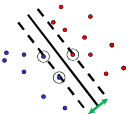


Support vector machines and kernels

Thurs April 19
Kristen Grauman
UT Austin



Last time

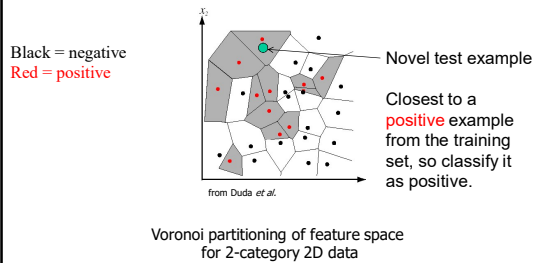
- Sliding window object detection wrap-up
 - Attentional cascade
 - Applications / examples
 - Pros and cons

Today

- Supervised classification continued
 - Nearest neighbors
 - 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

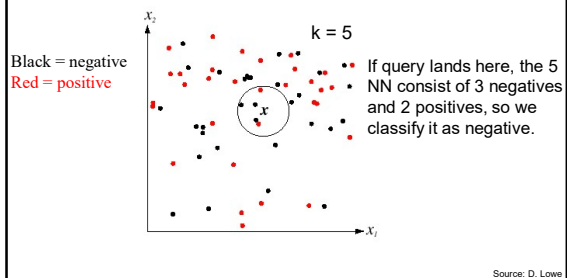
Nearest Neighbor classification

- Assign label of nearest training data point to each test data point

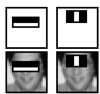


K-Nearest Neighbors classification

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



Three case studies



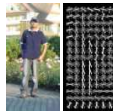
Boosting + face detection

Viola & Jones



NN + scene Gist classification

e.g., Hays & Efros



SVM + person detection

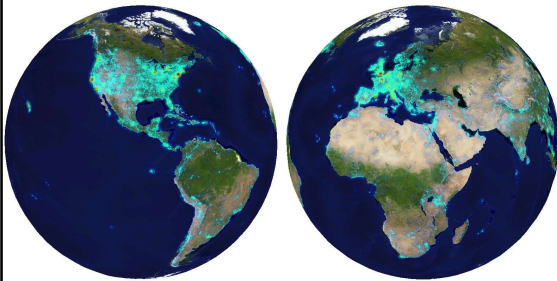
e.g., Dalal & Triggs

Where in the World?



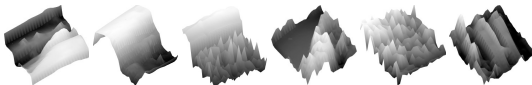
[Hays and Efros. **im2gps**: Estimating Geographic Information from a Single Image. CVPR 2008.]

6+ million geotagged photos
by 109,788 photographers



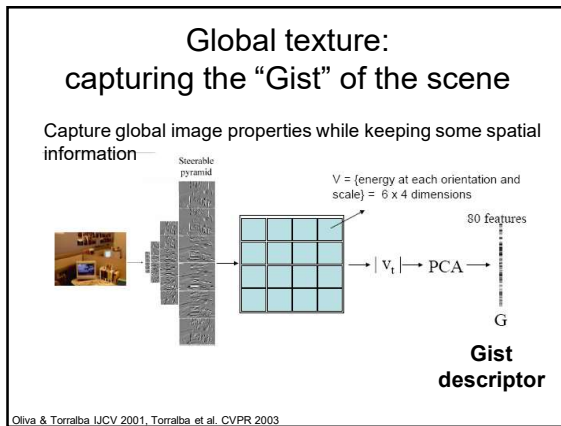
Annotated by Flickr users

Spatial Envelope Theory of Scene Representation Oliva & Torralba (2001)



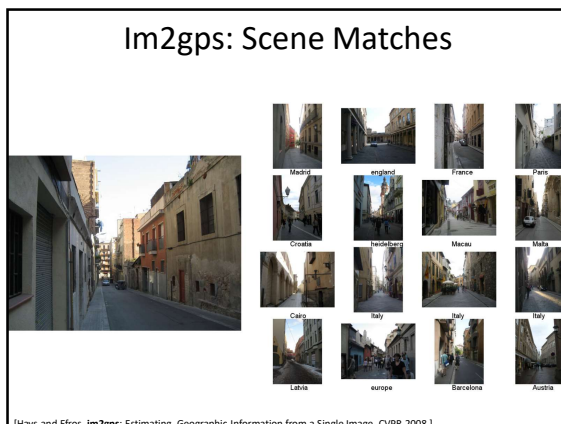
A scene is a single surface that can be
represented by global (statistical) descriptors

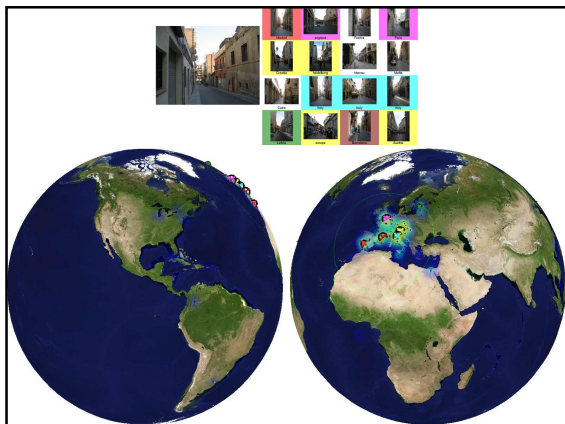
Slide Credit: Aude Oliva

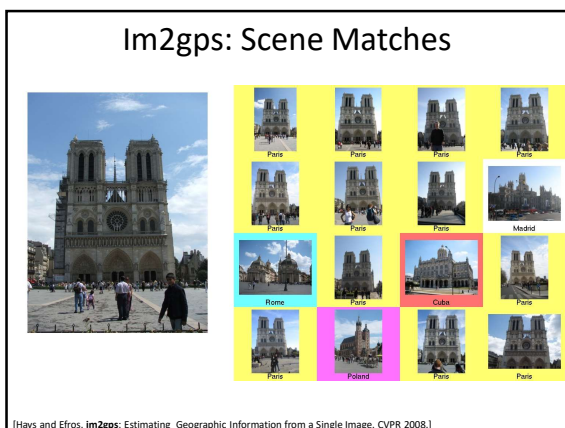


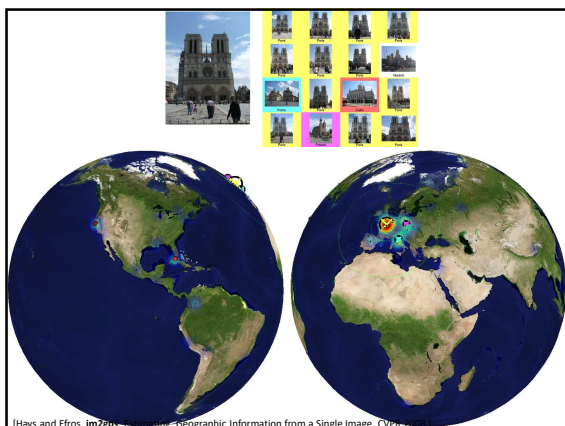
Which scene properties are relevant?

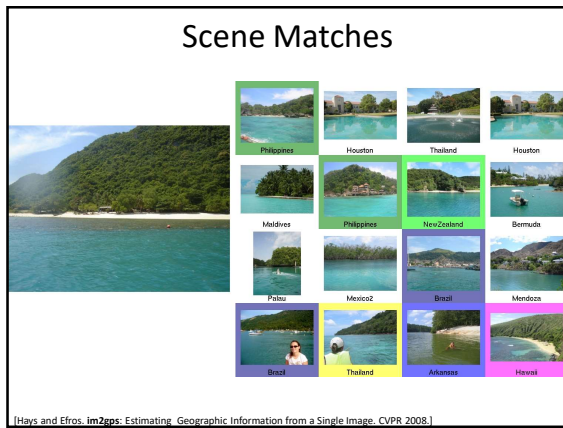
- **Gist scene descriptor**
- **Color Histograms** – $L \times A \times B \times 4 \times 14 \times 14$ histograms
- **Texton Histograms** – 512 entry, filter bank based
- **Line Features** – Histograms of straight line stats

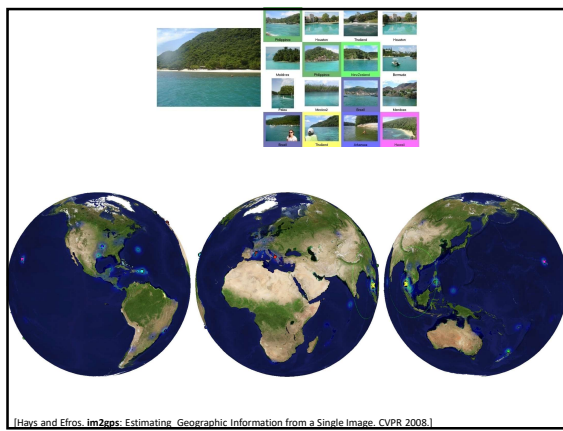


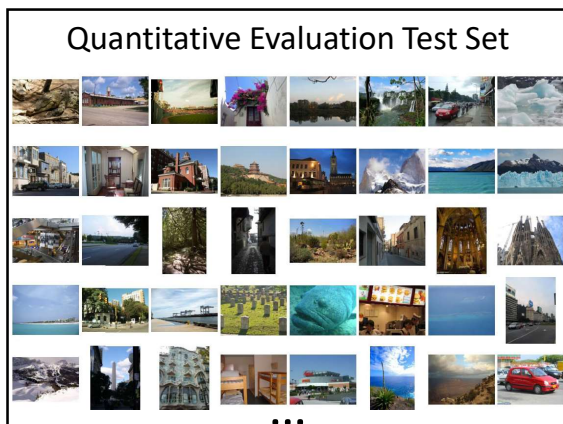


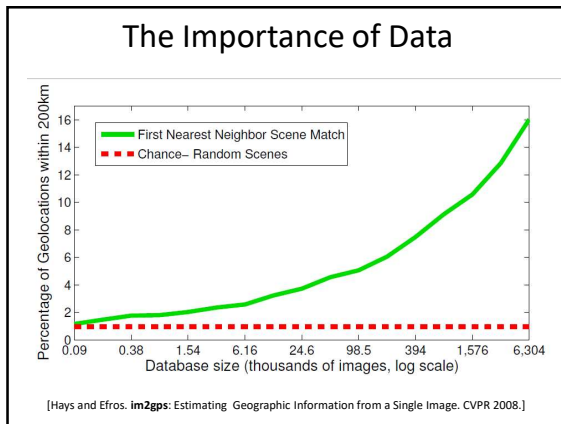










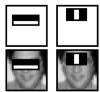


Nearest neighbors: pros and cons

- **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

Three case studies



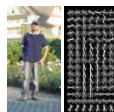
Boosting + face detection

Viola & Jones



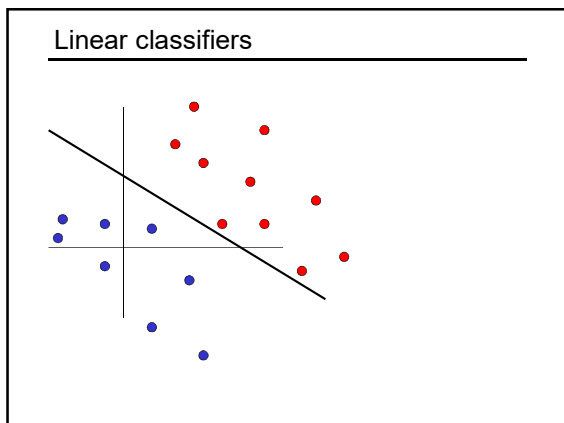
NN + scene Gist classification

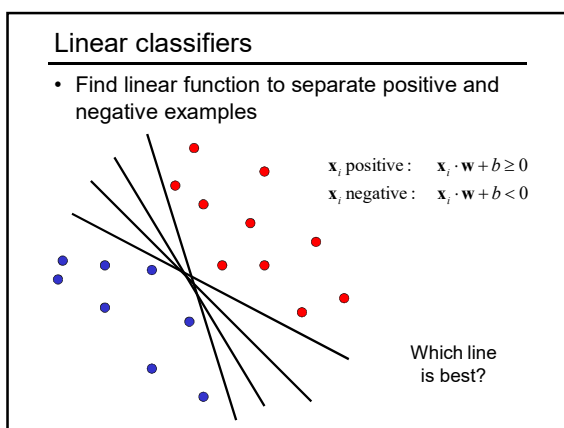
e.g., Hays & Efros

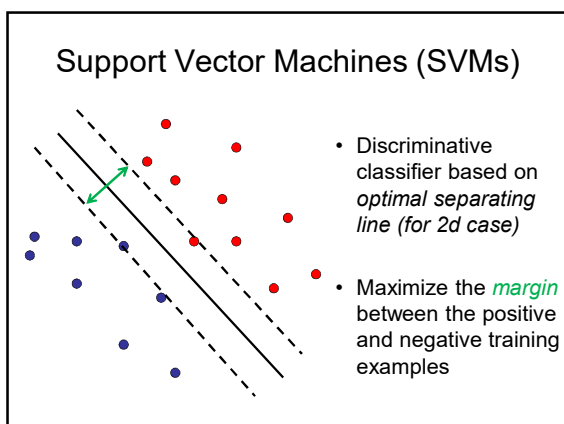


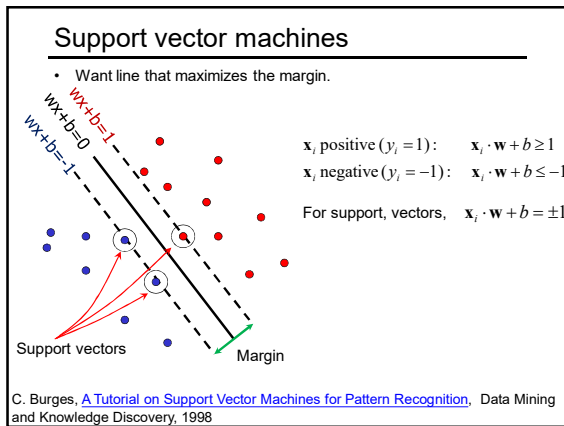
SVM + person detection

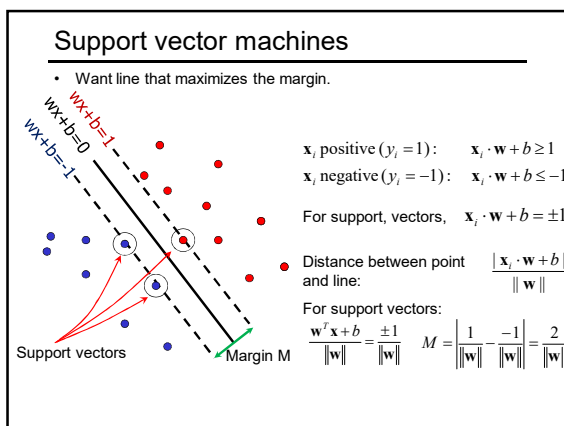
e.g., Dalal & Triggs











Finding the maximum margin line

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

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

Quadratic optimization problem:

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

Subject to $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$

Finding the maximum margin line

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

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:

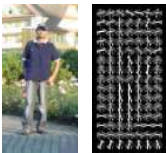
$$f(x) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

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

*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,

Person detection with HoG's & linear SVM's



- Histograms of oriented gradients (HoG): Map each grid cell in the input window to a histogram counting the gradients per orientation.

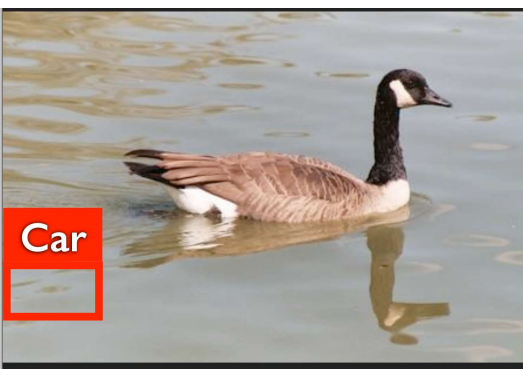
- Train a linear SVM using training set of pedestrian vs. non-pedestrian windows.

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.inria.fr/pubs/2005/DT05/>

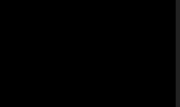


Carl Vondrick <http://web.mit.edu/vondrick/ihog/slides.pdf>


What information does HOG have?

Image






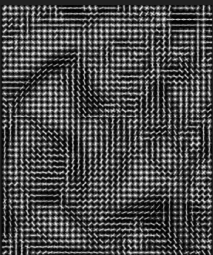
HOG



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

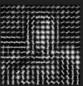
What information is lost?



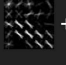


HOGgles: Visualizing Object Detection Features

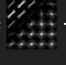
Method: Paired Dictionary




$= \alpha_1$



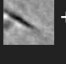
$+ \alpha_2$



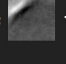
$+ \dots + \alpha_k$



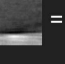
α_1




$+ \alpha_2$



$+ \dots + \alpha_k$

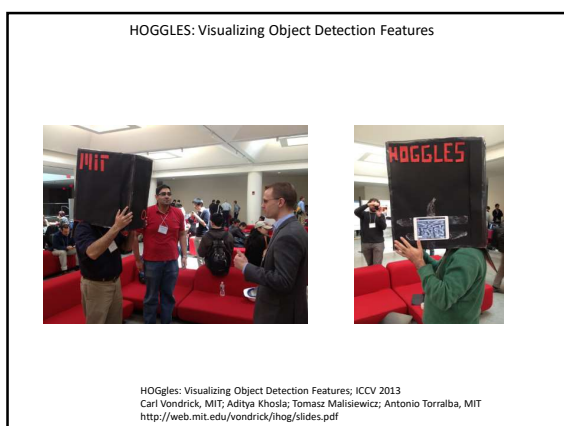


$=$



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





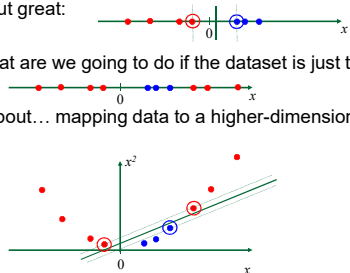


Questions

- What if the data is not linearly separable?

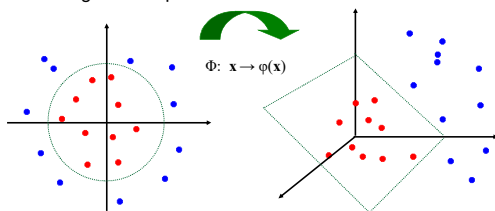
Non-linear SVMs

- 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 $\phi(\mathbf{x})$, define a **kernel function** K such that

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

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

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

“Kernel trick”: 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 $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$:

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

Examples of kernel functions

- Linear: $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$

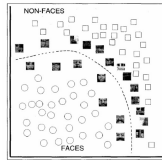
- Gaussian RBF: $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$

- Histogram intersection:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \sum_k \min(x_i(k), x_j(k))$$

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 Li and Leibe

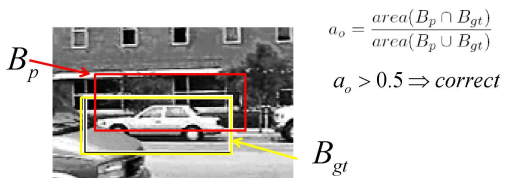
Scoring a sliding window detector



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

Kristen Grauman

Scoring a sliding window detector



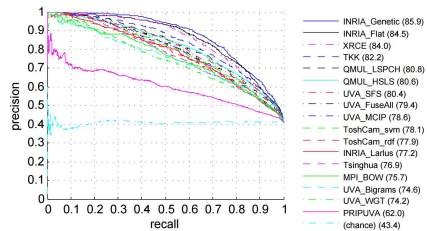
$$a_o = \frac{\text{area}(B_p \cap B_{gt})}{\text{area}(B_p \cup B_{gt})}$$

$$a_o > 0.5 \Rightarrow \text{correct}$$

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

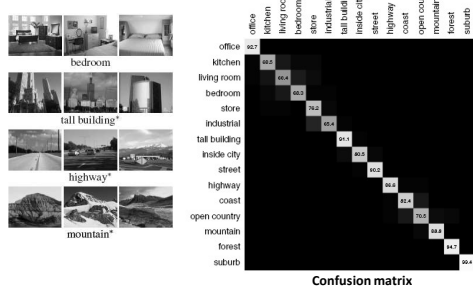
Scoring an object detector



- Detector produces *confidence scores*; plot precision vs. recall as a threshold on the confidence is varied.
- **Average Precision (AP)**: mean precision across recall levels.

Scoring a multi-class classifier

- Confusion matrix records probability of calling class i class j



Summary: This past week

- Object recognition as classification task
 - Boosting (face detection ex)
 - Support vector machines and HOG (person detection ex)
 - 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
 - Object proposals as alternative search
- Evaluation
 - Detectors: Intersection over union, precision recall
 - Classifiers: Confusion matrix