

Discriminative Sparse Neighbor Approximation for Imbalanced Learning

Chen Huang, Chen Change Loy, *Member, IEEE*, and Xiaoou Tang, *Fellow, IEEE*

Abstract—Data imbalance is common in many vision tasks where one or more classes are rare. Without addressing this issue conventional methods tend to be biased toward the majority class with poor predictive accuracy for the minority class. These methods further deteriorate on small, imbalanced data that has a large degree of class overlap. In this study, we propose a novel discriminative sparse neighbor approximation (DSNA) method to ameliorate the effect of class-imbalance during prediction. Specifically, given a test sample, we first traverse it through a cost-sensitive decision forest to collect a good subset of training examples in its local neighborhood. Then we generate from this subset several class-discriminating but overlapping clusters and model each as an affine subspace. From these subspaces, the proposed DSNA iteratively seeks an optimal approximation of the test sample and outputs an unbiased prediction. We show that our method not only effectively mitigates the imbalance issue, but also allows the prediction to extrapolate to unseen data. The latter capability is crucial for achieving accurate prediction on small dataset with limited samples. The proposed imbalanced learning method can be applied to both classification and regression tasks at a wide range of imbalance levels. It significantly outperforms the state-of-the-art methods that do not possess an imbalance handling mechanism, and is found to perform comparably or even better than recent deep learning methods by using hand-crafted features only.

Index Terms—Imbalanced learning, decision forest, discriminative sparse neighbor approximation, data extrapolation.

I. INTRODUCTION

DATA imbalance exists in many vision tasks ranging from low-level edge detection [1] to high-level facial age estimation [2] and head pose estimation [3]. For instance, in age estimation, there are often many more images of the youth than the old on the widely used FG-NET [2] and MORPH [4] datasets. In edge detection, various image edge structures [5] obey a power-law distribution, as shown in Figure 1. Without handling this imbalance issue conventional vision algorithms have a strong learning bias towards the majority class with poor predictive accuracy for the minority class, usually of equal or more interest (e.g. rare edges may convey the most important semantic information about natural images).

The insufficient learning for the minority class is due to the complete lack of representation by a limited number of or even no examples, especially in the presence of small datasets. For instance, FG-NET age dataset has 1002 images in total with only 8 images over 60 years old. Certain age classes of 60+ ages have no images at all. This reveals a bigger challenge on unseen data extrapolation from the few

All the authors are with the Department of Information Engineering, The Chinese University of Hong Kong, Hong Kong. E-mail: {chuang, cclloy, xtang}@ie.cuhk.edu.hk.

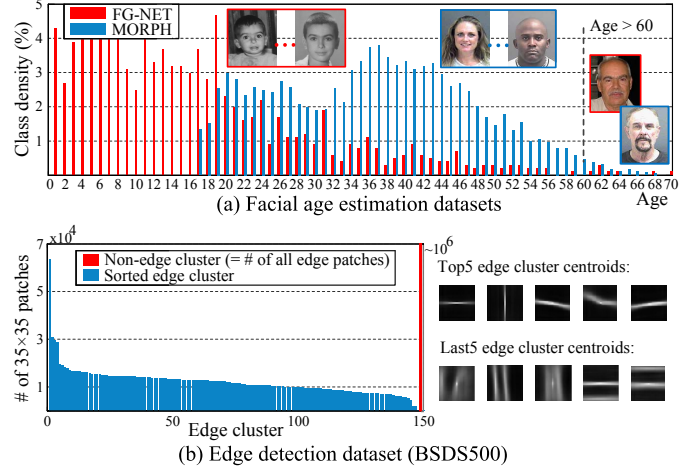


Fig. 1. Data imbalance in (a) age estimation and (b) edge detection. The given datasets characterize the underlying imbalanced distributions that can be seen as intrinsic in these problems.

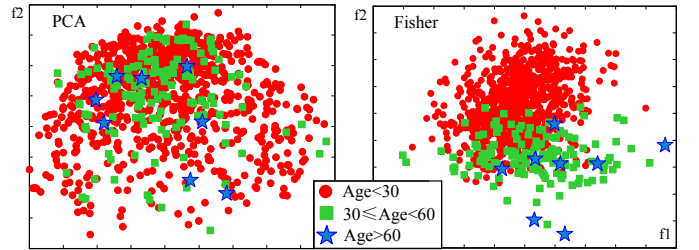


Fig. 2. 2D distributions of age data after PCA and LDA on FG-NET to highlight the class overlap issue.

minority class samples that usually have high variability. Even worse, the small imbalanced datasets can be accompanied by the class overlap problem. We plot the PCA and Fisher embeddings of FG-NET in Figure 2 to illustrate this problem. From the figure, it is evident that training a robust classifier or regressor capable of handling old ages is indeed a hard problem: (i) the corresponding minority class (blue star) contains insufficient samples for learning, (ii) these samples have high degree of variability which is hard to model, (iii) there is a severe class overlap between the rare samples and those from majority classes, further compounding the learning difficulty. Consequently, if we look into the local neighborhood of a minority class sample, it is very likely to be dominated by the majority class samples. Its weak local boundary would bias the prediction towards the majority class.

There are three common approaches to counter the neg-

ative impact of data imbalance: resampling [6], [7], cost-sensitive learning [8]–[10] and ensemble learning [11], [12]. Resampling approaches aim to make class priors equal by under-sampling the majority class or over-sampling the minority class (or both [6]). These methods can easily eliminate valuable information or introduce noise respectively. Cost-sensitive learning is often reported to outperform random resampling by adjusting misclassification costs, however the true costs are often unknown. An effective technique for further improvement is to resort to ensemble learning [13]. Chen *et al.* [11] combined bagging and weighted decision trees to generate a re-weighted version of random forest. We show in our experiments that the aforementioned strategies fall short in handling complex imbalanced data. Beyond empirical performance, the above approaches have two common drawbacks: 1) They are designed for either classification [6]–[8], [10]–[12] or regression [9] without a universal solution to both. 2) They have a limited ability to account for unseen appearances or extrapolate novel labels on the observed space. This is critical in the typical case of small imbalanced datasets where the minority class is under-represented by an excessively reduced number of or even no samples/labels.

In this paper we address the problems of data imbalance *and* unseen data extrapolation using a data-driven approach. The approach can be applied to *both* classification and regression scenarios. The key idea of our approach is intuitive – given a test sample, we first locate for it a ‘safe’ local neighborhood. This local neighborhood is formed by training samples, which are carefully mined so as to provide a relatively large coverage of minority class samples compared to the full training space. But overall, this space is tight and is less probable to be invaded by imposter samples¹. We show that this ‘safe’ local neighborhood can be constructed via a cost-sensitive decision forest. However, the local neighborhood may still be overwhelmed by majority classes especially when the minority ones are absolutely rare. Thus prediction by simple voting or averaging within it could easily smooth out the minority class samples. To this end, we further partition the local neighborhood into several discriminative but soft clusters with overlaps permitted. This process provides purer clusters eliminating the undesired class domination.

Subsequently, we propose a new Discriminative Sparse Neighbor Approximation (DSNA) method that allows robust prediction from our formed clusters. The clusters are all modelled as affine subspaces to account for unseen appearances in a similar spirit of [14]. The core of DSNA is a new cost function and a joint optimization approach to iteratively determine the best affine subspace that best approximates the test sample with the help of associated sparse neighbors. From the found neighbors and their approximating coefficients, we can transfer and combine their labels to achieve a robust prediction despite class-imbalanced issue. Figure 3 illustrates the effectiveness of DSNA in an age estimation example.

In summary, the main contributions of this paper are:

- A new discriminative sparse neighbor approximation

¹An imposter sample is defined as the one from a different class w.r.t. the test sample

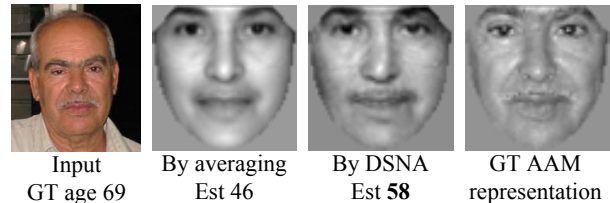


Fig. 3. A visualization of age estimation result when neither the testing appearance nor age label is observed during training. Averaging provides a crude way of estimating the face appearance (AAM, Active Appearance Model) from the nearest neighbors in the training set. The proposed discriminative sparse neighbor approximation (DSNA) provides a more robust estimation, thus an age value closer to the Ground Truth (GT).

(DSNA) method is proposed for unbiased predictions with preserved discriminative and extrapolative ability given class-imbalanced data.

- To facilitate robust predictions via DSNA, we formulate an effective way of constructing a safe local neighborhood through a cost-sensitive decision forest framework.
- The proposed method is applied to the vision tasks of age estimation (regression), head pose estimation (regression) and edge detection (classification) with varying degree of data imbalance and amount of data. It advances the state-of-the-art, sometimes considerably, across all tasks especially on highly imbalanced ones. It comes at only modest extra computational burden, showing its potential as a fast and general framework for imbalanced learning. Our results are particularly impressive when favorably compared to recent deep learning methods [15]–[22] as our method is built with no deeply learned features, but with the imbalance handling mechanisms absent in these deep models.

The rest of the paper is organized as follows. Section II briefly reviews related work on imbalanced learning and the considered vision tasks. Section III details the major components of the proposed framework of imbalanced learning. Section IV presents the results in ablation tests and three imbalanced vision tasks to highlight the benefit of each proposed component and their advantages over competing methods in different tasks. Section V concludes the paper.

II. RELATED WORK

Much effort for imbalanced learning in the machine learning community has been devoted to resampling approaches [7] that randomly under-sample the majority class or over-sample the minority. Other smart resampling techniques are also available (please refer to [7] for a comprehensive survey). Generally under-sampling may remove valuable information and over-sampling easily introduces noise with overfitting risks. Additionally, random over-sampling does not increase information by only replication, so it does not solve the fundamental ‘‘lack of data’’ issue. SMOTE [6], on the other hand, creates new examples by interpolating neighboring minority class instances. However, it is error-prone to interpolate noisy or borderline examples. Therefore under-sampling is often preferred to over-sampling [8], but is not suitable for small datasets (e.g. FG-NET) due to its caused information loss.

Cost-sensitive learning [8]–[10] as an alternative is closely related to resampling. Instead of manipulating samples at the data level, it adjusts misclassification costs at the algorithmic level and imposes heavier penalty on misclassifying the minority class. For example, Li and Lin [9] proposed RED-SVM to use the label-sensitive costs in the ordinal regression problem. In [23], a scaling kernel with the standard SVM is used to improve the classification on imbalanced datasets. Zadrozny *et al.* [10] combined cost sensitivity with ensemble approaches to further improve classification accuracy. Chen *et al.* [11] formed an ensemble of cost-sensitive decision trees by weighting the Gini criterion during the node splitting as well as final tree aggregation. We similarly grow cost-sensitive random trees but more generalized and principled ones, and propose a more discriminative and extrapolative “aggregation” scheme that proves necessary for complex imbalanced data.

The methods of [10], [11] already show the effectiveness of using classifier ensemble in the context of imbalanced data [11], [12]. Bagging and Boosting are the most popular ensemble strategies [13]. Boosting (e.g. [12], [24]) can easily embed the cost sensitivities in example weights according to the misclassification costs. Li *et al.* [24] further combined boosting with the training of an extreme learning machine. Generally boosting is vulnerable to noise and more prone to overfitting, which can be better addressed by Bagging [13]. Our method based on the improved random forest is essentially a Bagging-based method, thus shares this advantage.

Although much progress has been made on the vision tasks such as age estimation, head pose estimation and edge detection, relatively few works have studied the direct impacts of data imbalance on these tasks. Next we briefly discuss their representative works.

Age estimation: There are three main groups of age estimation methods: classification [2], [25], [26], regression [9], [27]–[30], and ranking [4], [31], [32] methods. OHRank [4], [32] surpasses previous classification- and regression-based methods by utilizing ordering information and cost sensitivities. However, the imbalance issue is neglected especially when designing the ordered binary classifiers at the youngest and oldest ages. Some recent state-of-the-arts focus on advanced feature extraction [25], [32], [33], including applying convolutional neural network (CNN) [15], [16] to automatically learn deep features instead of using hand-crafted ones. Unfortunately strong biases are still observed on imbalanced datasets, and we provide here an explicit solution to imbalanced learning with better results, using no deep features. Only three papers [30], [34], [35], as far as we know, consider data imbalance and sparseness when estimating ages. IsRCA [30] simply balances the number of used neighbors from each class to compute the similarity matrix for LPP (Locality Preserving Projection)-based dimensionality reduction. In [34], [35], the imbalance is only mitigated by leveraging adjacent labels in implicit ways, respectively via modeling cumulative attribute space and label distribution. We will show the advantage of our explicit solution and that of our built-in extrapolative mechanism for possible missing data/labels.

Head pose estimation: Methods for head pose estimation from 2D images can be categorized into two main groups: clas-

sification [36] and regression [37]–[41], with regression being more attractive for its continuous output. We refer readers to [42] for a comprehensive survey. Random forest is a popular choice for pose estimation in both classification [36] and regression [40] settings. It is also applied to depth images [43]. To our knowledge, the inherent imbalance in pose data [3] is seldom addressed again. Note on many pose datasets such as Pointing’04, the sparse data sampling (with typical pose intervals of $10^\circ+$) makes learning even more difficult.

Edge detection: State-of-the-art edge detection methods [1], [5], [44]–[49] mostly use engineered gradient features to classify edge pixels/patches. Recent CNN-based methods [17]–[22] achieve top results by learning deep features. Due to the large variety of edge structures, it is usually very hard to learn an ideal binary classifier to separate edges as one class from the non-edge class. Therefore some methods first cluster edge patches into compact subclasses (e.g. [5], [20]), and cast edge detection as a multi-way classification problem (i.e. to predict whether an input patch belongs to each edge subclass or the non-edge class) so as to implicitly solve the binary task. For the same reason, the numbers of “positive” and “negative” patches are commonly set to be equal to facilitate the binary goal. However, this results in a severe imbalance between each edge subclass and the dominant negative one (see Figure 1(b)), which is barely addressed properly by the above methods. Consequently, biased predictions tend to occur with low edge recall or damage of fine edge structures in those rare subclasses.

Another limitation of existing methods is that they cannot well predict the unseen edge structures from a novel class. For example, Sketch Tokens [5] only predict from a pre-defined set of edge classes based on random forest. Structured Edge (SE) detector [45] can model more subtle edge variations in a structured forest framework without the finite-class assumption, but still can only infer the edge structures observed during training. Although this problem is ameliorated by merging predicted structures while testing, it is in sharp contrast to our explicit DSNA method that empowers random forests to extrapolate.

III. METHODOLOGY

The proposed discriminative sparse neighbor approximation (DSNA) aims to provide unbiased predictions given a class-imbalanced dataset. More precisely, given a training set $\mathcal{D} = \{s_i = (\mathbf{x}_i, y_i)\}_{i=1}^N$, where $\mathbf{x}_i \in \mathbb{R}^D$ is the feature vector of sample s_i and y_i the label, our problem can be formulated as learning a function $F(\mathbf{x}) \rightarrow y$ to make unbiased predictions even in the presence of severely imbalanced and small datasets. The label $y \in \mathcal{C}$ refers to the class index (e.g. edge class) for classification or a numeric value (e.g. age and pose angle) for regression.

For a query \mathbf{q} , the key steps of DSNA are to draw a well localized neighborhood with a selective set of training data that is less probable to be invaded by imposter samples (i.e. from a different class w.r.t. the test sample), then follow the “divide and conquer” idea to perform a class-discriminative local clustering to obtain much purer clusters (modeled as affine subspaces) without undesired class domination, and finally choose the best cluster and use its member labels to predict.

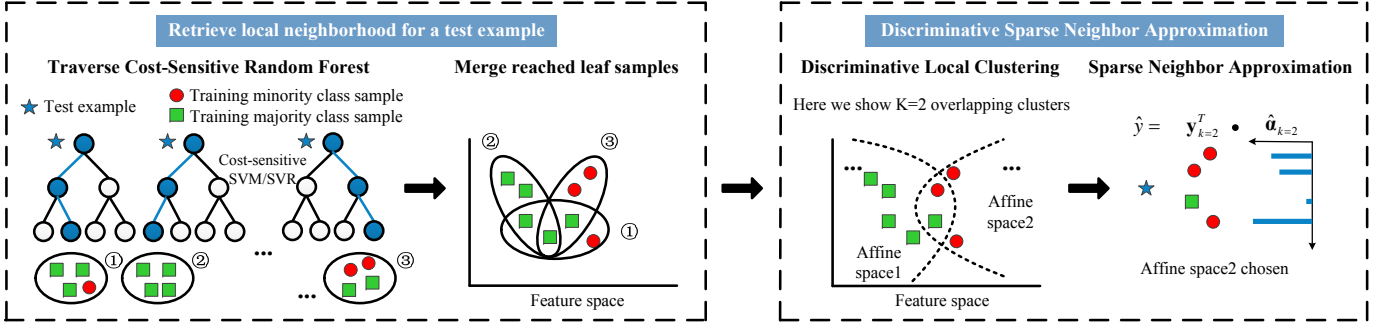


Fig. 4. The overall pipeline of our CS-RF-induced DSNA method.

The overall pipeline is shown in Figure 4. The pipeline begins with a cost-sensitive random decision Forest (CS-RF), which takes care of generating an initial good local neighborhood at leaf nodes, in order to reduce unnecessary distractions from majority classes. We retrieve all the leaf samples for a test instance, aiming to gain as more coverage of relevant minority samples as possible. The DSNA component starts with discriminative local clustering, and performs a sparse approximation to iteratively output unbiased predictions. To enable extrapolated prediction for unseen appearances, we model the found clusters as affine subspaces and extrapolate a prediction from them.

In the following, we first present the DSNA approach which is the key of this paper. We make the assumption that local neighborhood of training data is already available. We then describe the use of cost-sensitive random decision forest to obtain the local neighborhood.

A. Discriminative Sparse Neighbor Approximation

Discriminative local clustering - The first step of DSNA is to perform discriminative clustering within the local data neighborhood of a test sample. Suppose we have an initially retrieved local data neighborhood at hand, which can be noisy and class-imbalanced. This local neighborhood can be represented as

$$\mathcal{R} = \{s_i\}_{i=1}^M, \quad \mathcal{R} \subset \mathcal{D}, \quad M < N. \quad (1)$$

Intuitively, the samples in \mathcal{R} are close to the test sample based on some notions of metric or non-metric distance. Our objective is to separate the samples in \mathcal{R} based on their different class labels so as to pave the way for unbiased prediction of the test sample. We shall choose a clustering technique that possesses two desirable properties to achieve this goal: 1) It should generate discriminative clusters from one of which unbiased predictions can be made. 2) The found clusters should have adequate descriptiveness to account for unseen data patterns.

We achieve the aforementioned goal through a simple yet effective extension of K-means. It differs from the standard K-means in two respects. First, the inter-point distance $\tilde{d}(\mathbf{x}_i, \mathbf{x}_j)$ between \mathbf{x}_i and \mathbf{x}_j is label-aware:

$$\tilde{d}(\mathbf{x}_i, \mathbf{x}_j) = \begin{cases} d(\mathbf{x}_i, \mathbf{x}_j) * \mathbf{1}(y_i \neq y_j) & \text{for classification,} \\ d(\mathbf{x}_i, \mathbf{x}_j) * g(|y_i - y_j|) & \text{for regression,} \end{cases} \quad (2)$$

where $d(\cdot, \cdot)$ is the Euclidean distance, $\mathbf{1}(\cdot)$ is an indicator function, $g(y) = \tau y / (\max\{y\} - y + \text{eps})$ is a reciprocal increasing function with τ the trade-off parameter, and eps a small positive number to prevent overflow. The label-aware distance makes clustering discriminative by preferring the “same-class” data-pairs over those from different classes. In the extreme case, under classification scenarios for example, it forms clusters $\{\mathcal{L}_k\}_{k=1}^K$ each purely from one class even when the cluster members differ remarkably in appearances, which is suitable for classification.

Considering it is highly possible that the “pure” clusters in small imbalanced problems have limited samples, especially those mostly with the minority class samples, such clustering is actually not desirable for data/label extrapolation purposes. Hence, we allow cluster overlap by relaxing the cluster assignment of sample \mathbf{x}_i . Instead of assigning it solely to the nearest cluster centroid, we choose more than one centroids with distances slightly larger than the minimum distance in each K-means optimization iteration. This results in overlapping clusters each containing some “inter-class” samples. Such samples have complementary appearances to those “same-class” ones for enriching cluster representations.

Sparse neighbor approximation - The previous step generates K overlapping clusters $\{\mathcal{L}_k\}_{k=1}^K$ with their feature matrices $\{\mathbf{L}_k\}_{k=1}^K$ and labels $\{\mathbf{y}_k\}_{k=1}^K$. Our problem becomes how to discriminatively predict the label of a query \mathbf{q} and extrapolate to its possibly unseen appearance simultaneously.

To address this problem, we model each cluster by an affine hull model \mathcal{AH}_k [14] that is able to account for unseen data of different modes, and then choose the best prediction returned by them. Every single \mathcal{AH}_k covers all possible affine combinations of its belonging samples and can be parameterized as:

$$\mathcal{AH}_k = \{\mathbf{x} = \boldsymbol{\mu}_k + \mathbf{U}_k \mathbf{v}_k, k = 1, \dots, K\}, \quad (3)$$

where $\boldsymbol{\mu}_k = \sum_{\mathbf{x}_i \in \mathcal{L}_k} \mathbf{x}_i / |\mathcal{L}_k|$ is the centroid, \mathbf{U}_k is the orthonormal basis obtained from the SVD of centered \mathbf{L}_k , and \mathbf{v}_k is the coefficient vector.

Note that to predict the class label of query \mathbf{q} , we need to know which cluster the query should be assigned to. Thus Eq. 3 only provides a partial answer to our problem, since the cluster index k remains unknown. To this end, we formulate a joint optimization problem for simultaneously

finding the belonging cluster of the query and its affine hull approximation:

$$\begin{aligned} \min_{k, \mathbf{v}_k, \boldsymbol{\alpha}_k} & \|\boldsymbol{\mu}_k + \mathbf{U}_k \mathbf{v}_k - \mathbf{L}_k \boldsymbol{\alpha}_k\|_2 + \lambda \|\boldsymbol{\alpha}_k\|_1 + \gamma \|\boldsymbol{\alpha}_k - \bar{\boldsymbol{\alpha}}_k\|_1, \\ \text{s.t.} & \|\mathbf{q} - (\boldsymbol{\mu}_k + \mathbf{U}_k \mathbf{v}_k)\|_2 \leq \varepsilon, \end{aligned} \quad (4)$$

where $\varepsilon \geq 0$, and λ and γ are regularization parameters.

We explain the objective function as follows:

First term - This term approximates \mathbf{q} over the k^{th} cluster using the cluster's affine subspace as well as the feature matrix of associated member samples \mathbf{L}_k . This term is motivated by affine hull models [14] but differs significantly in the following aspects:

i) the affine space is class-aware. In particular, the affine space is learned from our class-discriminating cluster, and we solve for the best approximation among the clusters. A class-aware sparsity constraint is further imposed to promote discrimination (*Third term*).

ii) the affine space approximation benefits from the enriched descriptiveness of overlapping clusters.

Second term - This term constrains the loose affine approximation by imposing sparsity among the cluster samples. Thus a large drift is avoided when extrapolating \mathbf{q} on the affine subspace, because we constrain the affine subspace to be near to the observed samples using feature matrix \mathbf{L}_k .

Third term - This term regularizes the coefficient vector $\boldsymbol{\alpha}_k$ so it that focuses more on the ‘‘same-class’’ nearest neighbors, $\mathcal{N}_k = \{\mathbf{x}_i \in \mathcal{L}_k : \tilde{d}(\mathbf{x}_i, \mathbf{q}) \leq \varepsilon_k\}$, which are found by using the class-aware distances in Eq. 2. From our experiments, we empirically found that this term is useful to provide stable predictions. Formally, the $\bar{\boldsymbol{\alpha}}_k$ is estimated as:

$$\bar{\boldsymbol{\alpha}}_k = \sum_{\mathbf{x}_i \in \mathcal{N}_k} w_i \boldsymbol{\alpha}_i, \quad w_i \propto \exp(-\tilde{d}(\mathbf{x}_i, \mathbf{q})/h), \quad (5)$$

where h is the decay parameter, and $\boldsymbol{\alpha}_i$ is the representation coefficient of the i^{th} neighbor with the i^{th} element equal to one and the rest zero.

Eq. 4 can be solved by alternatively seeking the best affine approximation $\min_{k, \mathbf{v}_k} \|\mathbf{q} - (\boldsymbol{\mu}_k + \mathbf{U}_k \mathbf{v}_k)\|_2$ and the sparse neighbor approximation with two l_1 -norms:

$$\min_{\boldsymbol{\alpha}_k} \|\mathbf{q} - \mathbf{L}_k \boldsymbol{\alpha}_k\|_2 + \lambda \|\boldsymbol{\alpha}_k\|_1 + \gamma \|\boldsymbol{\alpha}_k - \bar{\boldsymbol{\alpha}}_k\|_1, \quad (6)$$

which can be efficiently solved by using the Augmented Lagrange Multiplier (ALM) method [50].

With the converged $\hat{\boldsymbol{\alpha}}_k$, the label for \mathbf{q} is finally predicted as $\hat{y} = \mathbf{y}_k^T \hat{\boldsymbol{\alpha}}_k$ for regression or by majority voting for classification (in this case we determine the nonzero entries of thresholded $\hat{\boldsymbol{\alpha}}_k$, and vote among the corresponding \mathbf{y}_k). The initial label to start the iterative process is set as the mean or majority vote of \mathbf{y}_k in the best-fit cluster.

B. Cost-Sensitive Random Decision Forest

Returning to the initial step of finding a ‘safe’ local neighborhood, we choose random decision forest for its efficiency and robustness. We first traverse a test example through every trained decision tree and retrieve the respective training samples \mathcal{R}_t stored at the leaf node. Traditional random forest

calculates either a class distribution for classification or a local mean for regression from each \mathcal{R}_t , and aggregates them as the final prediction. We face two fundamental problems by doing so: in the case of absolute rarity, each \mathcal{R}_t will still predominantly consist of majority classes that make simple aggregation biased to them; or \mathcal{R}_t may form pure but small disjuncts [7] of minority class samples leading to overfit.

Therefore, we instead merge all the retrieved leaf sample sets $\{\mathcal{R}_t\}$ into a single one $\mathcal{R} = \cup_t \mathcal{R}_t$, and treat \mathcal{R} as our initial local neighborhood in Eq. 1. Then DSNA is applied to it for prediction as described in Section III-A. Such a simple merging has the benefit of facilitating our data-driven prediction with as more coverage of relevant minority samples as possible. This can be easily realized thanks to the diversities between merged trees. In fact, random forest has the proved upper bound of generalization error given by [13]:

$$\epsilon \leq \rho(1 - m^2)/m^2, \quad (7)$$

where m is the strength of individual trees and ρ is the correlation between decision trees. Hence in order to maintain the low correlation and diversity among trees, we just keep the Bagging nature and feature randomness at internal nodes in the standard random forest.

To make the merged neighborhood less distracted by impostor samples, we focus on improving the strength m of each tree in the context of data imbalance by making the tree cost-sensitive. We have explored different cost-sensitive schemes, such as the re-weighting of nodes as in [11] and boosting of trees with class costs, but seen marginal effects. We finally came to a modified node splitting rule that can not only take into account the imbalanced distribution, but also can work seamlessly for both classification and regression.

Specifically, we first follow the standard Bagging procedure to grow an ensemble of random trees. Each tree recursively divides the input space into disjoint partitions in a coarse-to-fine manner. The key is to design good splitting functions. For a node j with local samples \mathcal{S}_j , a binary function $\phi_j : \mathbb{R}^{D'} \rightarrow \{0, 1\}$ is trained on some randomly sampled features ($D' = \sqrt{D}$) and splits into \mathcal{S}_j^l and \mathcal{S}_j^r to maximize the information gain:

$$\mathcal{I}(\mathcal{S}_j, \phi_j) = H(\mathcal{S}_j) - \left(\frac{|\mathcal{S}_j^l|}{|\mathcal{S}_j|} H(\mathcal{S}_j^l) + \frac{|\mathcal{S}_j^r|}{|\mathcal{S}_j|} H(\mathcal{S}_j^r) \right), \quad (8)$$

where $H(\cdot)$ denotes the class entropy. For regression, information gain can be replaced by the label variance as $H(\mathcal{S}) = \sum_y (y - \mu)^2 / |\mathcal{S}|$ where $\mu = \sum_y y / |\mathcal{S}|$. Training stops when a maximum depth is reached or if information gain or local sample size $|\mathcal{S}_j|$ falls below a fixed threshold.

The standard node splitting function ϕ_j is not necessarily optimal with respect to imbalanced data. To alleviate this problem, in both classification and regression scenarios, we incorporate a cost function $f(\cdot) \geq 0$ into ϕ_j that penalizes more heavily on the minority class. We describe in the following the cost function for classification trees and regression trees, respectively.

In classification trees, we first apply the widely used K-means technique [40], [45] to cluster \mathcal{S}_j into $\{\mathcal{S}_j^k\}_{k=1}^2$, and

then the splitting function ϕ_j that best preserves the two clusters is determined by a cost-sensitive version of linear SVM:

$$\min_{\mathbf{w}} \|\mathbf{w}\|_2 + C \sum_{k=1}^2 f(p_k) \sum_{\mathbf{x}_i \in \{\mathcal{S}_j^k\}} (\max(0, 1 - z_i \mathbf{w}^T \mathbf{x}_i))^2, \quad (9)$$

where $p_k = |\mathcal{S}_j^k|/|\mathcal{S}_j|$ denotes the cluster proportion, \mathbf{w} is the weight vector, C is a regularization parameter, and $z_i = 1$ if $\mathbf{x}_i \in \mathcal{S}_j^1$ and -1 otherwise. Each sample is finally sent to either \mathcal{S}_j^1 or \mathcal{S}_j^2 by $\text{sgn}(\mathbf{w}^T \mathbf{x}_i)$. The resulting splitting function is thus *learned* in a cost-sensitive manner instead of being chosen from some predefined splitting rules. Note the cost here is defined as a function of the cluster distribution rather than the targeted class distribution, but they will correlate well at the deeper tree depth with much purer nodes where Eq. 9 can better play its role.

In regression trees, we perform a cost-sensitive regression at each node \mathcal{S}_j using a weighted linear SVR:

$$\min_{\mathbf{w}} \|\mathbf{w}\|_2 + C \sum_{y \in \mathcal{C}} f(p_y) \sum_{\substack{y_i=y \\ \mathbf{x}_i \in \mathcal{S}_j}} (\max(0, |y_i - \mathbf{w}^T \mathbf{x}_i| - \varepsilon))^2, \quad (10)$$

where $\varepsilon \geq 0$, and we directly penalize the true label distribution $\{p_y = |\{y_i = y, \mathbf{x}_i \in \mathcal{S}_j\}|/|\mathcal{S}_j|\}$ as costs. The node then branches left if the numeric prediction $\{\mathbf{w}^T \mathbf{x}_i\}$ is smaller than the local mean of labels $\sum_{\mathbf{x}_i \in \mathcal{S}_j} y_i/|\mathcal{S}_j|$, otherwise branches right.

In practice, we use the cost transformation technique in [4] to solve the above weighted SVM/SVR. The cost function $f(\cdot)$ is defined by a reciprocal decreasing function as $f(p) = (1 - p)/p$. Obviously, $f(p)$ gives larger weights to the minority classes which proves effective to improve their prediction accuracies without losing the overall performance in our experiments. In addition, we use the inverse class frequencies to reweight the information gain (Eq. 8, as in [11]) to select the best D' random features in both classification and regression trees. The result is a CS-RF framework able to carve reasonably good local neighborhoods for both the majority and minority classes.

C. Convergence and Complexity

Our full algorithm is detailed in Algorithm 1. Similar to the affine hull (AH) method [14], Algorithm 1 can converge to a global solution. Compared with AH's non-asymptotic convergence rate of $O(1/t^2)$, our DSNA method converges even faster as shown in Figure 5. Typically DSNA converges within 10 iterations with lower objective values thanks to the introduced class discrimination as a guidance. In our experiments, we will visualize some converged examples with accurate predictions in different vision tasks. As there are usually dozens or hundreds of retrieved samples in the initial data neighborhood, DSNA typically takes only 0.5s to run on an Intel Core i7 4.0GHz CPU.

IV. EXPERIMENTS

We validate the effectiveness of our CS-RF-induced DSNA method in three vision tasks at various imbalance levels: the

Algorithm 1 : CS-RF-Induced DSNA

Input: Training set $\{(\mathbf{x}_i, y_i)\}_{i=1}^N$, trained CS-RF, query \mathbf{q} .

Initialization: to predict y of \mathbf{q}

- Merge for \mathbf{q} all its reached leaf samples to \mathcal{R} .
- Via **Discriminative Local Clustering**, obtain clusters $\{\mathcal{L}_k\}_{k=1}^K$, features $\{L_k\}_{k=1}^K$ and labels $\{y_k\}_{k=1}^K$.
- Set $y^{(0)}$ as the mean of y_k for regression or its majority vote for classification in $\mathcal{A}\mathcal{H}_k$ that best approximates \mathbf{q} .

Outer Loop: Iterate on $t = 1, \dots, T$ until convergence

- Update $\{k^{(t-1)}, \mathbf{v}_k^{(t-1)}\}$ by Eq. 3 as the ones that best approximate \mathbf{q} .
- Update the sparse coefficient estimate $\bar{\alpha}_k^{(t-1)}$ by Eq. 5.
- Update $\alpha_k^{(t-1)}$ via **Sparse Neighbor Approximation** by minimizing Eq. 6.
- Predict label $y^{(t)} = \mathbf{y}_k^T \alpha_k^{(t-1)}$ or by majority voting among y_k with nonzero coefficients.

Output: Converged label \hat{y} .

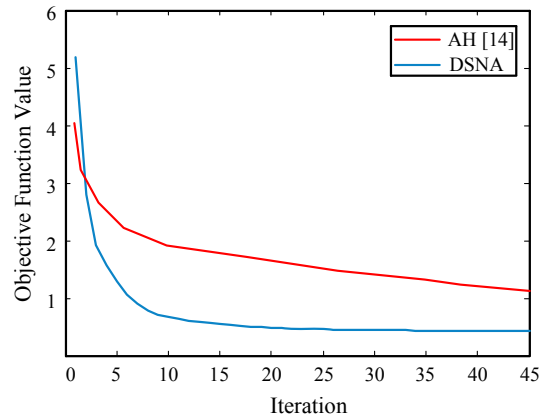


Fig. 5. Comparison of the convergence of unsupervised AH [14] and our DSNA given a query in the age estimation task.

high-level tasks of age estimation and head pose estimation (by regression) and the low-level task of edge detection (by classification).

A. Experimental Settings

Dataset settings: For age estimation, the FG-NET [2] and MORPH [4] datasets are used. FG-NET contains 1002 facial images of 82 subjects with ages in a range from 0 to 69. Algorithms are evaluated by the leave-one-person-out protocol. MORPH contains about 55000 images of more than 13000 subjects with ages between 16 and 77. We randomly split it into three disjoint subsets S1, S2 and S3 as in [16]. Algorithms repeat 1) training on S1, testing on S2+S3 and 2) training on S2, testing on S1+S3 with the average result reported. Both datasets are highly imbalanced (see Figure 1(a)) and class-overlapped. FG-NET further suffers from the issue of small data. For both, we use AAM [51] as the feature extractor, and Mean Absolute Error (MAE) as the evaluation metric.

For head pose estimation, poses should intrinsically admit an imbalanced distribution with much more near-frontal instances than the profile ones. Unfortunately, we are unable to obtain such datasets with ground truth labels (e.g. ‘‘Face

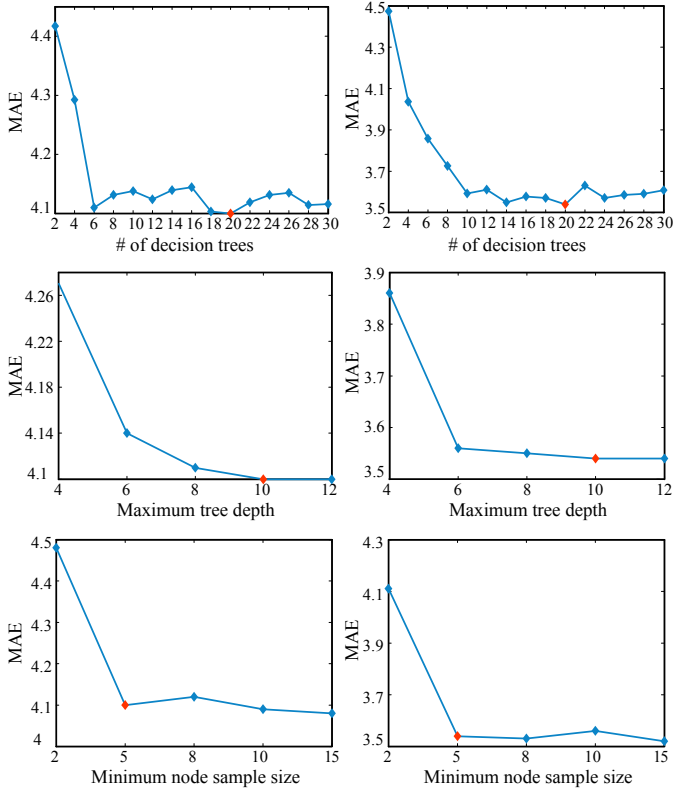


Fig. 6. Parameter sweeps for age estimation (left column) and head pose estimation (right column). Each row (from top to bottom) considers the parameter of tree number, maximum tree depth and minimum node sample size, respectively. The chosen parameter value is marked in red.

Pose” dataset [3]) for experiments. We adopt the popular Pointing’04 dataset instead that exhibits some imbalance in pitch angles. The dataset contains images from 15 subjects each with two series of 93 pose images. The pose is discretized into 9 pitch angles $\{\pm 90^\circ, \pm 60^\circ, \pm 30^\circ, \pm 15^\circ, 0^\circ\}$ and 13 yaw angles $\{\pm 90^\circ, \pm 75^\circ, \pm 60^\circ, \pm 45^\circ, \pm 30^\circ, \pm 15^\circ, 0^\circ\}$. However, when the pitch angles are $\{\pm 90^\circ\}$, the yaw is always $\{0^\circ\}$ (so $7 \times 13 + 2 = 93$ poses in total), leading to an imbalance ratio of 1:13 between $\{\pm 90^\circ\}$ pitch angles and others. We further test in the case where pitch angles are randomly removed to form a Gaussian-like distribution to mimic the real-world imbalanced distribution. As in [39], [40], evaluation of MAE is performed with 5-fold cross-validation using HOG features.

For edge detection, we use the BSDS500 [1] and NYUD (v2) [52] datasets, the latter for testing cross-dataset generalization. BSDS500 contains 200 training, 100 validation and 200 testing images. NYUD contains 1449 pairs of RGB and depth images. We follow [47] to use 60%/40% training/testing split (1/3 training data for validation) with the images reduced to 320×240 pixels. For cross-dataset testing, we only use RGB images on both datasets. We combine our method in classification mode with the structured edge detector [45] since it induces classification forest like us but operates on edge patches instead of pixels, which proves efficient in practice. We use the same multiple low-level features extracted from 32×32 image patches and apply non-maximal suppression prior to evaluation as in [45]. Edge detection accuracy is evaluated by: fixed contour threshold (ODS), per-image best threshold

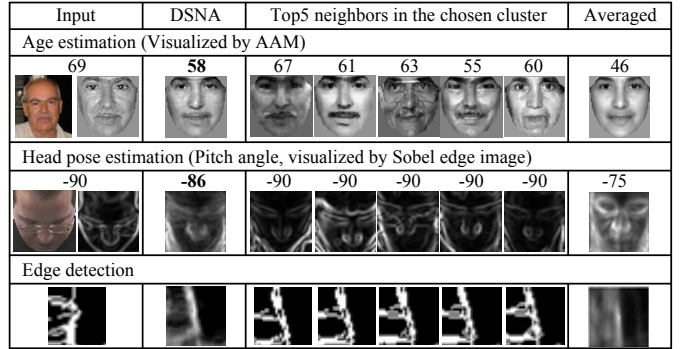


Fig. 7. Visualizations of both the DSNA converged result and simple averaged result on the retrieved samples by CS-RF. Results are shown for the minority class testing samples in all the three tasks.

TABLE I
ABLATION TEST AND COMPARISONS OF CS-RF AND DSNA IN AGE ESTIMATION (MAE ON FG-NET), HEAD POSE ESTIMATION (AVG. MAE ON POINTING’04) AND EDGE DETECTION (ODS ON BSDS500, HIGHER IS BETTER).

| Methods | RF | RF+ SMOTE [6] | RED- SVM [9] | WRF [11] | CS-RF | CS-RF+ AH [14] | CS-RF+ DSNA |
|---------|------|------------------|-----------------|-------------|-------------|-------------------|----------------|
| Age | 5.28 | 5.39 | 5.24 | – | 4.81 | 4.89 | 4.10 |
| Pose | 6.41 | 6.65 | 6.53 | – | 4.02 | 4.28 | 3.54 |
| Edge | 0.75 | 0.75 | – | 0.75 | 0.76 | 0.76 | 0.78 |

(OIS), and average precision (AP) [1].

Parameters: For age and head pose estimation, we empirically combine 20 cost-sensitive trees in our regression forest, and terminate splitting when the maximum depth 10 is reached or if the node sample size is fewer than 5. Figure 6 shows the robustness of these parameters across tasks. Evaluations are done by varying one parameter at a time, with others fixed. The chosen parameter value is marked in red. For edge detection, our method is combined with [45] and uses the same parameter settings.

Cross-validation is used to determine the trade-off parameter C for cost-sensitive SVM/SVR (Eq. 9 and 10), τ for biased distance (Eq. 2), λ and γ in Eq. 4. We select K for discriminative local clustering from 2 to 4.

B. Evaluation of the CS-RF and DSNA

We start with evaluating our key components of CS-RF and DSNA. CS-RF concerns about generating good local neighborhood, while DSNA makes unbiased and extrapolative prediction and is the major contribution of this paper.

Figure 7 visualizes the benefit of DSNA over simple averaged prediction in the three considered tasks. Clearly, given an appropriate local neighborhood, e.g. by CS-RF, DSNA can localize the correct mode (cluster) in it for the difficult minority class samples, thus making much more unbiased predictions than by simply averaging. More significantly, for age estimation on the small FG-NET dataset, there are very few elderly samples with many missing classes, but our DSNA extrapolates well from the limited data.

Table I quantifies the benefits of both CS-RF and DSNA against other competitive schemes in the three tasks. Note all

TABLE II
COMPARISONS OF AGE ESTIMATION RESULTS (MAE) ON FG-NET AND MORPH DATASETS.

| FG-NET | | | | | | | MORPH | |
|-------------------------|----------------|------------------------|-------------------|------------|------------|-------------|------------|-------------|
| RUN [31] | RED-SVM [9] | MTWGP [29] | BIF [33] | CPNN [35] | CSOHR [32] | CA-SVR [34] | KPLS [27] | KCCA [28] |
| 5.33 | 5.24 | 4.83 | 4.77 | 4.76 | 4.70 | 4.67 | 4.04 | 3.98 |
| Choi <i>et al.</i> [26] | MidFea-NS [15] | Han <i>et al.</i> [53] | RealAdaBoost [25] | OHRank [4] | lsRCA [30] | CS-RF+DSNA | MSCNN [16] | CS-RF+DSNA |
| 4.66 | 4.62 | 4.60 | 4.49 | 4.48 | 4.38 | 4.10 | 3.63 | 3.54 |

the RF variants in the left and middle columns—RF+SMOTE, WRF (Weighted RF) and CS-RF simply average tree predictions as in standard RF. They do not consider data extrapolation as AH (Affine Hull) [14] and DSNA do. We make the following observations: 1) The over-sampling method SMOTE shows no benefits over Bagging in standard RF since it can introduce undesirable noise (e.g. in age and pose cases). 2) Cost-sensitive learning, in the middle column, helps for these imbalanced tasks, and our CS-RF consistently outperforms RED-SVM and WRF. This suggests that simple weighting schemes in RED-SVM and WRF are not adequate in complex imbalanced tasks. In contrast, our CS-RF can be seen as an ensemble of cost-sensitive experts organized in hierarchical trees, with higher capability and robustness. Another advantage is that CS-RF provides a unified cost-embedded solution to both regression and classification. 3) The supervised DSNA combined with CS-RF leads to large improvements, whereas the unsupervised AH shows no improvements or even worse results. This emphasizes the importance of using supervisory information. DSNA uses such information intelligently by extrapolating from several discriminatively trained AH models with a class-aware constraint (Eq. 4).

C. Comparison with State-of-the-Arts

Age estimation: We compare with the state-of-the-arts on FG-NET and MORPH datasets in Table II. Our CS-RF+DSNA outperforms most methods by a large margin, and reduces the MAEs of the runner-up lsRCA and MSCNN on the two datasets by 6.4% and 2.5% respectively. The larger improvement on FG-NET is impressive because the dataset is very small and has missing class labels (old ages). This validates our competence in synthesizing novel labels on small imbalanced datasets. Note the mere cost-sensitive methods RED-SVM, CSOHR and OHRank all show their inferiority on this imbalanced dataset, necessitating the ability of extrapolation. The advanced features—Bio-Inspired Features (BIF), generalized BIF with scattering transform in [32] and feature selection by RealAdaBoost [25] also do not reach top in this task. In contrast, our method, using the AAM features only, even outperforms the deep feature-based MidFea-NS and MSCNN due to the handling of data imbalance. Compared with the indirect imbalance-handling methods, CPNN, CA-SVR and lsRCA, ours performs much better by introducing explicit mechanisms that are discriminative and extrapolative.

Figure 8 shows the MAEs per decade of CS-RF+DSNA and the competitive OHRank with public implementation. From the comparison, the benefits of our method become prominent at old ages with extremely limited samples. We attribute the

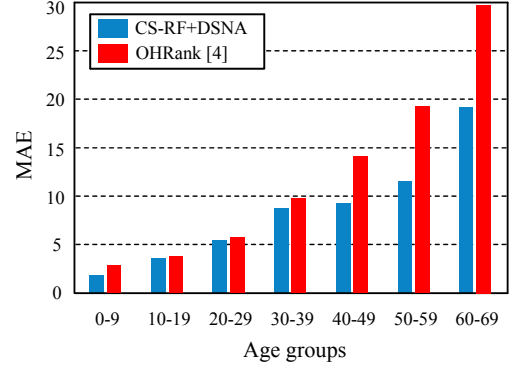


Fig. 8. MAEs at different age groups on the FG-NET dataset.

TABLE III
COMPARISON OF POSE ESTIMATION MAEs[°] ON POINTING'04.

| Method | Yaw | Pitch | Avg. |
|--------------------------|-------------|-------------|-------------|
| KPLS [39] | 6.56 | 6.61 | 6.59 |
| SLDML [41] | 6.31 | 6.71 | 6.51 |
| Fenzi <i>et al.</i> [38] | 5.94 | 6.73 | 6.34 |
| GLLiM [37] | 5.62 | 6.68 | 6.15 |
| KRF [40] | 5.29 | 2.51 | 3.90 |
| CS-RF+DSNA | 5.04 | 2.03 | 3.54 |

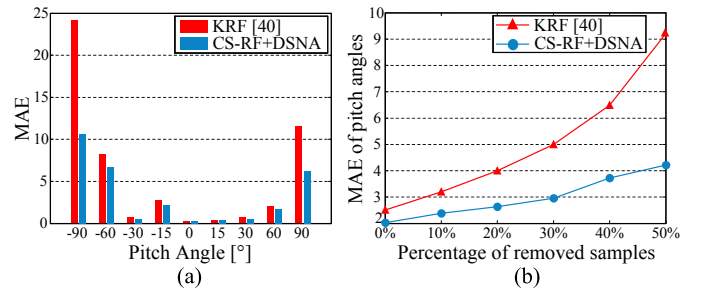


Fig. 9. Comparisons of imbalanced pitch angle estimation on Pointing'04. (a) MAEs at different pitch angles. (b) Average pitch MAEs with different percentages of samples removed.

performance gain to the extrapolative ability of the proposed DSNA.

Head pose estimation: Table III compares our method with the regression-based prior arts KPLS, SLDML, Fenzi *et al.*, GLLiM and KRF on Pointing'04 dataset. As mentioned in Section II, the sparse sampling of pose angle compounds the learning difficulty, especially for the imbalanced pitch angles. Our method performs best again for both pose angles, with a large margin for pitch. Figure 9(a) further compares our results

TABLE IV
COMPARISON OF EDGE DETECTION RESULTS ON THE BSDS500 DATASET.

| Method | ODS | OIS | AP |
|-----------------------------|-------------|-------------|-------------|
| ISCRA [48] | 0.72 | 0.75 | 0.46 |
| gPb-owt-ucm [1] | 0.73 | 0.76 | 0.73 |
| Sketch Tokens [5] | 0.73 | 0.75 | 0.78 |
| SCG [47] | 0.74 | 0.76 | 0.77 |
| PMI+sPb [46] | 0.74 | 0.77 | 0.78 |
| SE [45] | 0.75 | 0.77 | 0.80 |
| OEF [49] | 0.75 | 0.77 | 0.82 |
| SE+multi-ucm [44] | 0.75 | 0.78 | 0.76 |
| DeepNet [17] | 0.74 | 0.76 | 0.76 |
| N ⁴ -Fields [18] | 0.75 | 0.77 | 0.78 |
| DeepEdge [19] | 0.75 | 0.77 | 0.81 |
| DeepContour [20] | 0.76 | 0.78 | 0.80 |
| HFL [21] | 0.77 | 0.79 | 0.80 |
| HED [22] | 0.78 | 0.80 | 0.83 |
| CS-SE+DSNA | 0.77 | 0.79 | 0.81 |

at different pitch angles with those of KRF, the state-of-the-art regression forest-based method. Our method benefits from the proposed cost-sensitive (CS) RF and extrapolative DSNA, thus better handling the data imbalance at the rare $\pm 90^\circ$ poses. Specifically, we obtain a MAE of 8 degrees for those $\pm 90^\circ$ poses with only 24 training samples (hundreds of samples for other poses), which is 55.6% lower than that of KRF. The larger MAE at -90° is due to the higher variations of this angle compared to that of 90° . We finally compare with KRF in Figure 9(b) where pitch samples are randomly removed to form a Gaussian-like distribution aiming to mimic the real-world one. It can be observed that the performance of our method degrades more gracefully with the increase of removed data, showing a strong ability to handle small imbalanced data. **Edge detection:** In this task, severe imbalance exists between the positive edge classes and negative non-edge class. We refer our combined method with structured edge (SE) detector [45] as CS-SE+DSNA.

Table IV summarizes extensive comparisons with state-of-the-art methods on BSDS500. It is observed that CS-SE+DSNA outperforms all “shallow” methods (top cell) across all evaluation metrics, and also performs better than most deep models (bottom cell) and is comparable to the top HED. The results are impressive since our method only uses hand-designed features. It is worth mentioning that HED utilizes a deep CNN with as many as 16 layers for good performance. However our “shallow” method can still compete with such deep models due to the proposed mechanism for handling imbalance.

This advantage of our method also holds with respect to the compared non-deep methods. Among them, Sketch Tokens, SE and OEF are based on random forest similar to ours. The performance gain compared to their results can thus be attributed to the capability of correctly classifying imbalanced edge patches and generalizing to novel edge structures. Figure 10(a) confirms this standpoint by comparing CS-SE+DSNA with three related random forest-based methods, including DeepContour that applies random forest on top of deeply learned features. Clearly, CS-SE+DSNA is able to

TABLE V
COMPARISON OF EDGE DETECTION (TOP) AND CROSS-DATASET GENERALIZATION (BOTTOM) RESULTS ON THE NYU DATASET USING ONLY RGB IMAGES. TRAIN/TEST INDICATES THE TRAINING/TESTING DATASET USED.

| Method | ODS | OIS | AP |
|-----------------------------|-------------|-------------|-------------|
| gPb [1] (NYU/NYU) | 0.51 | 0.52 | 0.37 |
| SCG [47] (NYU/NYU) | 0.55 | 0.57 | 0.46 |
| SE [45] (NYU/NYU) | 0.60 | 0.61 | 0.56 |
| CS-SE+DSNA (NYU/NYU) | 0.62 | 0.63 | 0.60 |
| SE [45] (BSDS/NYU) | 0.55 | 0.57 | 0.46 |
| DeepContour [20] (BSDS/NYU) | 0.55 | 0.57 | 0.49 |
| CS-SE+DSNA (BSDS/NYU) | 0.57 | 0.58 | 0.51 |

produce cleaner results with preserved edge structures. In other words, it is capable of predicting the minority edges without jeopardizing the majority non-edges that make edge maps clean. Also, the computational overhead is modest as compared to SE.

To further validate the extrapolative ability of our method, we perform cross-dataset generalization tests in comparison to the competing methods with public results. The NYU/NYU results are used as baselines, see Table V. In both cases of NYU/NYU and BSDS/NYU testing, we find favorable performance, demonstrating a superior capability of generalization. Figure 10(b) shows the visual results.

V. CONCLUSION

We propose in this paper a principled method for handling data imbalance by discriminative sparse neighbor approximation. With this method, we are able to make unbiased predictions with preserved discriminative and extrapolative ability. Such predictions are made among the local data neighborhood retrieved by a cost-sensitive decision forest. Our method proves effective in diverse vision tasks at various imbalance levels, and substantially outperforms the state-of-the-arts including some deep learning methods that ignore the imbalance issue. Our method shows its great potential as an efficient and general purpose solution for imbalanced learning. Future works include pushing the framework deeper by using cascaded forests with multi-level predictions, to explore the extent to which we can achieve by simulating deep architectures.

REFERENCES

- [1] P. Arbelaez, M. Maire, C. Fowlkes, and J. Malik, “Contour detection and hierarchical image segmentation,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 33, no. 5, pp. 898–916, 2011.
- [2] X. Geng, Z.-H. Zhou, and K. Smith-Miles, “Automatic age estimation based on facial aging patterns,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 29, no. 12, pp. 2234–2240, 2007.
- [3] J. Aghajanian and S. Prince, “Face pose estimation in uncontrolled environments,” in *Proc. British Mach. Vis. Conf.*, 2009.
- [4] K.-Y. Chang, C.-S. Chen, and Y.-P. Hung, “Ordinal hyperplanes ranker with cost sensitivities for age estimation,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2011.
- [5] J. Lim, C. L. Zitnick, and P. Dollár, “Sketch tokens: A learned mid-level representation for contour and object detection,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2013.
- [6] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, “SMOTE: Synthetic minority over-sampling technique,” *J. Artif. Intell. Res.*, vol. 16, no. 1, pp. 321–357, 2002.

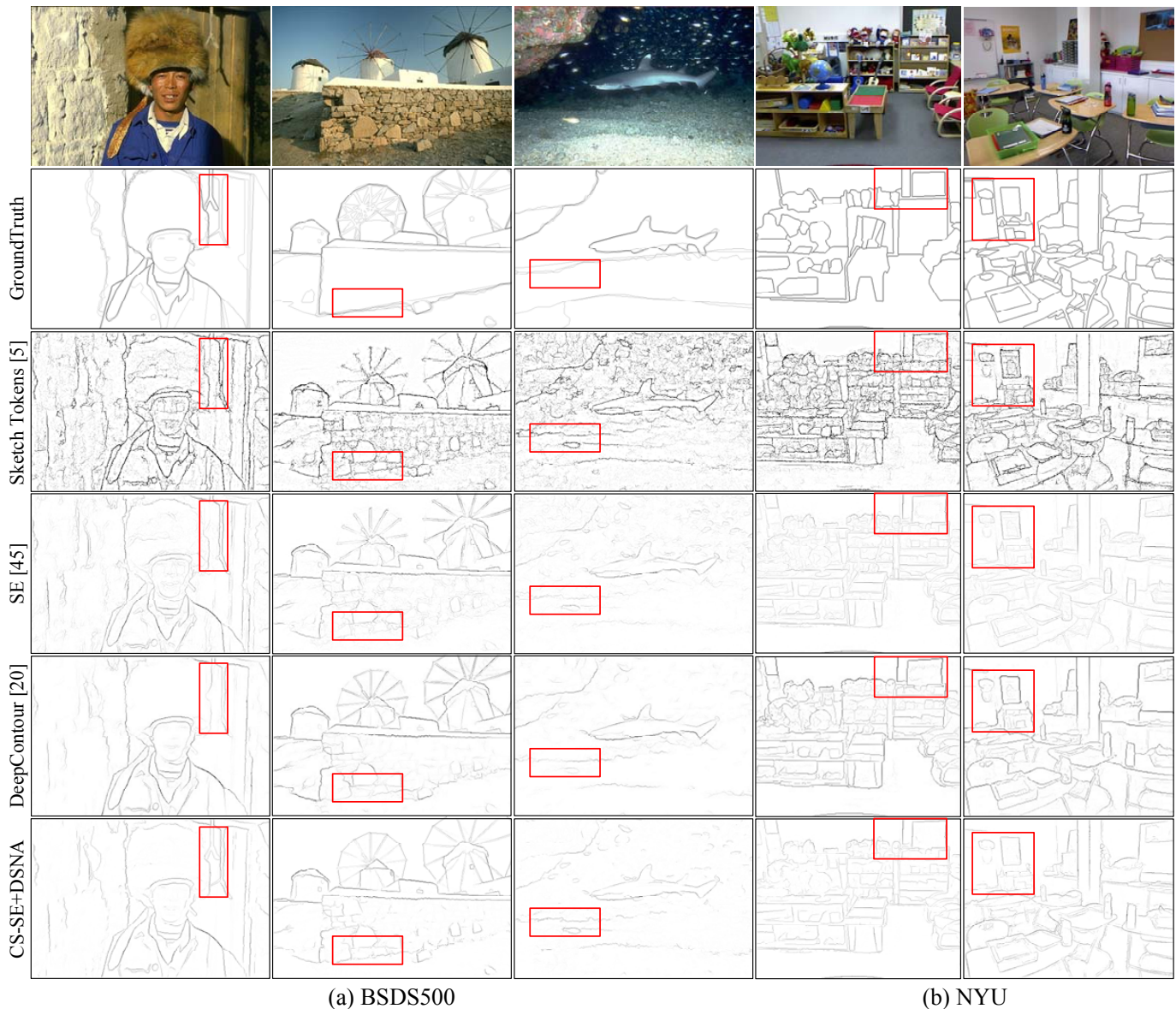
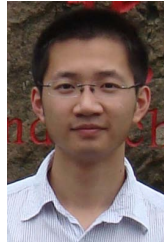


Fig. 10. Edge detection results on the (a) BSDS500 dataset and (b) NYU dataset with BSDS trained model.

- [7] H. He and E. A. Garcia, "Learning from imbalanced data," *IEEE Trans. Knowl. Data Eng.*, vol. 21, no. 9, pp. 1263–1284, 2009.
- [8] C. Drummond and R. C. Holte, "C4.5, class imbalance, and cost sensitivity: Why under-sampling beats over-sampling," in *Proc. Int. Conf. Mach. Learn. Workshops*, 2003.
- [9] L. Li and H. Lin, "Ordinal regression by extended binary classification," in *Neural Inf. Process. Syst.*, 2006.
- [10] B. Zadrozny, J. Langford, and N. Abe, "Cost-sensitive learning by cost-proportionate example weighting," in *Proc. IEEE Int. Conf. Data Mining*, 2003.
- [11] C. Chen, A. Liaw, and L. Breiman, "Using random forest to learn imbalanced data," University of California, Berkeley, Tech. Rep. 666, 2004.
- [12] K. M. Ting, "A comparative study of cost-sensitive boosting algorithms," in *Proc. Int. Conf. Mach. Learn.*, 2000.
- [13] L. Breiman, "Random forests," *J. Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001.
- [14] Y. Hu, A. Mian, and R. Owens, "Sparse approximated nearest points for image set classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2011.
- [15] S. Kong, Z. Jiang, and Q. Yang, "Learning mid-level features and modeling neuron selectivity for image classification," *arXiv preprint*, vol. arXiv:1401.5535, 2014.
- [16] D. Yi, Z. Lei, and S. Z. Li, "Age estimation by multi-scale convolutional network," in *Proc. Asian Conf. Comput. Vis.*, 2014.
- [17] J. J. Kivinen, C. K. I. Williams, and N. Heess, "Visual boundary prediction: A deep neural prediction network and quality dissection," in *Proc. Int. Conf. Artif. Intell. Stats.*, 2014.
- [18] Y. Ganin and V. S. Lempitsky, "N4-fields: Neural network nearest neighbor fields for image transforms," in *Proc. Asian Conf. Comput. Vis.*, 2014.
- [19] G. Bertasius, J. Shi, and L. Torresani, "Deepedge: A multi-scale bifurcated deep network for top-down contour detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015.
- [20] W. Shen, X. Wang, Y. Wang, and X. Bai, "Deepcontour: A deep convolutional feature learned by positive-sharing loss for contour detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015.
- [21] G. Bertasius, J. Shi, and L. Torresani, "High-for-low and low-for-high: Efficient boundary detection from deep object features and its applications to high-level vision," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015.
- [22] S. Xie and Z. Tu, "Holistically-nested edge detection," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2015.
- [23] Y. Zhang, P. Fu, W. Liu, and G. Chen, "Imbalanced data classification based on scaling kernel-based support vector machine," *Neural Comput. App.*, vol. 25, no. 3-4, pp. 927–935, 2014.

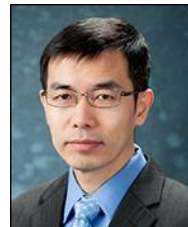
- [24] K. Li, X. Kong, Z. Lu, L. Wenyin, and J. Yin, "Boosting weighted ELM for imbalanced learning," *Neurocomputing*, vol. 128, pp. 15–21, 2014.
- [25] H. Ren and Z.-N. Li, "Age estimation based on complexity-aware features," in *Proc. Asian Conf. Comput. Vis.*, 2014.
- [26] S. E. Choi, Y. J. Lee, S. J. Lee, K. R. Park, and J. Kim, "Age estimation using a hierarchical classifier based on global and local facial features," *Pattern Recognit.*, vol. 44, no. 6, pp. 1262–1281, 2011.
- [27] G. Guo and G. Mu, "Simultaneous dimensionality reduction and human age estimation via kernel partial least squares regression," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2011.
- [28] —, "Joint estimation of age, gender and ethnicity: CCA vs. PLS," in *Proc. IEEE Int. Conf. Autom. Face and Gesture Recognit.*, 2013.
- [29] Y. Zhang and D.-Y. Yeung, "Multi-task warped gaussian process for personalized age estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2010.
- [30] W.-L. Chao, J.-Z. Liu, and J.-J. Ding, "Facial age estimation based on label-sensitive learning and age-oriented regression," *Pattern Recognit.*, vol. 46, no. 3, pp. 628–641, 2013.
- [31] S. Yan, H. Wang, T. Huang, Q. Yang, and X. Tang, "Ranking with uncertain labels," in *Proc. IEEE Int. Conf. Multimedia Expo.*, 2007.
- [32] K.-Y. Chang and C.-S. Chen, "A learning framework for age rank estimation based on face images with scattering transform," *IEEE Trans. Image Process.*, vol. 24, no. 3, pp. 785–798, 2015.
- [33] G. Guo, G. Mu, Y. Fu, and T. Huang, "Human age estimation using bio-inspired features," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2009.
- [34] K. Chen, S. Gong, T. Xiang, and C. Loy, "Cumulative attribute space for age and crowd density estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2013.
- [35] X. Geng, C. Yin, and Z.-H. Zhou, "Facial age estimation by learning from label distributions," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 10, pp. 2401–2412, 2013.
- [36] C. Huang, X. Ding, and C. Fang, "Head pose estimation based on random forests for multiclass classification," in *Proc. Int. Conf. Pattern Recognit.*, 2010.
- [37] V. Drouard, S. Ba, G. Evangelidis, A. Deleforge, and R. Horaud, "Head pose estimation via probabilistic high-dimensional regression," in *Proc. IEEE Int. Conf. Image Process.*, 2015.
- [38] M. Fenzi, L. Leal-Taixe, B. Rosenhahn, and J. Ostermann, "Class generative models based on feature regression for pose estimation of object categories," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2013.
- [39] M. Haj, J. Gonzalez, and L. Davis, "On partial least squares in head pose estimation: How to simultaneously deal with misalignment," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2012.
- [40] K. Hara and R. Chellappa, "Growing regression forests by classification: Applications to object pose estimation," in *Proc. Eur. Conf. Comput. Vis.*, 2014.
- [41] Y. Liu, Q. Wang, Y. Jiang, and Y. Lei, "Supervised locality discriminant manifold learning for head pose estimation," *Knowl. Based Syst.*, vol. 66, pp. 126–135, 2014.
- [42] E. Murphy-Chutorian and M. Trivedi, "Head pose estimation in computer vision: A survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2009.
- [43] G. Fanelli, M. Dantone, J. Gall, A. Fossati, and L. Gool, "Random forests for real time 3D face analysis," *Int. J. Comput. Vis.*, vol. 101, no. 3, pp. 437–458, 2013.
- [44] P. Arbeláez, J. Pont-Tuset, J. Barron, F. Marques, and J. Malik, "Multi-scale combinatorial grouping," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2014.
- [45] P. Dollár and C. L. Zitnick, "Fast edge detection using structured forests," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2015.
- [46] P. Isola, D. Zoran, D. Krishnan, and E. H. Adelson, "Crisp boundary detection using pointwise mutual information," in *Proc. Eur. Conf. Comput. Vis.*, 2014.
- [47] X. Ren and L. Bo, "Discriminatively Trained Sparse Code Gradients for Contour Detection," in *Neural Inf. Process. Syst.*, 2012.
- [48] Z. Ren and G. Shakhnarovich, "Image segmentation by cascaded region agglomeration," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2013.
- [49] S. Hallman and C. C. Fowlkes, "Oriented edge forests for boundary detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015.
- [50] D. Bertsekas, A. Nedic, and A. Ozdaglar, "Convex analysis and optimization," *Athena Scientific*, 2003.
- [51] T. Cootes, G. Edwards, and C. Taylor, "Active appearance models," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 23, no. 6, pp. 681–685, 2001.
- [52] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from RGBD images," in *Proc. Eur. Conf. Comput. Vis.*, 2012.
- [53] H. Han, C. Otto, and A. Jain, "Age estimation from face images: Human vs. machine performance," in *Proc. Int. Conf. Biometrics*, 2013.



Chen Huang received the Ph.D. degree in Electronic Engineering from Tsinghua University, Beijing, China, in 2014. He is currently a postdoctoral research associate in the Department of Information Engineering of the Chinese University of Hong Kong. His research interests include computer vision, pattern recognition and image processing, with focus on deep learning, face analysis and recognition.



Chen Change Loy received the PhD degree in Computer Science from the Queen Mary University of London in 2010. He is currently a Research Assistant Professor in the Department of Information Engineering, Chinese University of Hong Kong. Previously he was a postdoctoral researcher at Vision Semantics Ltd from 2011–2013. His research interests include computer vision and pattern recognition, with focus on face analysis, deep learning, and visual surveillance.



Xiaoou Tang received the B.S. degree from the University of Science and Technology of China, Hefei, in 1990, and the M.S. degree from the University of Rochester, Rochester, NY, in 1991. He received the Ph.D. degree from the Massachusetts Institute of Technology, Cambridge, in 1996. He is a Professor and the Chairman of the Department of Information Engineering. He worked as the group manager of the Visual Computing Group at the Microsoft Research Asia from 2005 to 2008. His research interests include computer vision, pattern recognition, and video processing. Dr. Tang received the Best Paper Award at the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) 2009. He is a program chair of the IEEE International Conference on Computer Vision (ICCV) 2009 and has served as an Associate Editor of IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI) and International Journal of Computer Vision (IJCV). He is a Fellow of IEEE.