# lecture 05 : Principles of sliding window detectors

What we'd like to do?

- Visual scene understanding
- What is in the image and where
- Object categories, identities, properties, activities, relations, ...

## Things vs Stuff

- Thing : object with specific shape
- Stuff : material defined by homogeneous or repetitive pattern. Has no specific spatial extent or shape
  $\hookrightarrow$ Segmentation method

## Recognition tasks

- Image classification : does the image contain an aeroplane?
- Object class detection / localization : where are the aeroplanes if any?
- Object class segmentation : which pixels are part of an aeroplane?
- Panoptic segmentation : Besides object segmentation, also background segmentation

## Challenges

- Background clutter : lot of stuff in the background
- Occlusions and truncation : partially seen objects from the camera
- Intraclass variation : a class can have a lot of different versions

  $\hookrightarrow$ object instant recognition : recognizes a specific model of an object

  $\times$ not its generic class

  category detection : recognizes generic classes
  it is harder to perform than instant detection

## So why bother?

- Spatial relationships for image understanding and retrieval ("a cat riding a skateboard")
- Visual question and answering : object grasping / tracking
  (collition prevention, face recognition

Sliding window detectors

Problem of background clutter : Solution

- Use sub window :
    - At correct position, no clutter is present
    - Slide window to detect objects
    - Change size of the window to search over scales

Detection by classification

- Basic component : binary classifier
    ↳ sliding window over window using CNNs is too slow
- Detect objects in clutter by searching
    - sliding window : exhaustive search over position and scale
        ↳ in practice, it is possible to use same window size over spatial pyramid
        ↳ more efficient

Window (image) classification

- Features usually engineered
- Classifier learned from data

$$\boxed{\text{image}} \longrightarrow \boxed{\begin{array}{c}\text{feature}\\\text{extraction}\end{array}} \longrightarrow \underset{x}{\left[\begin{array}{c}\vdots\\\vdots\end{array}\right]} \longrightarrow \boxed{\begin{array}{c}\text{classifier}\\ F(x)\end{array}} \longrightarrow \begin{array}{c}\text{Car/Non Car}\\ P(c\,|X) \propto F(x)\end{array}$$
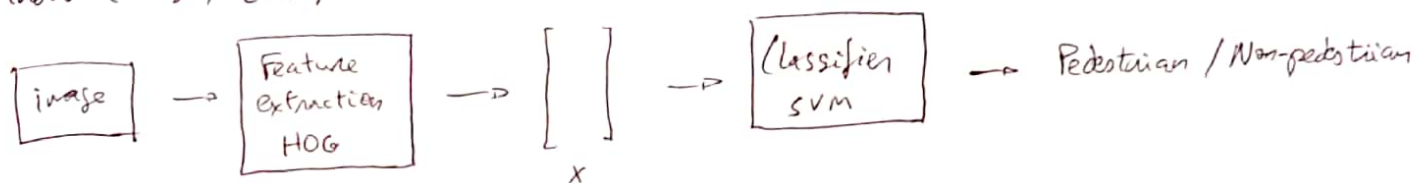
Problems with sliding windows

- Aspect ratio
- Granularity (finite grid)
- Partial occlusions
- Multiple responses ⟶ 
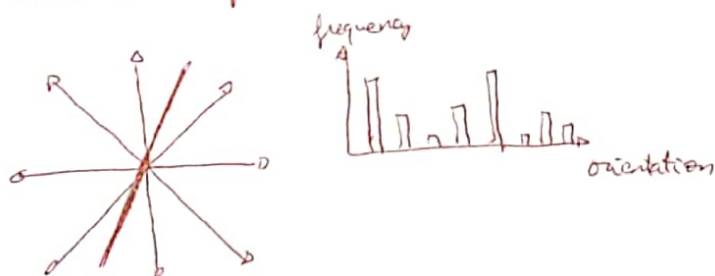
(Conference on CV and PR)

- Objective: detect (localize) standing humans in an image

- Sliding window classifier

- Train a binary SVM classifier on whether a window contains a standing person or not

- Histogram of Oriented Gradients (HOG) feature

- Although introduced for pedestrian detection, it has been successfully used with many other categories

## Window (image) classifier

image → Feature extraction HOG → [ ] _X_ → Classifier SVM → Pedestrian / Non-pedestrian

## Feature : HOG

- Tile 64 × 128 pixel window into 8 × 8 pixel cells

- Calculate gradient image (edge detection)

- Each cell represented by by histogram over 8 orientation bins (or sector)

frequency

orientation

- Adds a second level of overlapping spatial bins re-normalizing orientation histogram over larger spatial area

- Feature dimensions (approx) = $16 \times 8$ (for tiling) × 8 (orientations) × 4 (blocks) = 4096

- Similarity to CNN
  { CNN learns the filters automatically
  - Sum pooling
  - normalization

- OK job

# Linear classifier

- $f(x) = w^T x + b$
- It learns such weights $w$ and bias $b$ that

$$f(x_i) = \begin{cases} \geq 0 & y_i = 1 \\ < 0 & y_i = -1 \end{cases}$$

2D discriminant is a line
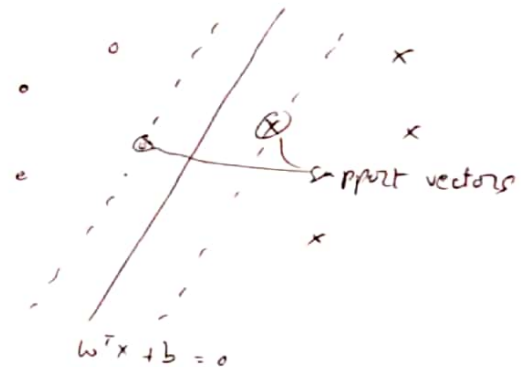3D " is a plane

# Linear separability

- The points might be linearly separable but with very narrow margin
- The large margin solution might be better, even constraints are violated
- In general there is a trade off between the margin and the number of mistakes on the training data

# Support Vector Machine

. It is a way to optimize this trade off
- Find a good trade off between the classification margin and the missclassified training points

$$f(x) = w^T x + b$$

$$\min_{w \in \mathbb{R}^d} \|w\|^2 + C \sum \max(0, 1 - y_i f(x_i))$$



support vectors

$w^T x + b = 0$

# Learned Model using HOG detector

└→ evidence it is a person

- Positive weights
- Negative weights

} Average over positive training data

└→ evidence it is not a person

$$f(x) = w^T x + b$$

└→ comes from $8 \times 8$

# what do negative weights mean?

$$W^T x > 0$$
$$(W_+ - W_-) x > 0$$
$$W_+ x > W_- x$$

- Complete system compete pedestrian / pillar / doorway models
- Discriminative models come with own background model
- Avoid detections on doorways by penalizing vertical edges

$$\text{Pedestrian model} > \text{Pedestrian background model}$$

## Problems when training a sliding window detector

- Inherently asymmetric problem: many more "non-object" than "objects"
- Classifier needs to have very low false positive rate
- "Non-object" category is very complex and needs a lot of data

## Optimizing approach: Bootstrapping

- 1. Pick a negative training set at random
- 2. Train classifier
- 3. Run on training data
- 4. Add false positives to training set
- 5. Repeat from step 2. (Retraining)

Data augmentation: With your available data you

- Flip
- Rotate
- Scale
- Crop
- Translate
- Apply Gaussian noise

} Jittered positive
$$\Uparrow$$
to the images
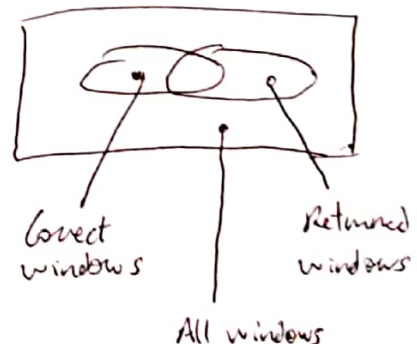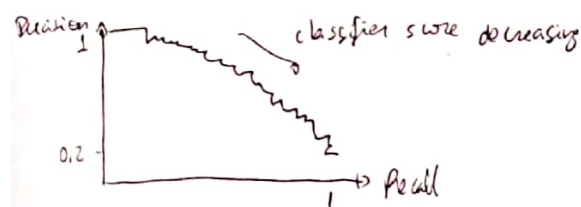$$\Downarrow$$

## !Window (image) first stage classification

Jittered positive and random negative $\rightarrow$ Feature extraction HOG $\rightarrow$ $\begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$ $x$ $\rightarrow$ linear SVM classifier $f(x) = w^T x + b$

## Precision - Recall curve

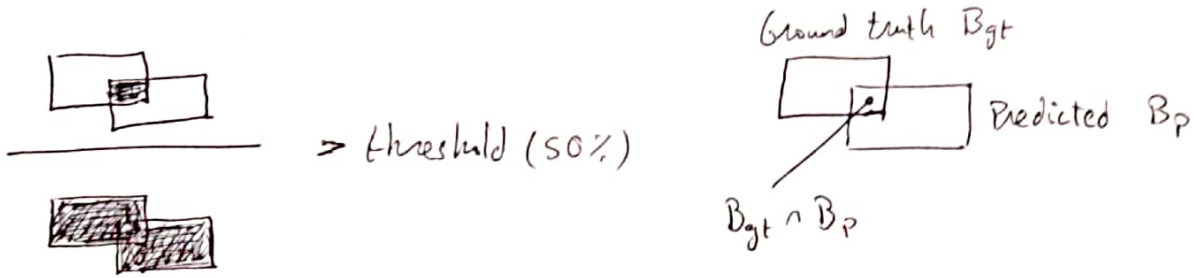- Precision: % of returned window that are correct
- Recall: % of correct windows that are returned



Precision / Recall curve (classifier score decreasing), axis Precision (1 to 0.2), Recall

Correct windows | Returned windows | All windows

# Evaluating the detected bounding boxes

- Area of overlap (AO) measure
- Correct detection if intersection over union larger than threshold



$$> \text{threshold } (50\%)$$

Ground truth $B_{gt}$

Predicted $B_p$

$B_{gt} \cap B_p$

$$AO(B_{gt}, B_p) = \frac{|B_{gt} \cap B_p|}{|B_{gt} \cup B_p|}$$

First training phase

Precision



0.4

0.2

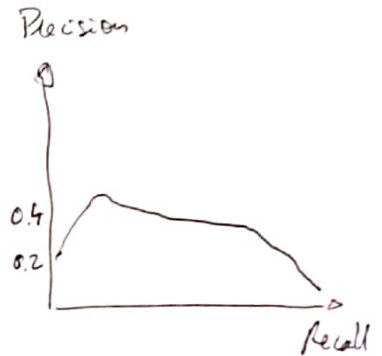Recall

Second training phase → Retrain using better data

- Find high scoring false positive detections
- Use them as hard negatives for next training round
- Cost = # training image · inference time per image

~~First training phase~~

## Accelerating sliding window search

- Sliding window search is slow since some many windows are needed
- $m \times n \times$ scale = 100 000 windows for $320 \times 240$ image
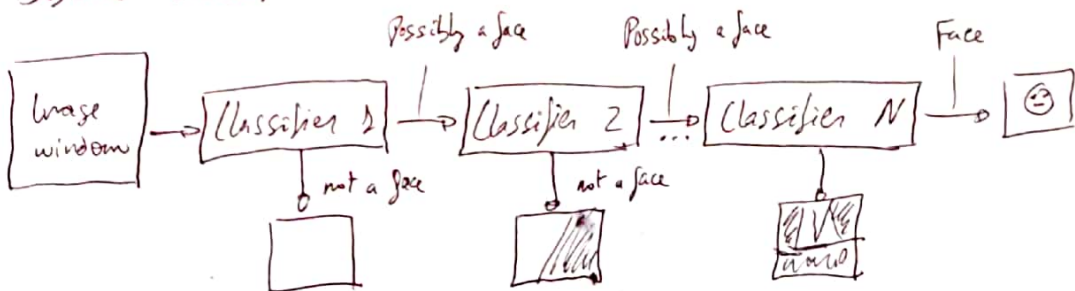- Most windows are negative
- It is possible to speed up the search

Second training phase

Precision



0.8 Recall

## Cascade classification



Possibly a face    Possibly a face    Face

Image window → Classifier 1 → Classifier 2 → ... → Classifier N → 😊

not a face    not a face

More complex, slower, lower false positive rate →

In the end, you spend much less computational cost

- Slow and expensive classifiers only applied to a few windows
- Controlling complexity vs speed: numbers of features, numbers of parts, ...

Scanned by CamScanner

Detection proposal : Hierarchically clustering superpixels

- larger homogeneous area is considered a unit in the image → superpixel

- Hierarchical segmentation : start with small and merge using cues

  Produces roughly 2000 regions per image with 95% of hitting relevant objects

Things to remember

- Detection by sliding window classification → Concept and components
- Multiple scales (and aspect ratios) to detect objects of different size
- Importance of hard negative mining (due to class imbalance)
- Cascade detectors → speed up inference
- Speed up training and inference by selecting sub-set of windows only