

# Interactively Test Driving an Object Detector: Estimating Performance on Unlabeled Data

Rushil Anirudh and Pavan Turaga  
School of Electrical, Computer and Energy Engineering  
School of Arts, Media and Engineering  
Arizona State University, Tempe.  
{ranirudh, pturaga}@asu.edu

## Abstract

In this paper, we study the problem of ‘test-driving’ a detector, i.e. allowing a human user to get a quick sense of how well the detector generalizes to their specific requirement. To this end, we present the first system that estimates detector performance interactively without extensive ground truthing using a human in the loop. We approach this as a problem of estimating proportions and show that it is possible to make accurate inferences on the proportion of classes or groups within a large data collection by observing only 5 – 10% of samples from the data. In estimating the false detections (for precision), the samples are chosen carefully such that the overall characteristics of the data collection are preserved. Next, inspired by its use in estimating disease propagation we apply pooled testing approaches to estimate missed detections (for recall) from the dataset. The estimates thus obtained are close to the ones obtained using ground truth, thus reducing the need for extensive labeling which is expensive and time consuming.

## 1. Introduction

Object detection is one of the most popular and practical applications of computer vision today. Even though it has been extensively studied for over two decades, it continues to be a challenging problem. The difficulty in solving the problem is attributed to several factors such as change in lighting, pose, viewpoint, size, shape and texture. However, there have been significant advances in building robust object detectors that are finally making them commercially viable in applications for home security, surveillance, online social networks, online shopping etc. In the near future, it is conceivable that even individual end-users will wish to pur-

This work was supported in part by the NSF-CCF-CIF award 1320267, ONR Grant N00014-12-1-0124 subaward Z868302 and by ASU startup grant.

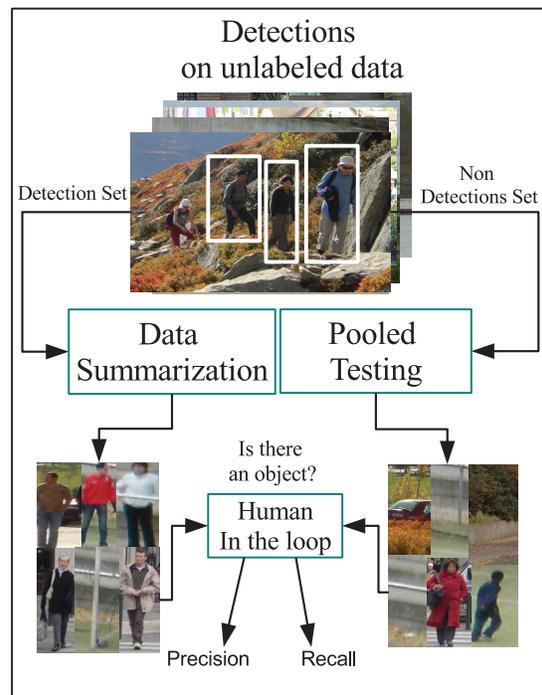


Figure 1: Our system estimates the performance of an object detector on unlabeled data using feedback from a human in the loop.

chase specific detectors for their specific use. For example VMX [2], is being developed as a “computer vision as a service” that allows one to outsource detectors and classifiers while focusing on their applications. In such a scenario, how does one allow a user to get a quick sense of how well the detector generalizes to their specific requirement, as opposed to interpreting the ROC curves of detectors on certain test-beds when the test-beds themselves may or may not be reflective of a user’s intended deployment scenario. Further, there has been an explosion of the number of detectors in the past decade promoted by challenges such as PASCAL VOC [12] where in the past three years alone, 43 methods com-

pleted the challenge (2010 – 22, 2011 – 13, 2012 – 8) [1]. Traditionally, rigorous evaluation of detectors requires extensive ground-truthing, which by its very nature restricts an end-user from judging how well the detector will generalize outside the test-bed. Theoretical bounds offered by statistical learning theory [23] offer guidance to algorithm developers, however the bounds themselves are generally quite weak and provide little useful information to end-users. The question we seek to answer in this paper is, how does one get a quick sense of a detector’s performance for a specific application without the need for extensive groundtruthing?

In this paper, we propose an interactive system by approaching this as a problem of smart sampling to estimate object/non-object proportions, a block diagram of the system is shown in fig 1. Let  $U$  be the set of all possible image patches among user provided images. Denote by  $D \subset U$ , a ‘detection-set’ which were reported as containing the object of interest by the detector and  $\bar{D}$ , the complement of the detection set. The challenge is to then sample informatively from  $D$  and  $\bar{D}$  in a way that quickly allows us to estimate class proportions. Proportions are useful because they are independent of the size of dataset, and two of the important measures in detection theory, precision and recall are expressed as fractions. So given a large collection of unlabeled images, which are the best samples to pick that preserve “class proportions”? The answer is not trivial, random sampling methods are accurate only when class proportions are nearly equal. Assuming this step is solved, the next step is to ground truth the chosen samples by asking a simple yes/no question to a human, based on which we can estimate the proportion of objects within the set of detections. Performance for a detector is always measured relative to ground truth, therefore in the absence of labels we seek to estimate “perceived performance” and show that this is indeed close to actual performance measured with ground truth when using the proposed sampling procedures. We adopt and extend data sampling strategies and present a coherent framework to estimate false positives (incorrect detection) and a probabilistic estimation approach known as group testing to determine the false negatives (missed detection). We treat the detector as a black box and operate directly on the detections.

The main idea is the following, first for different detection thresholds ( $\gamma$ ) we feed the set of detections to a data summarization algorithm that identifies a small number of samples based on preserving some notion of information. These samples are generated for each  $\gamma$  and presented to the human who gives them a 1/0 label for object/non-object. If the samples are chosen well, the class proportions are expected to be preserved giving us an accurate estimate of precision (portion of total detections that were true).

Recall is the other important measure, which is the portion of total objects that were detected. Once we estimate

precision, we only need to estimate the number of false negatives from  $\bar{D}$  to obtain recall. Unfortunately obtaining the complimentary set  $\bar{D}$  is not easy and it is further complicated by the fact that the objects are very similar to background in the feature space. This results in the failure of any feature based approaches such as summarization and clustering. To address this issue we use the theory of statistical group testing that is employed to study estimation of disease propagation among large populations of plants, animals or humans. Essentially, we pool in multiple images before ‘testing’ them with the human for the presence of an object. This information is used to obtain estimates on the number of false negatives within  $\bar{D}$ .

**Contributions** 1) We present the first system that allows a user to estimate detector performance without the need for extensive groundtruthing. 2) We achieve this by employing probabilistic techniques such as pooled testing and smart sampling methods such as summarization by keeping a human in the loop. This will allow users to “test-drive” object detectors in the future while minimizing time and effort required for labeling. 3) Our system also gives the best performance estimate compared to simpler baselines with the least amount of human effort.

## 2. Relevant Work

In this section we outline the important work that has been done in the different fields of research relevant to our system.

**Object detection** is the problem of putting a bounding box around an object within an image. Generally speaking, it involves running a window over the image at different scales to check for the highest match with some model representation of the object obtained from training data. In this work we look at pedestrian detection with three detectors, namely integral channel features (ChnFtrs) [9], Multiple features with color self similarity (MultiFtrs+CSS) [25] and aggregate channel features (ACF)[8], all of which have performed well on many standard datasets. For a detailed comparison of different pedestrian detectors on several datasets we refer the reader to [10]. There are several variations of each detector mainly varying in the choice of local descriptors or using extra information such as color, context etc. A popular feature used in most detectors is Histogram of Oriented Gradients (HOG) feature introduced by the Dalal-Triggs person detector [6]. Our system operates on the detections so it is applicable to any kind of object detector.

*Determining true positives with labels:* Once a bounding box has been predicted, an overlap measure determines if the detection is a true or false positive based on its overlap area with the ground truth. Like any detection problem, the main quantities that determine the quality of a detector are the number of false positives and false negatives. An ideal detector would have both quantities to be low over various

detection thresholds.

**Video summarization** attempts to identify the most concise representation of a video while maximizing ‘information’ content. Usually, algorithms for summarization differ in their definition of information keeping the main application in mind. Recent methods in summarization have employed high-level contextual features and user input to present summaries[16]. The summarization algorithm we seek for this system is a more general form of clustering and must be able to deal with image based features rather than video based ones such as flow, object tracks etc. *Precis* [20, 21] is one such algorithm that can work with image based features. *Precis* looks at maximizing ‘diversity’ and minimizing representational error among the data points.

**Human in the loop systems** Many difficult problems in vision including segmentation, summarization, and even learning have been shown to benefit by human-in-the-loop systems. The motivation for such systems is to actively minimize human effort in tasks that are otherwise laborious such as image or video annotations, active video segmentation etc [24]. Typically a human is asked a set of questions, the answers to which are used as feedback to the system to reduce the search space. For example, Branson et al. [3] develop a system that performs multi-object classification and uses pre-processing to identify a small set of important questions to ask the human, which then enables learning. This allows the authors to perform classification of subtly different classes which would have been very hard otherwise. Since we are working with unlabeled data, the role of the human is critical in determining the final estimated performance. Humans also have an extraordinary ability to identify objects using multiple cues [22], which can be taken advantage of for image pooling.

**Group testing** Pooled or group testing has its origins in the area of groundtruthing agricultural/biological samples, and its mathematical formulation has since been applied more broadly. The underlying premise is that individuals (blood samples, seed kernels etc.) are first pooled together into groups and the groups are tested for the presence of infection. When infection rate is low, there is an overwhelming evidence that pool testing can confer substantial benefits in terms of testing/labeling effort when compared to testing individually [14, 4]. A complete discussion of the assumptions, robustness of the estimates as a function of group size, errors in testing etc. can be found in [4]. These ideas have been applied to estimating disease propagation in large populations of plants, animals or humans [14].

**Image retrieval from large datasets** The techniques described here are very similar to methods that search for similar images within large data collections. For example hashing methods such as geometric min-hash [5], and tree based methods such as [17] could potentially be used to find the number of true detections in  $D$ , but this requires one to have

trained examples to learn the underlying data structure and we set out to work with unlabeled data with no prior training. It doesn’t help much for the negative set either since the feature descriptors are expected to be poor due to occlusions, shadows etc., which caused the detector to miss them the first time.

### 3. Approach & Implementation

In this section we will describe our system in detail and the theory behind estimating detector performance. The system presumes that a detector is used off-the-shelf, and is completely agnostic to the detector’s algorithm or features. A note on terminology used here, the term ‘class’ has no meaning in unlabeled data, but we refer to it as the label that would have been assigned by a human. Therefore, our estimate of quantities are ‘perceived’ but as we show, they are close to the true underlying values.

**Measuring Performance:** Knowledge of the number of false positives (FP) and false negatives (FN) provides sufficient information in most cases to judge the performance of a detector. Some applications may require the false positives per image (FPPI) to be restricted below a specification before implementing it into their system, but our current system does not allow us to estimate FPPI, so we will focus on the two main quantities FP, FN.

#### 3.1. Sampling from the detections set to estimate perceived precision

Interactively estimating the proportions of different classes (objects, non-objects here) from a large set of images essentially boils down to smart sampling, i.e. the problem would be solved if we were to pick the right proportion of samples from each class. Picking random samples is of little help in cases with skewed class proportions i.e., when the detector is operating at high precision. Given the fact that any detector based application is commercially viable only when the detectors are mature enough to have high precision and recall, the random sampling method is not suitable. One could use clustering algorithms such as K-means (K-medoids) and use the cluster centers as the samples, but this suffers from the same pitfalls as random sampling and does only marginally better.

Let the set of detections at a specific threshold  $\gamma$ , be denoted as  $D$ , define  $D_{fp}, D_{tp} \subset D$  as the set of false and true detections respectively and  $D_{fn} \subset \bar{D}$  as the set of missed detections or false negatives. Then  $|D_{fp}| \ll |D_{tp}|$  in most cases, where  $|\cdot|$  denotes the cardinality of a set. How does one sample informatively from  $D_{tp}$  &  $D_{fp}$  when they are not explicitly known? The first step is to extract a feature vector from each image patch, we use the histogram of oriented gradients for this purpose. The problem now is to identify the  $K_\gamma (\ll |D|)$  best samples, denoted by  $S$  from the data feature matrix  $X$  such that the proportion of

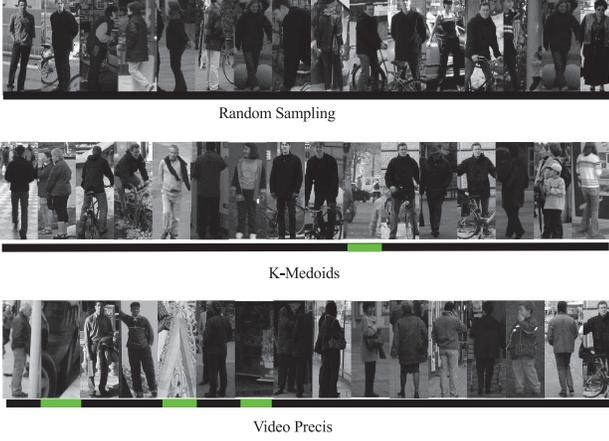


Figure 2: Random sampling and K-medoids overestimate the proportion of pedestrians because they fail to retrieve samples from the smaller class of non-objects. Precis’[21] samples give a much closer estimate to the ground truth. Here the non-objects are highlighted (best seen in color) for a fixed  $\gamma$  from the INRIA pedestrian dataset [6], the perceived precision for random, K-means, Precis are 1, 0.9375, 0.8125 respectively compared to the actual precision = 0.7.

elements in  $D_{tp}$  and  $D_{fp}$  are preserved. Since the number of false positives is far smaller than the number of true positives using algorithms such as K-medoids will result in picking centers that are mostly from  $D_{tp}$  since it only looks to optimize representational error, resulting in poor proportion estimates. Instead we use Precis [20, 21] to pick samples well distributed across each set. Precis solves for a set of samples  $S = \{F_i | i = 1, 2 \dots K_\gamma\}$ , that minimize representational error (‘coverage’) denoted by  $rep(S, X)$  and maximize diversity between chosen samples, denoted by  $div(S, X)$ . These quantities are given by

$$rep(S, X) = tr \left[ \sum_i \sum_{F_k \in V_i} (F_k - F_i)(F_k - F_i)^T \right], \quad (1)$$

$$div(S, X) = tr \left[ \sigma_i(F_i - \bar{F})(F_i - \bar{F})^T \right], \quad (2)$$

where  $V_i$  is the partition corresponding to  $F_i$ , i.e.  $V_i$  are the elements of  $X$  that are closer to  $F_i$  than any other element in  $S$ .  $\bar{F} = \frac{1}{K_\gamma} \sum_j F_j$  is the mean of  $S$ . This results much better sampling across  $D_{tp}$  &  $D_{fp}$  as shown in fig 2.

**Homogenizing the detection set** In general, samples chosen from  $D_{fp}$  are not well behaved in the feature space because they contain everything but the object of interest. This can throw off any metric based algorithm such as K-medoids due to outliers and Precis due to increased diversity etc., resulting in poor samples. To ensure better sample selection, we transform the data to a space where  $D_{fp}$  and  $D_{tp}$  are well separated. This transform, denoted by  $\mathcal{T}$ ,

needs to be estimated only once per object class and is obtained by minimizing a cost function using data points from  $O_T$ , a set of image patches with objects and  $O_F$  a set of image patches without objects. Since both of these sets are unknown for unlabeled data, we will approximate them by picking  $O_T = D_{\gamma_H}$  which is the detection set at a high confidence threshold  $\gamma_H$  and  $O_F = \bar{D}_{\gamma_L}$  the non-detection set at a low confidence threshold  $\gamma_L$ . This is a reasonable approximation since at  $\gamma_H$  one expects very few false positives and at  $\gamma_L$  one expects very few false negatives.

Typically, we want to ensure that  $d_{\mathcal{T}}(x_i, x_j) \leq u$ , for  $(x_i, x_j) \in O_T$  and  $d_{\mathcal{T}}(x_i, x_j) \geq \ell$  for  $x_i \in O_T$  &  $x_j \in O_F$  for some  $u, \ell > 0$ . The loss function(s) depend mainly on the Euclidean norm on the transformed space,  $X^T \mathcal{T}^T \mathcal{T} X$ . An example of such constraints are

$$c_1 = \max(0, d_{\mathcal{T}}(x_i, x_j) - u) \text{ for } (x_i, x_j) \in O_T, \quad (3)$$

$$c_2 = \max(0, \ell - d_{\mathcal{T}}(x_i, x_j)) \text{ for } x_i \in O_T \text{ \& } x_j \in O_F. \quad (4)$$

Other loss functions can be used depending on the application. A constraint on  $\mathcal{T}$  is imposed by a regularizing condition  $reg(\mathcal{T})$  (such as  $\mathcal{T} \geq 0$ ). A general cost function for  $m$  loss functions is given as follows:

$$J_1(\mathcal{T}) = reg(\mathcal{T}) + \lambda \sum_i^m c_i. \quad (5)$$

This is very similar to the metric learning problem and other examples of loss functions and regularizers can be found in [15]. In this paper, we use Information Theoretic Metric Learning (ITML) [7] to find the optimal  $\mathcal{T}$ . The estimated transform is independent of the particular samples and is estimated once per object class. Once we have the optimal  $\mathcal{T}$ , we solve for the best set of samples,  $S$  using Precis which minimizes the cost function:

$$J_2(S) = \alpha rep(S, \mathcal{T}X) + (1 - \alpha) div(S, \mathcal{T}X). \quad (6)$$

Finally, the set of samples are shown to a human who marks false positives. Typically  $K_\gamma$  is chosen to be about 5 – 10% of the detections at  $\gamma$ . From this, estimated precision is calculated as :  $\hat{P}_\gamma = 1 - \frac{\widehat{FP}_\gamma}{K_\gamma}$  where  $\widehat{FP}_\gamma$  is the perceived number of false positives at  $\gamma$ .

### 3.2. Pooled testing of samples from the detection complement

Metric based approaches such as clustering, Precis and random sampling severely underestimate the number of false negatives due multiple reasons: 1)The proportions of objects to non-objects are much smaller ( $\sim 1 - 5\%$  as compared to at least 10% for false positives) and 2) the fact that the detector missed them, indicate that they are very hard to distinguish from the background in the feature space. For

this reason, we use probabilistic methods such as pooled testing that allows us to quickly estimate the number of false negatives within the image population. Pooled testing is used often to estimate disease propagation within large populations of plants and animals. It is especially useful during the outbreak of an epidemic where there is a need to quickly estimate the extent of damage before taking action.

In the estimation problem, the primary goal is to estimate  $p$ , the individual probability of infection by testing pools of size  $s$ , of individual samples for the infection. There are different ways of testing pools, we use the inverse binomial sampling method [19] which provides the fastest estimates. In this method, suppose that  $T$  pools need to be tested before  $n$  positive samples have been identified, then the maximum likelihood estimate of  $p$  is given by

$$\hat{p} = 1 - \left(1 - \frac{n}{T}\right)^{\frac{1}{s}}. \quad (7)$$

Better estimates of  $p$  that take into account unequal pool size, inaccurate testing etc. can be found in [4, 14].

Before estimating the number of false negatives using (7), we need to accurately obtain  $\bar{D}$ , where the image patches are such that each one contains at most a single object, this is not easy since objects can be of different scales within the same image and there may be multiple objects with occlusions. Assuming we have an accurate  $\bar{D}$ , our problem is to estimate the number of patches with objects in them. This is a catch-22 situation since it is the detection problem all over again, therefore we resort to probabilistic methods like pooled testing. In order to obtain  $\bar{D}$ , we make the assumptions that most objects in an image are of approximately the same size and that the objects are not occluding each other. This allows us to use bounding boxes obtained from the detector as an indicator of scale for an object within the image. Next the bounding box is used to divide the image into smaller patches, excluding regions previously detected. This generates  $\bar{D}$  which contains a large set of image patches with a small number of objects (usually  $\sim 1 - 5\%$ ). In cases of multiple detections, we choose the box with the average height and width.

**Approaches for Image Pooling** The images need to be pooled together in a way that their individual characteristics (such as edges) are preserved. There are multiple ways to achieve this, and can be chosen specific to the application. In applying this to pedestrian detection we looked at three techniques namely, a) spatial averaging of two images b) Gradient domain averaging using the Frankot-Chellappa algorithm [13] and c) Poisson image blending [18]. Gradient domain techniques are appealing to use in this context since they can be tuned to preserve high frequency content (such as edges, corners) and leave out background. Concatenation

of multiple images can also be effective instead of pooling them, but one must consider the trade-off between the time it takes for humans to analyze a single pooled image versus multiple individual images, currently there is no research suggesting either way and a thorough evaluation of this is left for future work. A qualitative comparison of these methods is shown in fig:3. Based on our experience with multiple trials, we found that the simple spatial average pooling was sufficient for seeking interactive feedback from a user. Thus, we report results via the simple spatial average pooling in experiments.

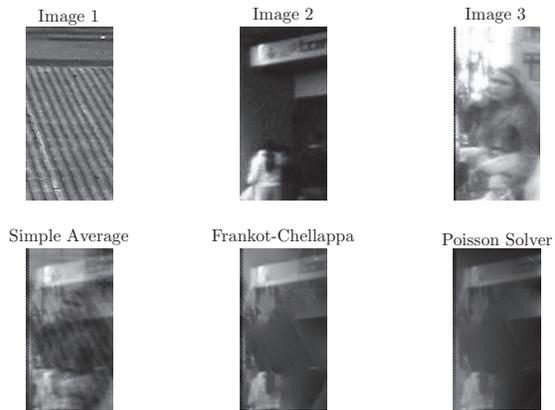


Figure 3: **Image Pooling** - This example shows a simple averaging scheme performing perceptually better than gradient domain methods, the original images are shown in the first row. Although gradient domain methods are attractive for image fusion, they may produce artifacts if they are not tuned correctly and work poorly on low-res images.

## 4. Experiments

In this section we demonstrate our system and its performance estimates compared to performance obtained using ground truth. We show results on the INRIA Pedestrian dataset [6] and the ETH Pedestrian Dataset [11], which are popularly used for pedestrian detection. The INRIA dataset is relatively simple with 2-4 pedestrians on average per image and of roughly the same scale. ETH dataset is more complex, with 10-12 pedestrians per frame and of varying sizes, it is acquired using a moving camera on a crowded street. We estimate the performance of three detectors namely Aggregate Channel Features (ACF)[8], Integral Channel Features (ChnFtrs)[9], MultiFtr+CSS [25] on these datasets.

### 4.1. Interactively test-driving pedestrian detectors

There has been great success in the field of pedestrian detection with a large number of excellent detectors proposed in the past few years. A good evaluation of the state of the art methods is available in [10], the authors have

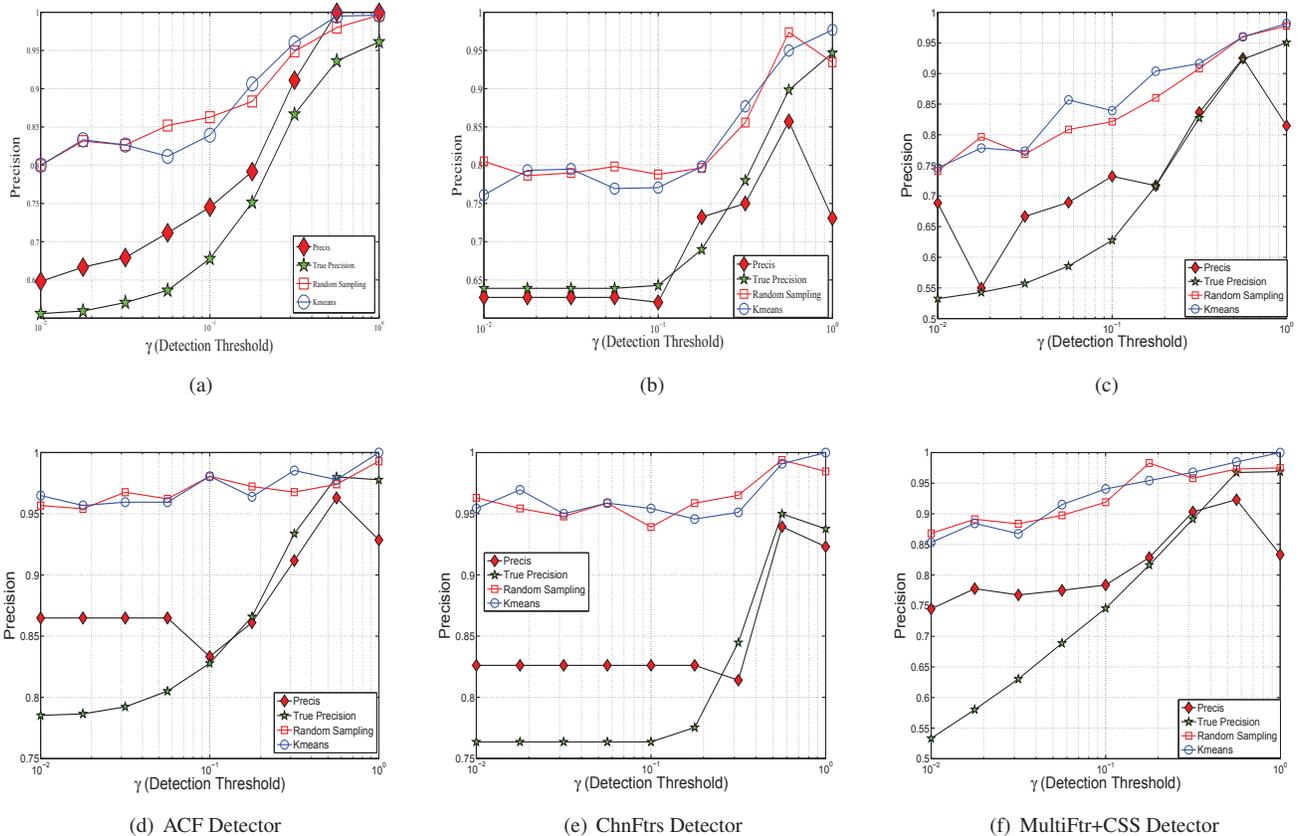


Figure 4: Comparing the precision of ACF Detector [8], the ChnFtrs Detector [9] and MultiFtr+CSS Detector [25] on the INRIA Pedestrian Dataset [6] (4(a),4(b),4(c)) and ETH Dataset [11] (4(d), 4(e), 4(f)) using just 10% of the total number of detections at each  $\gamma$ . The values for K-means and random sampling are averaged over 10 iterations as compared to a single iteration of Precis.

also made available the detection outputs for several detectors and datasets on the web, we choose 3 detectors and 2 datasets to demonstrate the proof of concept.

**Note on using labels to generate results** Samples obtained from a given sampling algorithm (random, K-medoids or homogenized Precis) are shown to a human to identify the number of false positives, using which the ‘perceived’ precision rate is calculated as explained earlier. Since we work with labeled datasets, we simulate the process with the given label which enabled much faster results for the sake of testing and proof of concept. However multiple duplicate detections around an object are considered false positives when evaluating with ground truth where only the detection with the highest overlap is chosen as the true detection and the rest are considered false. It is impossible to tell duplicates apart by looking at the detections from unlabeled data. To account for this during run time, we consider all the duplicate detections to be true detections even though they are actually false detections.

We used the detections for two datasets and estimated precision for different  $\gamma$ , the results for estimates based on

only 10% of the total number of detections are shown in fig 4. It is seen that the proposed homogenized Precis outperforms K-means and random sampling significantly. However, in most cases there is still a positive bias in the estimate which we attribute to the duplicate false positives that are impossible to distinguish without true labels.

The K-means and random sampling estimates are averages over 10 iterations, because they tend to be noisy, meaning this would require the human to look at 10 times the number of images used by Precis before getting a comparable estimate. In picking  $K_\gamma$ , smaller is lesser work for the human, and we also found that there is not a significant difference in performance in most cases.

## 4.2. Failure cases

Figure 5 shows the performance estimate using the Dalal and Triggs detector [6]. Precis has a poor estimate in this case compared to random sampling and K-means. This is because the general performance of the detector is far lower than state of the art detectors today, which results in many more false positives. Precis by construction is sen-

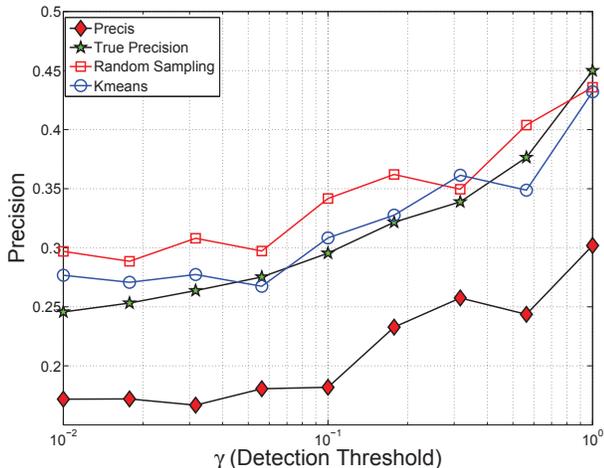


Figure 5: Precs fails to estimate precision accurately for detectors that perform poorly, but this is hardly an issue since poor detectors are never preferred in applications. The precision estimates are shown for the HOG+SVM detector [6] on the INRIA dataset.

sitive to diversity and often tends to pick many more false positives than pedestrians resulting in under estimating precision severely. Fortunately, this condition is not limiting since a detector is not used commercially until it is high performing.

### 4.3. Estimating the number of false negatives

Recall is the portion of objects that were correctly detected and can be completely described by  $\gamma, P_\gamma, N_\gamma, |D_{fn}|$  where  $\gamma$  is the detection threshold,  $P_\gamma$  is the precision,  $N_\gamma$  is the number of detections and  $|D_{fn}|$  is the number of false negatives or missed detections. All the quantities are known except  $|D_{fn}|$  so we will focus on estimating the number of false negatives, which gives us all the information to calculate Recall.

Estimating the number of false negatives is the second part of evaluating a detector. This is a much harder problem because a negative set is not clearly defined and could be extremely diverse. In obtaining  $\bar{D}_\gamma$  for different  $\gamma$ , the image patches are usually of low resolution (approx.  $50 \times 50$  pixels), and are pooled together before being presented to a human. We found that because of the poor resolution of the images, smaller pool size proved most effective. We used pool size  $s = 2$  in our experiments and randomly chose two images from the detection complement each time to pool them together. To use pooled testing we use the inverse binomial sampling approach where we fix  $n$ , the number of objects and make  $T$  a variable, from (7). So the user is asked to ground truth pooled samples until they find  $n$  objects. The choice of  $n$  affects the estimate in a way that smaller  $n$  gives a quicker estimate but is more noisy. We fixed  $n = 2$  and found that it was reliable and reasonably quick (in most

cases, one can find 2 images with pedestrians before about 2% of images have been seen). Using equation 7 we can get the the ratio  $\beta$  of objects within  $\bar{D}$ , which gives the number of false negatives as  $|D_{fn}| = \beta \times |\bar{D}|$ . Estimated results are shown in fig 6 for all the detectors and both the datasets. It is observed that the estimated results are poorer in the ETH dataset, this is because the dataset is complex with several pedestrians per frame and of constantly varying scale (the images are taken from a moving camera). Obtaining an accurate  $\bar{D}$  for constantly varying object scales is out of the scope of this work and we look towards addressing these issues in the future.

## 5. Discussion and Conclusion

In this paper, we introduce the problem of interactively test-driving object detectors in a user-centric manner, without the need for extensive ground-truthing. We presented the first system that is interactively able to provide estimates of the performance of an object detector without annotations. Although the problem is far from solved, the results shown in this work are promising. Using smart sampling techniques such as summarization and pooled testing, we show that the estimates obtained are close to those obtained using annotations and the effort required is significantly lesser than our baseline techniques. This points towards interactively evaluating classifiers, detectors etc. easily when purchasing them for tailored applications in the future.

## References

- [1] The PASCAL Visual Object Classes homepage. <http://pascallin.ecs.soton.ac.uk/challenges/VOC/>. Accessed: 16-11-2013. 2
- [2] The VMX project: Computer vision for everyone. <https://www.kickstarter.com/projects/visionai/vmx-project-computer-vision-for-everyone>. Accessed: 02-01-2014. 1
- [3] S. Branson, C. Wah, F. Schroff, B. Babenko, P. Welinder, P. Perona, and S. Belongie. Visual recognition with humans in the loop. In *ECCV*, pages 438–451, 2010. 3
- [4] C. L. Chen and W. H. Swallow. Using group testing to estimate a proportion, and to test the binomial model. *Biometrics*, 46(4):pp. 1035–1046, 1990. 3, 5
- [5] O. Chum, M. Perdoch, and J. Matas. Geometric minhashing: Finding a (thick) needle in a haystack. In *CVPR*, 2009. 3
- [6] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. *IEEE Conference on Computer Vision and Pattern Recognition*, pages 886–893, 2005. 2, 4, 5, 6, 7
- [7] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon. Information-theoretic metric learning. In *ICML*, 2007. 4
- [8] P. Dollár, S. Belongie, and P. Perona. The fastest pedestrian detector in the west. In *BMVC*, pages 1–11, 2010. 2, 5, 6
- [9] P. Dollár, Z. Tu, P. Perona, and S. Belongie. Integral channel features. In *BMVC*, pages 1–11, 2009. 2, 5, 6

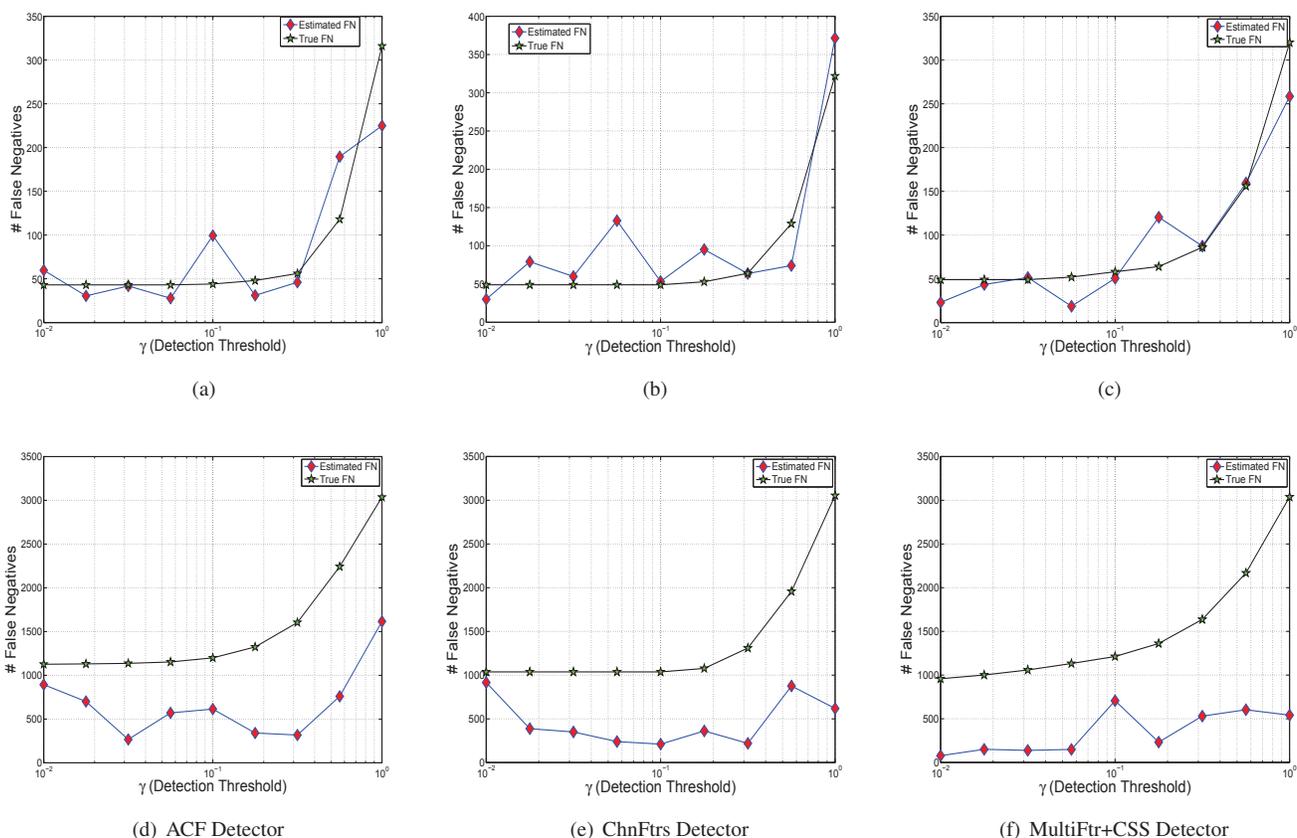


Figure 6: False negative estimates obtained via pooled testing on the INRIA 6(a), 6(b), 6(c) & ETH datasets 6(d), 6(e), 6(f). The negative bias in the estimates for the ETH dataset is attributed to its complexity as it has multiple pedestrians per frame and constantly varying scale.

- [10] P. Dollár, C. Wojek, B. Schiele, and P. Perona. Pedestrian detection: An evaluation of the state of the art. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2012. 2, 5
- [11] A. Ess, B. Leibe, K. Schindler, and L. van Gool. A mobile vision system for robust multi-person tracking. In *CVPR*, June 2008. 5, 6
- [12] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2):303–338, June 2010. 1
- [13] R. Frankot and R. Chellappa. A method for enforcing integrability in shape from shading algorithms. *IEEE Trans. Pattern Anal. Mach. Intell.*, 10(4):439–451, 1988. 5
- [14] G. Hepworth. Improved estimation of proportions using inverse binomial group testing. *Journal of Agricultural, Biological, and Environmental Statistics*, 2013. 3, 5
- [15] B. Kulis. Metric learning: A survey. *Foundations and Trends in Machine Learning*, 2013. 4
- [16] Y. J. Lee, J. Ghosh, and K. Grauman. Discovering important people and objects for egocentric video summarization. In *CVPR*, pages 1346–1353, 2012. 3
- [17] D. Nistér and H. Stewénus. Scalable recognition with a vocabulary tree. In *CVPR (2)*, 2006. 3
- [18] P. Pérez, M. Gangnet, and A. Blake. Poisson image editing. In *ACM SIGGRAPH 2003 Papers*, SIGGRAPH '03, New York, NY, USA, 2003. ACM. 5
- [19] N. Pritchard and J. Tebbs. Estimating disease prevalence using inverse binomial pooled testing. *Journal of Agricultural, Biological, and Environmental Statistics*, 2011. 5
- [20] N. Shroff, P. K. Turaga, and R. Chellappa. Video précis: Highlighting diverse aspects of videos. *IEEE Transactions on Multimedia*, 12(8):853–868, 2010. 3, 4
- [21] N. Shroff, P. K. Turaga, and R. Chellappa. Manifold précis: An annealing technique for diverse sampling of manifolds. In *NIPS*, 2011. 3, 4
- [22] E. S. Spelke. Principles of object perception. *Cognitive Science*, 1990. 3
- [23] V. Vapnik. *Statistical learning theory*. Wiley, 1998. 2
- [24] S. Vijayanarasimhan and K. Grauman. What’s it going to cost you?: Predicting effort vs. informativeness for multi-label image annotations. In *CVPR*, 2009. 3
- [25] S. Walk, N. Majer, K. Schindler, and B. Schiele. New features and insights for pedestrian detection. In *CVPR*, 2010. 2, 5, 6