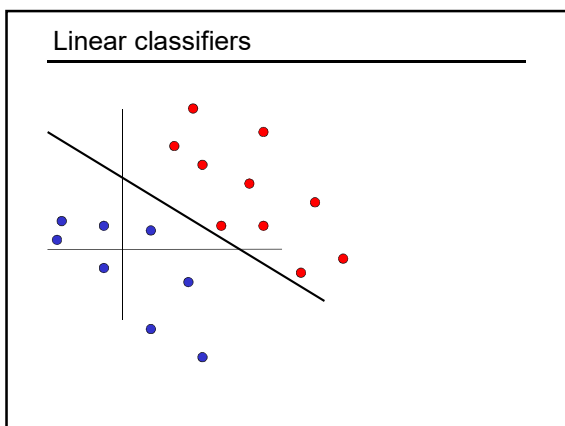


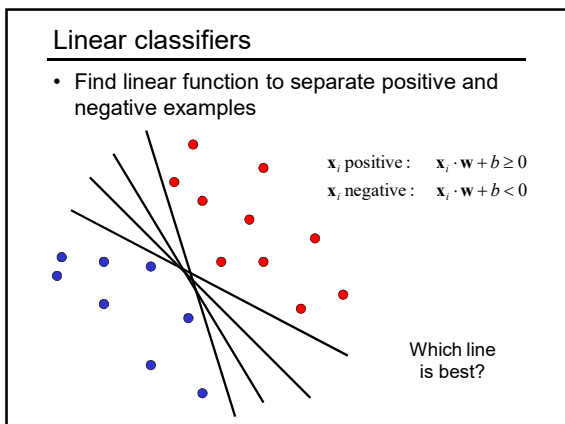


- ### Nearest neighbors: pros and cons
- **Pros:**
 - Simple to implement
 - Flexible to feature / distance choices
 - Naturally handles multi-class cases
 - Can do well in practice with enough representative data
 - **Cons:**
 - Large search problem to find nearest neighbors
 - Storage of data
 - Must know we have a meaningful distance function
- Kristen Grauman

Window-based models: Three case studies

 <p>Boosting + face detection</p> <p>Viola & Jones</p>	 <p>NN + scene Gist classification</p> <p>e.g., Hays & Efros</p>	 <p>SVM + person detection</p> <p>e.g., Dalal & Triggs</p>
---	---	---





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.

$w \cdot x + b = 1$
 $w \cdot x + b = 0$
 $w \cdot x + b = -1$

x_i positive ($y_i = 1$): $x_i \cdot w + b \geq 1$
 x_i negative ($y_i = -1$): $x_i \cdot w + b \leq -1$
 For support vectors, $x_i \cdot w + b = \pm 1$

Support vectors Margin

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

Support vector machines

- Want line that maximizes the margin.

$w \cdot x + b = 1$
 $w \cdot x + b = 0$
 $w \cdot x + b = -1$

x_i positive ($y_i = 1$): $x_i \cdot w + b \geq 1$
 x_i negative ($y_i = -1$): $x_i \cdot w + b \leq -1$
 For support vectors, $x_i \cdot w + b = \pm 1$

Distance between point and line: $\frac{|x_i \cdot w + b|}{\|w\|}$
 For support vectors:
 $\frac{w^T x + b}{\|w\|} = \frac{\pm 1}{\|w\|}$ $M = \left| \frac{1}{\|w\|} - \frac{-1}{\|w\|} \right| = \frac{2}{\|w\|}$

Support vectors Margin M

Support vector machines

- Want line that maximizes the margin.

x_i positive ($y_i = 1$): $x_i \cdot w + b \geq 1$
 x_i negative ($y_i = -1$): $x_i \cdot w + b \leq -1$
 For support vectors, $x_i \cdot w + b = \pm 1$
 Distance between point and line: $\frac{|x_i \cdot w + b|}{\|w\|}$
 Therefore, the margin is $2 / \|w\|$

Support vectors Margin M

Finding the maximum margin line

- Maximize margin $2/\|w\|$
- Correctly classify all training data points:
 - x_i positive ($y_i = 1$): $x_i \cdot w + b \geq 1$
 - x_i negative ($y_i = -1$): $x_i \cdot w + b \leq -1$

Quadratic optimization problem:

Minimize $\frac{1}{2} w^T w$

Subject to $y_i(w \cdot x_i + b) \geq 1$

Finding the maximum margin line

- Solution: $w = \sum_i \alpha_i y_i x_i$

learned weight

Support vector

C. Burges, A Tutorial on Support Vector Machines for Pattern Recognition, Data Mining and Knowledge Discovery.

Finding the maximum margin line

- Solution: $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$
 $b = y_i - \mathbf{w} \cdot \mathbf{x}_i$ (for any support vector)
 $\mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$
- Classification function:
 $f(x) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$
 $= \text{sign}\left(\sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b\right)$

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

C. Burges, A Tutorial on Support Vector Machines for Pattern Recognition, Data Mining and Knowledge Discovery.

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

Abstract

We study the question of feature sets for robust visual object recognition, adapting linear SVM based human detection as a test case. After reviewing existing edge and gradient based descriptors, we show experimentally that grids of Histograms of Oriented Gradient (HOG) descriptors significantly outperform existing feature sets for human detection. We study the influence of each stage of the computation on performance, concluding that fine-scale gradients, fine orientation binning, relatively coarse spatial binning, and high-quality local contrast normalization in overlapping descriptor blocks are all important for good results. The new approach gives near-perfect separation on the original MIT pedestrian database, so we introduce a more challenging dataset containing over 1800 annotated human images with a large range of pose variations and backgrounds.

We briefly discuss previous work on human detection in §2, give an overview of our method §3, describe our data sets in §4 and give a detailed description and experimental evaluation of each stage of the process in §5-6. The main conclusions are summarized in §7.

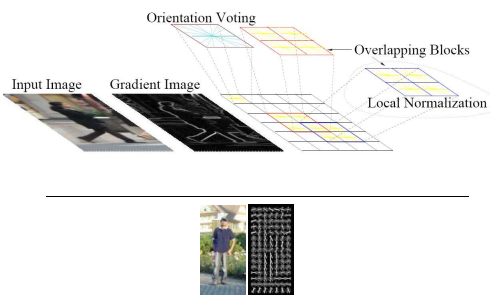
2 Previous Work

There is an extensive literature on object detection, but here we mention just a few relevant papers on human detection [18, 17, 22, 16, 20]. See [6] for a survey. Papageorgiou *et al* [18] describe a pedestrian detector based on a polynomial SVM using rectified Haar wavelets as input descriptors, with a parts (subwindow) based variant in [17]. Desouster *et al* give an optimized version of this [2]. Gavrilu & Phalonen [8] take a more direct approach, extracting edge images and matching them to a set of learned exemplars using chamfer distance. This has been used in a practical real-time pedestrian detection system [7]. Viola *et al* [22] build an efficient

1 Introduction

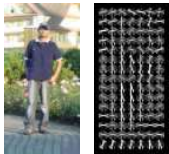
- CVPR 2005
- 18,317 citations

HoG descriptor



Dalal & Triggs, CVPR 2005


Person detection with HoG's & linear SVM's



- Map each grid cell in the input window to a histogram counting the gradients per orientation.
- Train a linear SVM using training set of pedestrian vs. non-pedestrian windows.

Dalal & Triggs, CVPR 2005

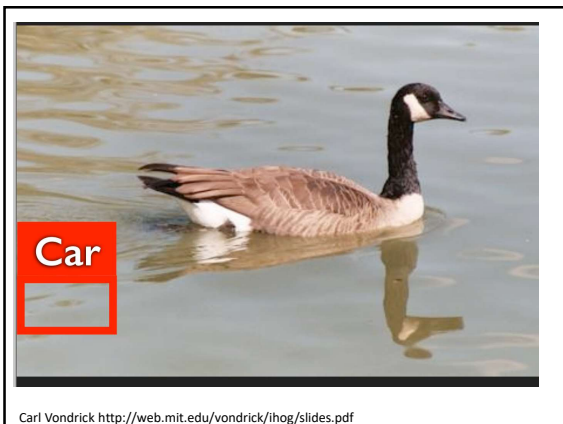
Person detection with HoGs & linear SVMs

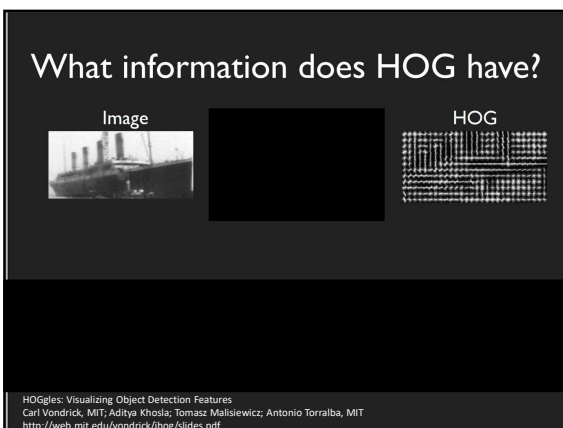


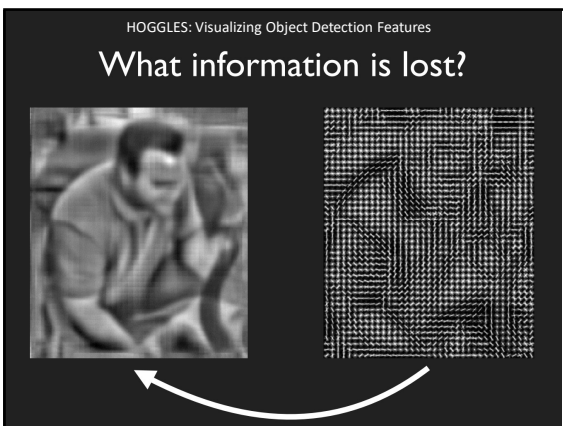
- Histograms of Oriented Gradients for Human Detection, [Navneet Dalal](#), [Bill Triggs](#). International Conference on Computer Vision & Pattern Recognition - June 2005 <http://lear.inrialpes.fr/pubs/2005/DT05/>

Understanding classifier mistakes









HOGgles: Visualizing Object Detection Features

Method: Paired Dictionary

HOGgles: Visualizing Object Detection Features
Carl Vondrick, MIT; Aditya Khosla; Tomasz Malisiewicz; Antonio Torralba, MIT
<http://web.mit.edu/vondrick/hog/slides.pdf>

HOGgles: Visualizing Object Detection Features

A microscope to view HOG

HOGgles: Visualizing Object Detection Features;
Carl Vondrick, MIT; Aditya Khosla; Tomasz Malisiewicz;
Antonio Torralba, MIT
<http://web.mit.edu/vondrick/hog/slides.pdf>

HOGgles: Visualizing Object Detection Features

HOGgles: Visualizing Object Detection Features; ICCV 2013
Carl Vondrick, MIT; Aditya Khosla; Tomasz Malisiewicz; Antonio Torralba, MIT
<http://web.mit.edu/vondrick/hog/slides.pdf>



Questions

- What if the data is not linearly separable?

Non-linear SVMs

- Datasets that are linearly separable with some noise work out great:
- But what are we going to do if the dataset is just too hard?
- How about... mapping data to a higher-dimensional space:

Non-linear SVMs: feature spaces

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

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

Nonlinear SVMs

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

$$K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$$
- This gives a nonlinear decision boundary in the original feature space:

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

"Kernel trick": Example

2-dimensional vectors $x = [x_1 \ x_2]$;
 let $K(x_i, x_j) = (1 + x_i^T x_j)^2$
 Need to show that $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$:

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

Examples of kernel functions

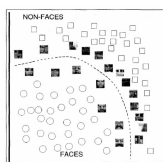
- Linear: $K(x_i, x_j) = x_i^T x_j$

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

- Histogram intersection: $K(x_i, x_j) = \sum_k \min(x_i(k), x_j(k))$

SVMs for recognition

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



Kristen Grauman

Questions

- 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

Kristen Grauman

SVMs: Pros and cons

- Pros
 - 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 small training sample sizes
- Cons
 - No “direct” multi-class SVM, must combine two-class SVMs
 - Can be tricky to select best kernel function for a problem
 - Computation, memory
 - During training time, must compute matrix of kernel values for every pair of examples
 - Learning can take a very long time for large-scale problems

Adapted from Lene Logstrup

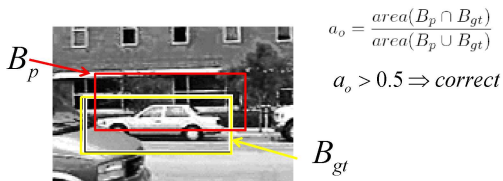
Scoring a sliding window detector



If prediction and ground truth are *bounding boxes*, when do we have a correct detection?

Kristen Grauman

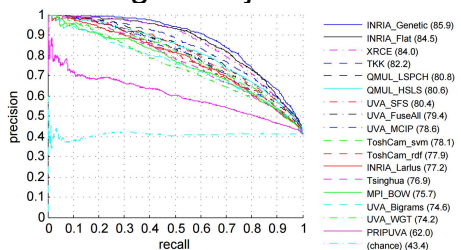
Scoring a sliding window detector



We'll say the detection is correct (a "true positive") if the intersection of the bounding boxes, divided by their union, is > 50%.

Kristen Grauman

Scoring an object detector



- If the detector can produce a *confidence score* on the detections, then we can plot its precision vs. recall as a threshold on the confidence is varied.
- **Average Precision (AP)**: mean precision across recall levels.

Summary: This week

- Object recognition as classification task
 - Boosting (face detection ex)
 - Support vector machines and HOG (person detection ex)
 - Pyramid match kernels
 - Hoggles visualization for understanding classifier mistakes
 - Nearest neighbors and global descriptors (scene rec ex)
- Sliding window search paradigm
 - Pros and cons
 - Speed up with attentional cascade
- Evaluation
 - Detectors: Intersection over union, precision recall
 - Classifiers: Confusion matrix

Recall: Examples of kernel functions

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

- Kernels go beyond vector space data
- Kernels also exist for "structured" input spaces like sets, graphs, trees...

Discriminative classification with sets of features?

- Each instance is unordered set of vectors
- Varying number of vectors per instance

Slide credit: Kristen Grauman

Partially matching sets of features

Optimal match: $O(m^3)$
 Greedy match: $O(m^2 \log m)$
 Pyramid match: $O(m)$

$X = \{\bar{x}_1, \dots, \bar{x}_m\}$ $Y = \{\bar{y}_1, \dots, \bar{y}_n\}$ ($m = \text{num pts}$)

$\min_{\pi: X \rightarrow Y} \sum_{x_i \in X} \|x_i - \pi(x_i)\|$ hate matching kernel that makes it practical to compare large sets of features based on their partial correspondences.

[Previous work: Indyk & Thaper, Bartal, Charikar, Agarwal & Varadarajan, ...]

Slide credit: Kristen Grauman

Pyramid match: main idea

Feature space partitions serve to "match" the local descriptors within successively wider regions.

descriptor space

$X = \{\bar{x}_1, \dots, \bar{x}_m\}$ $Y = \{\bar{y}_1, \dots, \bar{y}_n\}$

Slide credit: Kristen Grauman

Pyramid match: main idea

$\mathcal{I}(H_X, H_Y) = \sum_j \min(H_X(j), H_Y(j)) = 3$

Histogram intersection counts number of possible matches at a given partitioning.

Slide credit: Kristen Grauman

Pyramid match

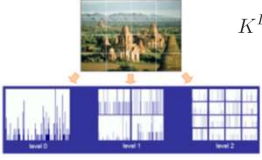
$$K_{\Delta}(X, Y) = \sum_{i=0}^L 2^{-i} \left[\underbrace{\mathcal{I}(H_X^{(i)}, H_Y^{(i)})}_{\text{measures difficulty of a match at level } i} - \underbrace{\mathcal{I}(H_X^{(i-1)}, H_Y^{(i-1)})}_{\text{number of newly matched pairs at level } i} \right]$$

- For similarity, weights inversely proportional to bin size (or may be learned)
- Normalize these kernel values to avoid favoring large sets

[Grauman & Darrell, ICCV 2005] Slide credit: Kristen Grauman

Spatial pyramid match

- Make a pyramid of bag-of-words histograms.
- Provides some loose (global) spatial layout information




$$K^L(X, Y) = \sum_{m=1}^M \kappa^L(X_m, Y_m)$$

Sum over PMKs computed in *image coordinate space*, one per word.

[Lazebnik, Schmid & Ponce, CVPR 2006]

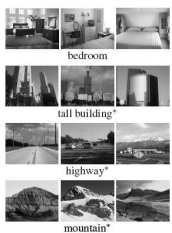
Spatial pyramid match

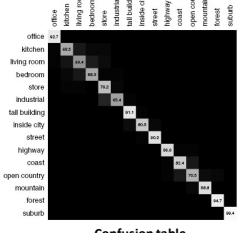
- Can capture **scene** categories well---texture-like patterns but with some variability in the positions of all the local



Spatial pyramid match

- Can capture **scene** categories well---texture-like patterns but with some variability in the positions of all the local pieces.
- Sensitive to global shifts of the view





Confusion table

Summary: This week

- Object recognition as classification task
 - Boosting (face detection ex)
 - Support vector machines and HOG (person detection ex)
 - Pyramid match kernels
 - Hoggles visualization for understanding classifier mistakes
 - Nearest neighbors and global descriptors (scene rec ex)
- Sliding window search paradigm
 - Pros and cons
 - Speed up with attentional cascade
- Evaluation
 - Detectors: Intersection over union, precision recall
 - Classifiers: Confusion matrix
