

Correction of Noisy Labels via Mutual Consistency Check

Sahey Bhadra^{a,b,*}, Matthias Hein^b

^a*Max-Planck Institute for Informatics, Saarbrücken, Germany*

^b*Saarland University, Saarbrücken, Germany*

Abstract

Label noise can have severe negative effects on the performance of a classifier. Such noise can either arise by adversarial manipulation of the training data or from unskilled annotators frequently encountered in crowd sourcing (e.g. Amazon mechanical turk). Based on the assumption that an expert has provided some fraction of the training data, where labels can be assumed to be true, we propose a new pre-processing method to identify and correct noisy labels via a mutual consistency check using a Parzen window classifier. While the resulting optimization problem turns out to be a combinatorial problem, we design an efficient algorithm for which we provide approximation guarantees. Extensive experimental evaluation shows that our method performs similar and often much better than existing methods for the detection of noisy labels, thus leading to a boost in performance of the resulting classifiers.

Keywords: Noisy Annotation, Label Correction, Mutual Consistency, Parzen Window Estimation, Non-Convex Optimization and Spannogram Framework

1. Introduction

Labeled training data is essential for supervised classification. As annotating large datasets can be time consuming, nowadays often crowd sourcing [1] (e.g.

*Corresponding author

Email addresses: sahely@mpi-inf.mpg.de (Sahey Bhadra), hein@cs.uni-saarland.de (Matthias Hein)

URL: <http://www.mpi-inf.mpg.de/~sahely/> (Sahey Bhadra), <http://www.ml.uni-saarland.de/people/hein.htm> (Matthias Hein)

4 Amazon mechanical turk) is used as a quick and cost effective solution. However,
5 annotations acquired by crowd sourcing are generally contaminated with noise.
6 Often it happens that some annotators do not understand the task correctly
7 and thus provide wrong labels. Even more severe is adversarial manipulation
8 of the training data to change the classifier in a “maximal” way. It is obvious
9 that these different types of label noise can have a significant adverse effect on
10 the classification performance. In the literature many negative results [2, 3, 4]
11 have been shown regarding hardness of learning under adversarial or malicious
12 noise.

13 Being able to cope with such label noise is therefore an important practical
14 problem which has recently attracted a lot of attention. One can identify two
15 major directions among prior work. The first one attempts to correct mislabeled
16 examples during model building [5, 6, 7, 8, 9], while the second one applies
17 noise filtering as a pre-processing step prior to the model building [10, 11, 12].
18 Typically, pre-processing techniques tend to be less prone to over-fitting as they
19 are independent of the final classifier.

20 Our approach follows the second direction for label correction and involves
21 maximization of a global consistency criterion which predicts the label of a
22 training data point based on its neighbors using a Parzen window type approach.
23 This, in spirit is close to the work of [10]. However, unlike their approach,
24 we enforce hard decisions, that is either a label is wrong or right. For model
25 selection we assume that a small fraction of the training annotations has been
26 provided by an expert for which we assume that all the labels are correct.

27 Similar to other approaches for label noise correction [10, 5, 6, 7], the re-
28 sulting optimization approach is non-convex. In our case, it is an NP-hard
29 combinatorial problem. However, it turns out that in the particular setting we
30 are working, we can develop an algorithm based on the Spannogram technique
31 [13, 14], for which we can provide quite tight approximation guarantees. More-
32 over, we show in the experiments that our optimization technique outperforms
33 standard methods based on sequential linearization, which are often employed
34 in machine learning. Thus we think that, modifications of the Spannogram

35 technique could also be of potential interest in other areas of machine learning
 36 which deal with some other combinatorial problem.

37 We show on a large number of datasets with different types of noise that
 38 the proposed method is able to detect noisy labels with high precision and
 39 recall. This even holds up to a point where more than 40% of the training
 40 data-points are noisy. Finally, we show that, two different classifiers (SVM and
 41 Parzen window) trained using data pre-processed by our technique outperform
 42 or atleast performs similarly to the classifiers trained using data pre-processed
 43 by other existing techniques and also their robust counterpart.

44 2. Related Work

45 Most of the existing work for detecting noisy labels as pre-processing tech-
 46 niques such as [2, 15, 12] uses some kind of local learning. They learn a few
 47 local classifiers from sub-sampled training data and then try to detect the cor-
 48 rect label of a training data point via a majority vote among the local classifiers.
 49 The main problem of such greedy approaches is that they examine each training
 50 data point individually without using the mislabeling information of the sub-
 51 sampled set. When the observed labels of a large portion of the sub-sampled
 52 set are noisy then the local classifiers based on them can be completely wrong
 53 and will in the worst case insert even more noise.

Given the training dataset $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ such that $\mathbf{x}_i \in \mathcal{X} \subset \mathbb{R}^d$ and
 $y_i \in \mathcal{Y} = \{1, -1\}$ a global approach to the problem has been put forward by [10],
 where they assume that each labeled data point x_i is noisy with an unknown
 probability $p_i \in [0, 1]$. Hence the expected class label of the i^{th} data point is
 $E[y_i] = (1 - p_i)y_i + p_i(-y_i) = (1 - 2p_i)y_i$. Given a non-negative kernel function
 $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}_+$, the predicted value of \mathbf{x}_i based on a Parzen window type
 classifier is

$$f(\mathbf{x}_i) = \frac{\sum_{j=1}^n (1 - 2p_j)y_j K(\mathbf{x}_i, \mathbf{x}_j)}{\sum_{j=1}^n K(\mathbf{x}_i, \mathbf{x}_j)}.$$

54 Finally, [10] suggests a criterion to find p_i in order to maximize the label
 55 consistency with respect to the expected labels i.e $\sum_i E[y_i]f(\mathbf{x}_i)$ for the whole

56 dataset. This leads to the following optimization problem

$$\max_{0 \leq p_i \leq 1} \langle (1 - 2\mathbf{p}), Q^{norm}(1 - 2\mathbf{p}) \rangle - C \|\mathbf{p}\|_1 \quad (1)$$

57 where $Q_{ij}^{norm} = \frac{y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)}{\sum_j K(\mathbf{x}_i, \mathbf{x}_j)}$. The regularization term $\|\mathbf{p}\|_1$ is added in order
 58 to enforce sparsity in the solution. Moreover, for a non-negative kernel function
 59 without regularization term the optimization problem will result into a trivial
 60 solution where

$$p_i = \begin{cases} 0 & \text{if } y_i = 1 \\ 1 & \text{otherwise} \end{cases} \quad \text{or} \quad p_i = \begin{cases} 1 & \text{if } y_i = 1 \\ 0 & \text{otherwise} \end{cases}$$

61 Note that, Q^{norm} is not necessarily a negative semi definite matrix and hence
 62 (1) is a non-convex problem.

63 **Notation.** $\|\cdot\|_p = l_p$ norm, $|\mathbf{I}|$ =cardinality of set \mathbf{I} , the eigenvalues λ_j of Q
 64 with corresponding eigenvector \mathbf{v}_j are in decreasing order, $\mathbf{1}$ is the vector of all
 65 ones, $[z] =$ the largest integer $\leq z$, $\langle \mathbf{x}, \mathbf{v} \rangle = \sum_i x_i v_i$ and e_i is a vector with 1
 66 in the i^{th} position and 0 else.

67 3. Label Noise Detection by Mutual Consistency Check

Given a training dataset $\hat{\mathcal{D}} = \{\mathbf{x}_i, y_i^t\}_{i=1}^n$, where $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i^t \in \{-1, 1\}$
 is the true label of the i -th data point, the Parzen window classifier [16], f^t is
 defined as

$$f^t(\mathbf{x}_i) = \frac{\sum_{j=1}^n y_j^t K_h(\mathbf{x}_i, \mathbf{x}_j)}{\sum_j K_h(\mathbf{x}_i, \mathbf{x}_j)}.$$

68 Typically, K_h is chosen to be the Gaussian kernel, $K_h(x, y) = \exp(-\frac{\|x-y\|_2^2}{2h})$,
 69 which is non-negative and positive semi-definite. Here h is the bandwidth of
 70 the kernel. Note that f^t takes values in $[-1, 1]$. Hence the loss in terms of true
 71 labels can be defined as $L(y_i^t, f^t(\mathbf{x}_i)) = 1 - y_i^t f^t(\mathbf{x}_i)$. Finally, we use a weighted
 72 loss over $\hat{\mathcal{D}}$ where the idea is that we penalize errors more in regions of high
 73 density (where we have a lot of nearby points and thus can be more sure that

74 the simple Parzen window classifier is correct) than in regions of low density,

$$\text{Loss}(f^t, \hat{\mathcal{D}}) = \frac{1}{n} \sum_{i=1}^n w_i L(y_i^t, f^t(\mathbf{x}_i)), \quad (2)$$

75 where the weight is $w_i = \frac{1}{n} \sum_{j=1}^n K_h(\mathbf{x}_i, \mathbf{x}_j)$. After proper rescaling of w_i , it is a
 76 consistent¹ density estimator [17], i.e., $\frac{w_i}{h^d} = \frac{1}{nh^d} \sum_{j=1}^n K_h(\mathbf{x}_i, \mathbf{x}_j)$ is a consistent
 77 density estimator: $\frac{w_i}{h^d} \rightarrow p(x_i)$ if $h \rightarrow 0$, $n \rightarrow \infty$, and $nh^d \rightarrow \infty$.

78 In this paper we assume that we do not know the true labels and the observed
 79 labels may be noisy in the sense that some of the given labels are different from
 80 the true labels. More precisely, given noisy training data $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^n$ we
 81 assume that the annotation y_i is a perturbed version of the true label y_i^t . For
 82 binary classification problems, the most intuitive noise model is $y_i = \eta_i y_i^t$, where
 83 $\eta_i \in \{1, -1\}$. $\eta_i = 1$ indicates a correctly observed label and $\eta_i = -1$ indicates
 84 a noisy label. The goal is to find y^t and hence η by minimizing the global loss
 85 in (2) with respect to η . We have

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n w_i L(y_i^t, f^t(\mathbf{x}_i)) &= \frac{1}{n^2} \sum_{i,j=1}^n K_h(\mathbf{x}_i, \mathbf{x}_j) - \frac{1}{n^2} \sum_{i,j=1}^n y_i^t y_j^t K_h(\mathbf{x}_i, \mathbf{x}_j) \\ &= \frac{1}{n^2} \sum_{i,j=1}^n K_h(\mathbf{x}_i, \mathbf{x}_j) - \frac{1}{n^2} \sum_{i,j=1}^n \eta_i y_i \eta_j y_j K_h(\mathbf{x}_i, \mathbf{x}_j) \end{aligned} \quad (3)$$

86 Here $\frac{1}{n^2} \sum_{i,j=1}^n K_h(\mathbf{x}_i, \mathbf{x}_j)$ is constant. Hence in order to maximize mutual
 87 consistency in the labels of all training data-points, we have to minimize the
 88 loss (3) and hence maximize the following optimization problem

$$\eta^* = \arg \max_{\eta \in \{-1, 1\}^n} \sum_{i,j=1}^n \eta_i \eta_j y_i y_j K_h(\mathbf{x}_i, \mathbf{x}_j). \quad (4)$$

89 Unlike greedy approaches, above optimization problem detects all the noisy
 90 points simultaneously.

91 Note that, the optimization problem (4) has a trivial solution $\eta_i = y_i$ and
 92 hence is of no interest. The problem becomes challenging as we have two kinds
 93 of extra information on the problem and we encode them as constraints. First,

¹Note that, consistent has been mentioned earlier in a different sense.

94 we make the natural assumption that a small part of the data set is annotated
 95 by experts and the corresponding labels are always correct. Thus the labels
 96 given by experts are fixed during the optimization. Let us define \mathbf{I}_E as the set of
 97 indices corresponding to data points annotated by expert annotators. Then $\eta_i =$
 98 $1, \forall i \in \mathbf{I}_E$. The second kind of extra information are the fractions of data points
 99 with noisy labels given by ρ_+ and ρ_- for positive and negative annotated classes
 100 respectively. Hence $\rho_+ = \frac{|\{i|y_i=1 \text{ and } y_i^t=-1\}|}{n_+}$ and $\rho_- = \frac{|\{i|y_i=-1 \text{ and } y_i^t=1\}|}{n_-}$ where
 101 $n_+ = |I_+| = |\{i|y_i = 1\}|$ and $n_- = |I_-| = |\{i|y_i = -1\}|$. The fractions of noisy
 102 labels ρ_+ and ρ_- can be estimated similar to [7], but in this paper we estimate
 103 them by cross-validation using the knowledge of the expert labels. The final
 104 constraint set of our problem (4) is defined as

$$\mathcal{E} = \{\eta \mid \eta \in \{1, -1\}^n, \langle \mathbf{1}, \eta_{\mathbf{I}_+} \rangle = n_+ - 2\delta_+, \langle \mathbf{1}, \eta_{\mathbf{I}_-} \rangle = n_- - 2\delta_-, \eta_{\mathbf{I}_E} = 1\} \quad (5)$$

105 where $\delta_+ = \lfloor \rho_+ n_+ \rfloor$ and $\delta_- = \lfloor \rho_- n_- \rfloor$. Hence our Label Noise Detection
 106 (**LND**) method solves the following optimization problem.

$$\mathbf{LND}: \quad \eta^* = \arg \max_{\eta \in \mathcal{E}} \langle \eta, Q\eta \rangle \quad (6)$$

107 where, $Q_{ij} = y_i y_j K_h(\mathbf{x}_i, \mathbf{x}_j)$. Along with the constraint (5) the solution of (6)
 108 is no more trivial and we present an efficient algorithm to solve this problem in
 109 Section 4.2.

110 4. Algorithms

111 The optimization problem **LND** is a combinatorial problem and NP-hard in
 112 general. We present two possible ways to solve it approximately. The first one
 113 linearizes the objective function in each step and then solves the corresponding
 114 constrained maximization problem in closed form and the second one is an
 115 algorithm using the Spannogram technique for which we provide approximation
 116 guarantees.

117 4.1. Sequential Linearization

118 Algorithm 1 linearizes the objective function in the k -th step as $\langle \eta, Q\eta^k \rangle$
 119 and then solves the corresponding maximization problem, $\max_{\eta \in \mathcal{E}} \langle \eta, Q\eta^k \rangle$. It

120 turns out that the solution is a simple projection onto the constraint set \mathcal{E} .

121 **LND** solved using Algorithm 1 is denoted as **LND**_{slp}.

122 **Lemma 1.** $\arg \max_{\eta \in \mathcal{E}} \langle \eta, \mathbf{v} \rangle \equiv \Pi_{\mathcal{E}}(\mathbf{v})$ where $\Pi_{\mathcal{E}}$ is projection on \mathcal{E} .

123 *Proof.* $\forall \eta \in \mathcal{E}$, $\|\eta\|^2 = n$ is constant and hence

$$\Pi_{\mathcal{E}}(\mathbf{v}) = \arg \min_{\eta \in \mathcal{E}} \|\eta - \mathbf{v}\|^2 \equiv \arg \max_{\eta \in \mathcal{E}} \langle \eta, \mathbf{v} \rangle.$$

124 □

125 Moreover, it turns out that the projection onto the discrete set \mathcal{E} can be
126 easily computed in closed form using Algorithm 2 .

127 **Theorem 1.** *Algorithm 2 computes $\Pi_{\mathcal{E}}(\mathbf{v})$.*

Proof. We give a proof by contradiction. Let us assume that η , the outcome from Algorithm 2, is not equal to $\Pi_{\mathcal{E}}(\mathbf{v})$ and there exists an $\eta^* \neq \eta$ such that $\eta^* = \Pi_{\mathcal{E}}(\mathbf{v})$. Hence there are at-least two indices (j, l) such that $y_j = y_l$ but $\eta_j^* = -\eta_j$, $\eta_l^* = -\eta_l$ and $\eta_l = -\eta_j$. Without loss of generality let us assume that $v_l \geq v_j$ and hence according to Algorithm 2 $\eta_l = 1$, $\eta_j = -1$ and hence according to above assumption $\eta_l^* = -1$ and $\eta_j^* = 1$. Assuming $\eta_i^* = \eta_i \forall i \neq j, l$, we get

$$\|\eta - \mathbf{v}\|^2 - \|\eta^* - \mathbf{v}\|^2 = -2(\eta_l v_l + \eta_j v_j) + 2(\eta_l^* v_l + \eta_j^* v_j) = 4(v_j - v_l).$$

128 Now, as $v_l \geq v_j$, $4(v_j - v_l) \leq 0$ or $\|\eta - \mathbf{v}\|^2 \leq \|\eta^* - \mathbf{v}\|^2$. Hence η must be equal
129 to $\Pi_{\mathcal{E}}(\mathbf{v})$. □

130 Finally, we can show that Algorithm 1 leads to monotonic ascent.

131 **Lemma 2.** *Algorithm 1 at each step provides a feasible η^{k+1} with monotonically
132 increasing function value if K_h is positive definite.*

Proof. If K_h is positive definite, then Q is also positive definite. Using the first-order condition of convex functions, we get $\langle \eta, Q\eta \rangle \geq \langle \eta^k, Q\eta^k \rangle + \frac{1}{2} \langle Q\eta^k, \eta - \eta^k \rangle$. Thus the algorithmic scheme maximizes a lower bound on the objective. Moreover, as η^k is feasible, we get

$$\langle \eta^{k+1}, Q\eta^{k+1} \rangle \geq \langle \eta^k, Q\eta^k \rangle + \frac{1}{2} \langle Q\eta^k, \eta^{k+1} - \eta^k \rangle \geq \langle \eta^k, Q\eta^k \rangle.$$

Algorithm 1 LND_{slp}

Initialization: Randomly take η^0 in \mathcal{E} and $k := 0$.

Output: η^k

repeat

Iteration k: $\eta^{k+1} = \Pi_{\mathcal{E}}(Q\eta^k)$.

until $\frac{\langle \eta^{k+1}, Q\eta^{k+1} \rangle - \langle \eta^k, Q\eta^k \rangle}{\langle \eta^k, Q\eta^k \rangle} \leq \epsilon$

Algorithm 2 $\Pi_{\mathcal{E}}(\mathbf{v})$

Initialization: $\eta = \mathbf{1}$ and $I_{\eta} = \emptyset$.

Output: η

$I_{\eta_+} \leftarrow \{i | \mathbf{v}_{i \in \{i | y_i = 1\}} \leq \mathbf{v}_+^{[\delta_+]}, \text{ where } \mathbf{v}_+^{[\delta_+]} = \delta_+ \text{-th smallest element of } \mathbf{v}_{i \in \{i | y_i = 1\}}\}$

$I_{\eta_-} \leftarrow \{i | \mathbf{v}_{i \in \{i | y_i = -1\}} \leq \mathbf{v}_-^{[\delta_-]}, \text{ where } \mathbf{v}_-^{[\delta_-]} = \delta_- \text{-th smallest element of } \mathbf{v}_{i \in \{i | y_i = -1\}}\}$

$I_{\eta} = \{I_{\eta_+}, I_{\eta_-}\}$

$\eta_{I_{\eta}} = -1$

133

□

134 *4.2. Algorithm based on low rank approximation*

135 The LND_{slp} (Algorithm 1) has the problem that it can get stuck in local
136 optima without any approximation guarantees. On the other-hand, finding the
137 global optimum of (6) is NP-hard [18]. In this paper we propose an algorithm
138 based on the Spannogram framework [13, 14] which allows to solve **LND** with
139 a certain approximation guarantee.

140 **Efficient use of \mathbf{I}_E :** Using ($\eta_i = 1, \forall i \in I_E$),

$$\langle \eta, Q\eta \rangle = \sum_{i,j \in I_E} Q_{ij} + 2 \sum_{i \notin I_E, j \in I_E} \eta_i Q_{ij} + \sum_{i,j \notin I_E} \eta_i \eta_j Q_{ij}.$$

141 Hence (6) can be solved by solving

$$\eta^{s*} = \arg \max_{\eta \in \mathcal{E}_s} \langle \eta_s, Q_s \eta_s \rangle \quad (7)$$

142 where, $Q_s = \begin{bmatrix} \langle \mathbf{1}, Q_{i \in I_E} \ j \in I_E \mathbf{1} \rangle & \langle \mathbf{1}, Q_{i \in I_E} \ j \notin I_E \rangle \\ Q_{i \notin I_E} \ j \in I_E \mathbf{1} & Q_{i \notin I_E} \ j \notin I_E \end{bmatrix}$ and the new feasibility set
143 is defined as, $\mathcal{E}_s = \{\eta | \eta_i \in \{+1, -1\}, |\{i | \eta_i = -1 \text{ and } y_i = 1\}| = \delta_+, |\{i | \eta_i =$
144 $-1 \text{ and } y_i = -1\}| = \delta_- \text{ and } \eta_1 = 1\}$. Finally we can get back the solution of
145 (6) by assigning $\eta_{i \in I_E}^* = \eta_1^{s*} = 1$ and $\eta_{i \notin I_E}^* = \eta_{j \geq 2}^{s*}$. As structure of problem (7)
146 is the same as that of (6), henceforth we will consider $Q = Q_s$ and $\mathcal{E} = \mathcal{E}_s$.

147 For a positive semi-definite(PSD) kernel K_h , Q is also PSD hence using
148 eigenvalue decomposition, (6) is equivalent to

$$\arg \max_{\eta \in \mathcal{E}} \langle \eta, Q_n \eta \rangle, \quad \text{where} \quad Q_n = \sum_{j=1}^n \lambda_j \mathbf{v}_j \mathbf{v}_j^T \quad (8)$$

149 where n is number of training data points. A low rank approximation of (6)
150 and (8) is given by

$$\arg \max_{\eta \in \mathcal{E}} \langle \eta, Q_r \eta \rangle, \quad \text{where} \quad Q_r = \sum_{j=1}^r \lambda_j \mathbf{v}_j \mathbf{v}_j^T \quad (9)$$

where Q_r is the low rank approximation of Q and ideally $r \ll n$. This in turn
is equivalent to

$$\arg \max_{\eta \in \mathcal{E}} \|\mathbf{V}_r \eta\|_2^2, \quad \text{where} \quad \mathbf{V}_r = [\sqrt{\lambda_1} \mathbf{v}_1, \dots, \sqrt{\lambda_r} \mathbf{v}_r]^T.$$

151 Let \mathbf{c} be a $r \times 1$ unit length vector, i.e., $\|\mathbf{c}\|_2 = 1$. Using Cauchy-Schwarz
152 inequality, we get $\langle \mathbf{c}, \mathbf{V}_r \eta \rangle^2 \leq \|\mathbf{V}_r \eta\|_2^2$; with equality, if and only if, \mathbf{c} is co-linear
153 to $\mathbf{V}_r \eta$. Hence (9) is equivalent to

$$\arg \max_{\eta \in \mathcal{E}} \max_{\mathbf{c}: \|\mathbf{c}\|_2=1} \langle \mathbf{c}, \mathbf{V}_r \eta \rangle^2. \quad (10)$$

For a given \mathbf{c}^* , defining $\mathbf{v}_r^{c^*} = \mathbf{V}_r^T \mathbf{c}^*$, (10) is equivalent to

$$\arg \max_{\eta \in \mathcal{E}} \langle \mathbf{c}^*, \mathbf{V}_r \eta \rangle^2 = \arg \max_{\eta \in \mathcal{E}} \left| \langle \mathbf{v}_r^{c^*}, \eta \rangle \right|.$$

154 **Lemma 3.** $\arg \max_{\eta \in \mathcal{E}} |\langle \mathbf{v}_r^c, \eta \rangle| \in \mathbf{S}_{(r, \mathcal{E})}^c$, where $\mathbf{S}_{(r, \mathcal{E})}^c = \{\Pi_{\mathcal{E}}(\mathbf{v}_r^c), \Pi_{\mathcal{E}}(-\mathbf{v}_r^c)\}$.

155 *Proof.* We have, $\arg \max_{\eta \in \mathcal{E}} |\langle \mathbf{v}_r^c, \eta \rangle| \in \{\arg \max_{\eta \in \mathcal{E}} \langle \mathbf{v}_r^c, \eta \rangle, \arg \max_{\eta \in \mathcal{E}} \langle (-\mathbf{v}_r^c), \eta \rangle\}$. Now,
156 using Lemma 1 we get, $\arg \max_{\eta \in \mathcal{E}} |\langle \mathbf{v}_r^c, \eta \rangle| \in \{\Pi_{\mathcal{E}}(\mathbf{v}_r^c), \Pi_{\mathcal{E}}(-\mathbf{v}_r^c)\} = \mathbf{S}_{(r, \mathcal{E})}^c$. \square

157 Now for solving (10) the remaining part is to find \mathbf{c}^* such that

$$\mathbf{c}^* = \arg \max_{\mathbf{c}: \|\mathbf{c}\|^2=1} \max_{\eta \in \mathbf{S}_{(r, \mathcal{E})}^c} \langle \mathbf{c}, \mathbf{V}_r \eta \rangle^2. \quad (11)$$

158 For $r = 1$, that is $c \in \mathbb{R}$, there are only two feasible \mathbf{c} such that $\|\mathbf{c}\|^2 = 1$ and
 159 they are given by $\mathbf{c} = \pm 1$ and hence instead of solving (11) one can directly
 160 solve (10) by

$$\arg \max_{\eta \in \mathcal{E}} \max_{\mathbf{v}_r^c \in \pm \sqrt{\lambda_1} \mathbf{v}_1^T} \langle \mathbf{v}_r^c, \eta \rangle^2.$$

161 Similarly, for $r \geq 2$ instead of solving (11) from an infinitely large set of
 162 feasible \mathbf{c} we find a finite set of potential \mathbf{c} denoted as \mathcal{C}_r such that

$$\arg \max_{\eta \in \mathcal{E}} \max_{\mathbf{c}: \|\mathbf{c}\|^2=1} \langle \mathbf{c}, \mathbf{V}_r \eta \rangle^2 \in \cup_{\mathbf{c} \in \mathcal{C}_r} \arg \max_{\eta \in \mathbf{S}_{(r, \mathcal{E})}^c} \langle \mathbf{c}, \mathbf{V}_r \eta \rangle^2.$$

163 Please note that for a fixed \mathbf{c} , $\arg \max_{\eta \in \mathcal{E}} \langle \mathbf{c}, \mathbf{V}_r \eta \rangle^2$ is solved using $\Pi_{\mathcal{E}}$ (Lemma
 164 3) where $\Pi_{\mathcal{E}}$ (or Algorithm 2) uses the sorted order of elements of \mathbf{v}_r^c . Again
 165 an $\eta \in \mathcal{E}$ can be a potential solution of (9) only if there exist a \mathbf{c} for which
 166 $\eta \in \Pi_{\mathcal{E}}(\mathbf{v}_r^c)$. Hence it is enough to build \mathcal{C}_r which contains all \mathbf{c} which generate
 167 a different sorted order of elements of \mathbf{v}_r^c .

168 4.2.1. The Spannogram framework

169 The key idea of the Spannogram framework [13, 14] is the introduction of
 170 spherical coordinates. For any $r \geq 2$ this transformation can be done by using
 171 $r - 1$ phase variable $\Phi = [\phi_1, \dots, \phi_{r-1}] \in [[-\frac{\pi}{2}, \frac{\pi}{2}]^{(r-2)}, [-\pi, \pi]]$ by expressing \mathbf{c}
 172 without loss of generality as

$$\mathbf{c} = \begin{bmatrix} \sin(\phi_1) \\ \cos(\phi_1) \sin(\phi_2) \\ \cos(\phi_1) \cos(\phi_2) \sin(\phi_3) \\ \dots \\ \cos(\phi_1) \cos(\phi_2) \dots \sin(\phi_{r-1}) \\ \cos(\phi_1) \cos(\phi_2) \dots \cos(\phi_{r-1}) \end{bmatrix} \quad (12)$$

173 which is a vector of unit norm and for all ϕ it produces all $r \times 1$ unit vectors.
 174 Under this transformation \mathbf{v}_r^c can be expressed in terms of ϕ as

$$\begin{aligned} \mathbf{v}(\Phi) = & \sin(\phi_1)[\sqrt{\lambda_1}\mathbf{v}_1] + \cos(\phi_1)\sin(\phi_2)[\sqrt{\lambda_2}\mathbf{v}_2] + \dots \\ & + \cos(\phi_1)\cos(\phi_2)\dots\cos(\phi_{r-1})[\sqrt{\lambda_r}\mathbf{v}_r] \end{aligned} \quad (13)$$

where each element of $[\mathbf{v}(\Phi)]_i$ is continuous function of $r - 1$ variables Φ . Calculating $\Pi_{\mathcal{E}}(\mathbf{v}_r^c)$ for a fixed vector \mathbf{c} is equivalent to finding the relative sorting of the n surfaces $[\mathbf{v}(\Phi)]_{i=1,\dots,n}$ for corresponding Φ . If the relative ranking of i -th and j -th elements of $\mathbf{v}(\Phi_1)$ and $\mathbf{v}(\Phi_2)$ change, i.e, $[\mathbf{v}(\Phi_1)]_i > [\mathbf{v}(\Phi_1)]_j$ but $[\mathbf{v}(\Phi_2)]_i < [\mathbf{v}(\Phi_2)]_j$, then there must be a $\Phi_3 \in [\Phi_1, \Phi_2]$ such that $[\mathbf{v}(\Phi_3)]_i = [\mathbf{v}(\Phi_3)]_j$. Hence it is enough to collect all the points where any two of these surfaces intersect. Please note that, the solution of $[\mathbf{v}(\Phi)]_i = [\mathbf{v}(\Phi)]_j$ is not a single point (i.e., a single vector \mathbf{c}) but a $(r - 1)$ dimensional space of solutions denoted as

$$\Phi_{i,j} = \{\phi \mid [\mathbf{v}(\phi)]_i = [\mathbf{v}(\phi)]_j\}.$$

175 Since $\Pi_{\mathcal{E}}$ and hence local optimum changes only if the local ranking changes,
176 the intersection points defined by all $\Phi_{i,j}$ sets are the only points of interest.
177 For the vectors in this space $\Phi_{i,j}$, there are again some critical ϕ 's where both
178 the i -th and j -th elements are greater than δ_+ -th (δ_- -th) largest or less than
179 δ_+ -th (δ_- -th) smallest element of $\mathbf{v}(\phi)$ and still sorted order of element of $\mathbf{v}(\phi)$
180 changes at ϕ . This happens when both i -th and j -th elements of $\mathbf{v}(\phi)$ become
181 equal to the l -th element of $\mathbf{v}(\phi)$. This new $(r - 3)$ dimensional sub-space is
182 denoted as $\Phi_{i,j,l}$. At this point, the intersection points defined by all $\Phi_{i,j,l}$ sets
183 are the only points of interest. In this manner we can find Φ_{i_1,i_2,\dots,i_r} which
184 contains all Φ such that

$$[\mathbf{v}(\Phi)]_{i_1} = [\mathbf{v}(\Phi)]_{i_2} = \dots = [\mathbf{v}(\Phi)]_{i_r}.$$

185 Finally, the intersection points defined by all such Φ_{i_1,i_2,\dots,i_r} sets are the
186 only points of interest. Hence only \mathbf{c} corresponding to $\phi \in \Phi_{i_1,i_2,\dots,i_r}$ need to
187 be checked and can be obtained by solving the system of $r - 1$, linear equations

Algorithm 3 Computing \mathcal{C}_r corresponding to \mathbf{V}_r

Initialization: $\mathcal{C}_r = \emptyset$

Output: \mathcal{C}_r

for all class y do

for all $\binom{n_y}{r}$ subsets $\{i_1, \dots, i_r\} \subset \{i | i \in \mathbf{I}_y / \mathbf{I}_E\}$ **do**

$$\mathbf{c} = \text{nullspace} \left(\begin{bmatrix} e_{i_1}^T - e_{i_2}^T \\ \dots \\ e_{i_1}^T - e_{i_r}^T \end{bmatrix} \mathbf{V}_r^T \right)$$

sort elements of $\mathbf{v}_r^c = \mathbf{V}_r^T \mathbf{c}$

if $\exists(i \in \{i_1, \dots, i_r\})$ such that $\mathbf{v}_{r_i}^c$ (or $-\mathbf{v}_{r_i}^c$) is equal to δ_y -th smallest element of \mathbf{v}_r^c (or $-\mathbf{v}_r^c$) **then**

$$\mathcal{C}_r \leftarrow \mathcal{C}_r \cup \{\mathbf{c}\}$$

end if

end for

end for

$$\begin{bmatrix} e_{i_1}^T - e_{i_2}^T \\ \dots \\ e_{i_1}^T - e_{i_r}^T \end{bmatrix} \mathbf{v}(\Phi) = \begin{bmatrix} e_{i_1}^T - e_{i_2}^T \\ \dots \\ e_{i_1}^T - e_{i_r}^T \end{bmatrix} \mathbf{v}_r^c = \begin{bmatrix} e_{i_1}^T - e_{i_2}^T \\ \dots \\ e_{i_1}^T - e_{i_r}^T \end{bmatrix} \mathbf{V}_r^T \mathbf{c} = 0. \quad (14)$$

188 where e_i is a vector of all zeros except 1 in the i -th position.

189 Please note that, we need only the sorted order of elements corresponding to
 190 each annotation separately. Algorithm 3 finds such \mathbf{c} considering elements of \mathbf{v}_r^c
 191 corresponding to all positively and negatively annotated data points separately.
 192 Again, we will consider those \mathbf{c} for which at least one of these r elements of \mathbf{v}_r^c
 193 with equal values is equal to δ_+ -th (δ_- -th) smallest (or largest) element of \mathbf{v}_r^c .

194 Finally, **LND** solved by Algorithm 4 is denoted as **LND** $_r$ for r -rank approx-
 195 imation.

196 **Similarity with other algorithms:** The Spannogram framework is inspired
 197 by the work of [13, 14], where sparse PCA and the densest subgraph problem
 198 have been studied. Both of these problems are quadratic maximization prob-

Algorithm 4 LND_r

Input: \mathcal{E}, Q, r **Output:** η_r^* Compute $\mathbf{V}_r = [\sqrt{\lambda_1}\mathbf{v}_1, \dots, \sqrt{\lambda_r}\mathbf{v}_r]^T$.Build \mathcal{C}_r (Algorithm 3)Build $\mathbf{S}_{(r,\mathcal{E})} = \cup_{\mathbf{c} \in \mathcal{C}_r} \mathbf{S}_{(r,\mathcal{E})}^{\mathbf{c}}$ (Lemma 3) $\eta_r^* = \arg \max_{\eta \in \mathbf{S}_{(r,\mathcal{E})}} \langle \eta, Q\eta \rangle$

199 lems, like **LND** in (6). Unlike **LND** they need sparse solutions while **LND**
200 has a $\{1, -1\}$ constraint on the variables which requires a fundamental modifi-
201 cation in the theoretical analysis. [14] seems to be the closest to the proposed
202 algorithm where the Spannogram technique is used to maximize a quadratic
203 function with $\{0, 1\}$ constraint. For projecting onto the cardinality constraint
204 [14] needs a sorted order of elements of a vector similar to \mathbf{v}_r^c and to get differ-
205 ent possible rankings they use the Spannogram technique. For our problem the
206 projection onto \mathcal{E} also has a closed form solution depending on the sorted order
207 of elements. This similarity between both the problems motivate us to extend
208 the Spannogram framework for solving **LND**. We show that the Spannogram
209 type algorithms can also be used when the required solution is not sparse. The
210 proposed algorithm can also be easily extended to solve similar problems with
211 other integer constraints on η in place of $\{1, -1\}$ constraint by assigning $\eta_i \geq \eta_j$
212 when $\mathbf{v}_{r_i}^c > \mathbf{v}_{r_j}^c$ and again $\eta_i \geq \eta_j$ when $-\mathbf{v}_{r_i}^c > -\mathbf{v}_{r_j}^c$.

213 **Complexity of proposed algorithm:** For a rank- r approximation we have to
214 solve a set of $\binom{n_+}{r}$ and $\binom{n_-}{r}$ equations to find \mathcal{C}_r . Each of these equation sets will
215 add one \mathbf{c} in \mathcal{C}_r and add at-most $2\binom{r}{\frac{r}{2}}$ candidates in $\mathbf{S}_{(r,\mathcal{E})}^{\mathbf{c}}$. Hence, $|\mathbf{S}_{(r,\mathcal{E})}|$ is less
216 than $2\binom{r}{\frac{r}{2}} \left(\binom{n_+}{r} + \binom{n_-}{r} \right)$ or $O\left(\binom{\max\{n_+, n_-\}}{r}\right)$. Considering the complexity of
217 sorting $O(n \log n)$ and of final matrix multiplication $\langle \eta, Q\eta \rangle$ for every $\eta \in \mathbf{S}_{r,\mathcal{E}}$,
218 the time complexity of the proposed algorithm is $O(\max\{n_+, n_-\}^r (n^2 + n \log n))$.
219 For our experiment we have used mostly $r \leq 2$ so that the search space has been
220 reduced to $O(n^2)$ and time complexity is $O(n^4)$.

221 **5. Approximation Guarantees**

222 We have the following different problems,

$$\begin{aligned} OPT^* &= \max_{\eta \in \mathcal{E}} \langle \eta, Q\eta \rangle, & \eta^* &= \arg \max_{\eta \in \mathcal{E}} \langle \eta, Q\eta \rangle, \\ OPT_r^* &= \max_{\eta \in \mathbf{S}_{(r, \mathcal{E})}} \langle \eta, Q\eta \rangle, & \eta_r^* &= \arg \max_{\eta \in \mathbf{S}_{(r, \mathcal{E})}} \langle \eta, Q\eta \rangle, \end{aligned}$$

223 Here $\mathbf{S}_{(r, \mathcal{E})} = \cup_{\mathbf{c} \in \mathcal{C}_r} \mathbf{S}_{(r, \mathcal{E})}^{\mathbf{c}}$.

224 **Lemma 4.** If $O_g = \langle \eta_g, Q\eta_g \rangle$, where $\eta_g = \Pi_{\mathcal{E}}(q_1)$ and q_1 denotes the first row
225 of Q . Then

$$OPT^* \geq \max \left\{ O_g, \lambda_1 \langle \Pi_{\mathcal{E}}(\mathbf{v}_1), \mathbf{v}_1 \rangle^2, \lambda_1 \langle \Pi_{\mathcal{E}}(-\mathbf{v}_1), \mathbf{v}_1 \rangle^2 \right\}. \quad (15)$$

226 *Proof.* The first part follows using $OPT^* = \max_{\eta \in \mathcal{E}} \langle \eta, Q\eta \rangle \geq \langle \eta_g, Q\eta_g \rangle$, while
227 the second and third parts inside the max in (15) can be proven as follows.

$$\begin{aligned} OPT^* &\geq \max_{\eta \in \mathcal{E}} \langle \eta, Q_1\eta \rangle = \max_{\eta \in \mathcal{E}} \lambda_1 \left(\langle \eta, \mathbf{v}_1 \rangle^2 \right) \quad (\text{as } Q \text{ is PSD}) \\ &= \lambda_1 \max \left\{ \langle \Pi_{\mathcal{E}}(\mathbf{v}_1), \mathbf{v}_1 \rangle^2, \langle \Pi_{\mathcal{E}}(-\mathbf{v}_1), \mathbf{v}_1 \rangle^2 \right\}. \end{aligned} \quad (16)$$

228 □

229 **Theorem 2.** $OPT_r^* \geq (1 - \epsilon_r)OPT^*$, where

$$\epsilon_r \leq \frac{n(\lambda_{r+1} - \lambda_n)}{\max \left\{ O_g, \lambda_1 \langle \Pi_{\mathcal{E}}(\mathbf{v}_1), \mathbf{v}_1 \rangle^2, \lambda_1 \langle \Pi_{\mathcal{E}}(-\mathbf{v}_1), \mathbf{v}_1 \rangle^2 \right\}}.$$

Proof. We decompose the quadratic form in (6) in two parts

$$\langle \eta, Q\eta \rangle = \langle \eta, Q_r\eta \rangle + \langle \eta, Q_{r^c}\eta \rangle$$

230 where $Q_r = \sum_{i=1}^r \lambda_i \mathbf{v}_i \mathbf{v}_i^T$ and $Q_{r^c} = \sum_{i=r+1}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^T$. By defining, $\eta_r =$
231 $\arg \max_{\eta \in \mathbf{S}_{(r, \mathcal{E})}} \langle \eta, Q_r\eta \rangle$ and hence $OPT_r = \langle \eta_r, Q_r\eta_r \rangle = \max_{\eta \in \mathbf{S}_{(r, \mathcal{E})}} \langle \eta, Q_r\eta \rangle$,

$$\begin{aligned} OPT_r^* &\geq \langle \eta_r, Q\eta_r \rangle \quad (\text{as } \eta_r \in \mathbf{S}_{(r, \mathcal{E})}) = \langle \eta_r, Q_r\eta_r \rangle + \langle \eta_r, Q_{r^c}\eta_r \rangle \\ &\geq OPT_r + n \min_{\eta: \|\eta\|_2=1} \langle \eta, Q_{r^c}\eta \rangle \geq OPT_r + n\lambda_n \\ \Rightarrow OPT_r &\leq OPT_r^* - n\lambda_n \end{aligned} \quad (17)$$

232 Using $\max_{\eta \in \mathcal{E}} \langle \eta, Q_{r^c} \eta \rangle \leq n \max_{\eta: \|\eta\|_2=1} \langle \eta, Q_{r^c} \eta \rangle = n\lambda_{r+1}$ and (17) we get,

$$\begin{aligned}
 OPT^* &\leq \max_{\eta \in \mathcal{E}} \langle \eta, Q_r \eta \rangle + \max_{\eta \in \mathcal{E}} \langle \eta, Q_{r^c} \eta \rangle \leq OPT_r + n\lambda_{r+1} \\
 &\quad \left(\max_{\eta \in \mathcal{E}} \langle \eta, Q_r \eta \rangle = \max_{\eta \in \mathbf{S}_{(r, \mathcal{E})}} \langle \eta, Q_r \eta \rangle \text{ as } Q_r \text{ is a rank-}r \text{ matrix} \right) \\
 \Rightarrow OPT^* &\leq OPT_r^* - n\lambda_n + n\lambda_{r+1} \quad (\text{using (17)}) \tag{18}
 \end{aligned}$$

233 One can write (18) as $OPT_r^* \geq \left(1 - \frac{n(\lambda_{r+1} - \lambda_n)}{OPT^*}\right) OPT^*$. Now the lower
 234 bound of OPT^* from Lemma 4 completes the proof.

235 □

236 Note that, when a significant number of the training data points are labeled
 237 by expert annotators the first row and the first column of Q will have signifi-
 238 cantly larger values than other elements of Q . This makes the largest eigenvalue
 239 (λ_1) much greater than other eigenvalues and hence ϵ_r becomes relatively small
 240 as is also shown in Section 6.

241 6. Experiments

242 This section presents our experimental setup and the results. We apply
 243 the proposed **LND** on a variety of data-sets contaminated with different kind
 244 of label noise. The objectives of our experiments are: (1) to illustrate the
 245 improvements in classification accuracy after correcting the label noise, (2) to
 246 compare the performance of **LND** against other existing methods in terms of
 247 ability to detect wrongly annotated data points, and (3) to prove superiority
 248 of the proposed Algorithm 4 for solving the non-convex optimization problem
 249 **LND** against the popular sequential linearization based algorithm **LND**_{slp}.

250 6.1. Experimental Setup

251 **Data-sets:** To evaluate the performance of our method, we use 8 data-sets
 252 from [19] described in Table 1. All the experiments are repeated ten times on
 253 ten random partitions and then the average performance is reported.

254 The percentage of training data-set labeled by experts is fixed at 10% for
 255 most of our experiments. In our experiments we select all the parameters for

Table 1: Datasets

Name	# Data-set	Train-Test Split
Mushrooms	8124 x 112	(50%, 50%)
Svmguide1	7089 x 4	As in [19]
Fourclass	862 x 2	(50%, 50%)
Australian	619 x 14	(50%, 50%)
WDBC	569 x 30	(50%, 50%)
Heart	270 x 13	(50%, 50%)
Adult	32561 x 123	As in [19]
Covertypes	581012 x 54	(1%, 9%)

256 various methods using 5-fold cross-validation and hence consider at-least 5 data
 257 points from each class for each data-set to be labeled by expert annotators.
 258 Considering the size of our smallest data-set (Heart), we fix this number at 10%
 259 for our experiments. But to study the effect of amount of available expert labels
 260 on performance of the proposed method, we repeated all experiments with less
 261 number (1% , 2% and 5%) of expert labels for larger data-sets (Mushrooms,
 262 Svmguide1 and Covertypes). The rest of the training data-set, other than the
 263 portion annotated by experts, is contaminated by noise according to the follow-
 264 ing three different noise models.

265 **Noise models:** The idea behind these different noise models is that they reflect
 266 real life scenarios of label noise.

267 *Boundary or margin noise (M):* Here we try to mimic the situation where
 268 annotators are confused about the correct label of a data point if it is close to
 269 the decision boundary. Hence to simulate margin noise we first train a support
 270 vector machine classifier (**SVM**) and get the margin γ . Then we flip the labels
 271 of 60% of the data points lying inside the margin. To increase the noise level
 272 we widen the margin by changing the parameter C of the **SVM**.

273 *Biased annotator noise (BA):* Here we are simulating the situation when a
 274 fraction of annotators are biased towards some unknown classifier and hence

275 they annotate all the data points according to that specific classifier. To simu-
 276 late this kind of noise we fix a random classifier and then change the labels of
 277 randomly chosen data points according to the outcome of that random classifier.

278 *Adversarial noise (A)*: We follow the method described in [8] to insert adver-
 279 sarial noise by flipping labels of those data points which have maximum impact
 280 on the classifier.

281 To illustrate the effect of noisy labels on various classifiers, we repeat all the
 282 experiments by varying the number of mislabeled examples from 5% to 45% of
 283 the training data-set and use a **SVM** and a Parzen window **PW** classifier as
 284 the final classifiers. We also study the case where label noise is present only in
 285 one class.

286 **Methods compared**: We compare the proposed **LND** with **KBDMS1**[10].
 287 In case of **KBDMS1**, we use two regularization terms $C_+ \|p_{\mathbf{I}_+}\|_1$ and $C_- \|p_{\mathbf{I}_-}\|_1$
 288 instead of $C \|p\|_1$ so that it can handle class conditional noise, where both C_+ and
 289 C_- are tuned in the range of $2^{\{-5:1:5\}}$. For learning a classifier from the outcome
 290 of the noisy label detection method **KBDMS1**, we flip the label of the i -th data
 291 point when $p_i > 0.5$. We also compare **LND** with some intuitive and simple
 292 approaches, **SubSVM** [12] and **SubPW** where we learn 50 local classifiers
 293 on the sub-sampled training data-set with a sample size of $\log_2 n$. We correct
 294 the label of every training data point except those annotated by the experts,
 295 using majority vote from predictions of these 50 sub learners. We keep the
 296 parameter values equal for all the 50 sub learners. We study the impact of noise
 297 correction methods against a nominal **SVM**, **PW** and robust counter parts of
 298 **SVM** such as **SVM** with class conditional cost (**CSVM**) [9] and Robust SVM
 299 (**RSVM**) [8]. Moreover, we compare the obtained classification accuracy with
 300 the results of classifiers trained with correctly annotated labels (**True-SVM** or
 301 **True-PW**) and trained with data points annotated by the expert annotators
 302 (**Expert-SVM** or **Expert-PW**).

303 The parameters ρ_+/ρ_- for all kinds of **LND** are chosen by cross-validation
 304 from $\{0 : 0.05 : 0.5\}$ using knowledge of expert labels. For cross-validation, the
 305 training data-set is divided in such a way that the data points labeled by expert

306 annotators are equally distributed over all the partitions and the validation error
 307 is calculated by considering the mis-classification error **only** for the data points
 308 labeled by experts. All other parameters like $C, C_+, C_- \in 2^{\{-10:2:10\}}$ for all
 309 **SVM** classifiers and $\mu \in \{0 : 0.1 : 0.5\}$ of **RSVM** are chosen also by 5-fold
 310 cross-validation in a similar way. The bandwidth h of the Gaussian kernel is
 311 tuned independently for all the algorithms varying in the range of $2^{\{-5:1:5\}}$. We
 312 are not able to compare the proposed method with **ROD** [6] as it is not scalable
 313 beyond only hundred data points.

314 We compare the performances of all the proposed algorithms **LND**_{slp}, **LND**_r
 315 for $r = 1, 2$ and **LND**_{slp1} to verify importance of the proposed Algorithm 4.
 316 **LND**_{slp1} uses **LND**_{slp} with output from **LND**₁ as the starting point. The
 317 results for **LND**_{slp} correspond to the best local maximum obtained from 100
 318 random initializations.

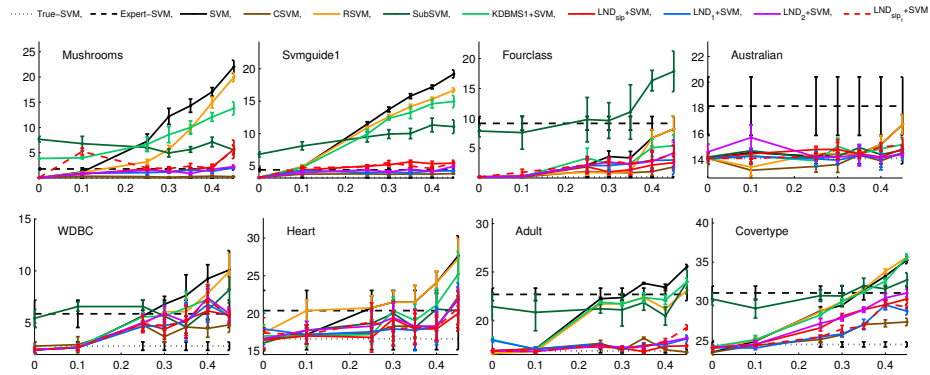
319 **Computational setting:** All experiments are done using Intel(R)-Xeon(R)
 320 (2.67GHz) processor with 36 GB RAM. For **SVM** and **SubSVM** we use Lib-
 321 SVM [19] with its Matlab interface. While all other algorithms are implemented
 322 with Matlab(R2013a). ².

323 6.2. Results

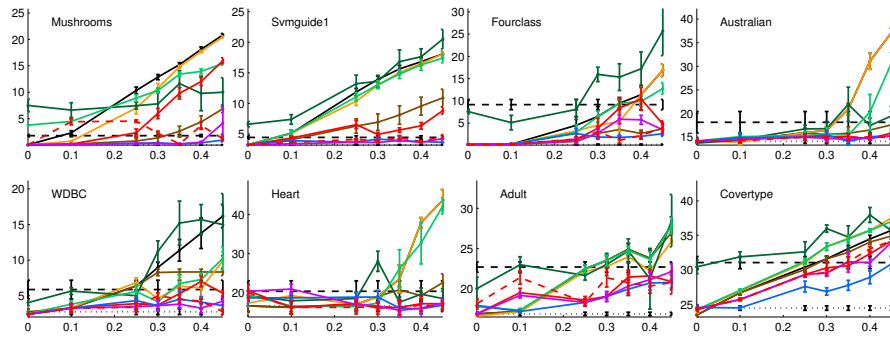
324 **Classification accuracy:** In this section we study how the correction of label
 325 noise by different algorithms influence the test error of the final classifiers (**SVM**
 326 and **PW**). Experimental results in Figure 1,2,3,4³ indicate the fact that for both
 327 **SVM** and **PW** classifiers, correcting annotations using our model helps to get
 328 better test errors. By increasing the number of noisy labels in the training
 329 data-set test errors of **SVM** and **PW** classifiers increase heavily while after
 330 pre-processing with the proposed method (**LND**₁ , **LND**₂ and **LND**_{slp1}) the
 331 test error increases with lower rate and in most of the cases remain very close to

²Both the code for LNDs and the used data-sets with label noise are available at
<http://www.ml.uni-saarland.de/code/LabelNoiseCorrection/LabelNoiseCorrection.htm>

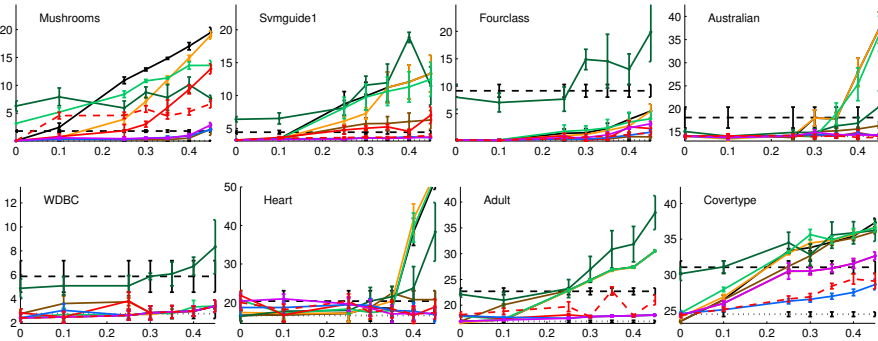
³ Each plot shows comparison of performance of various methods on each data-set (written
 on the plot).



(a) Boundary Noise

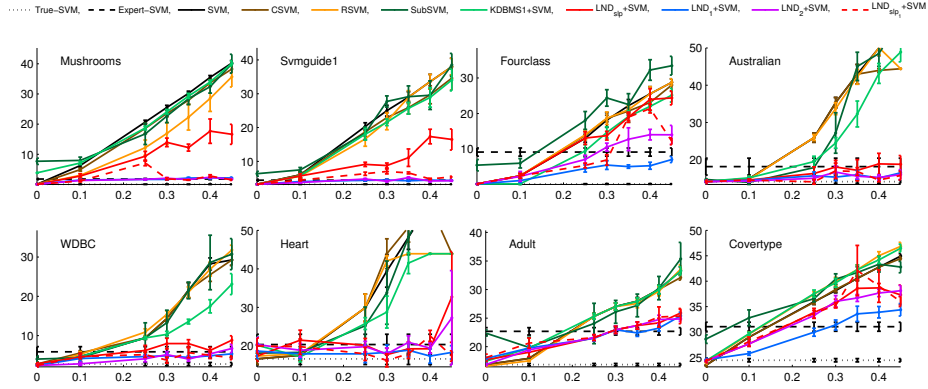


(b) Biased Annotator Noise

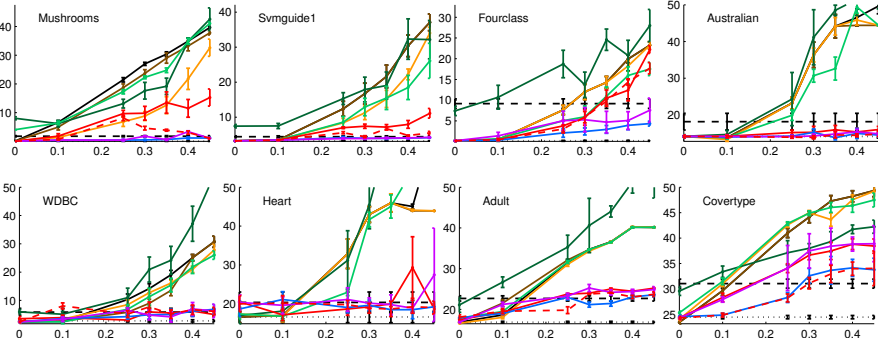


(c) Adversarial Noise

Figure 1: Comparison of different label noise correction methods and robust SVMs in terms of classification error of the **SVM** classifiers trained with corrected labels. In each plot the y-axis shows the test error (in %) and the x-axis shows the fraction of noisy labels present in the majority class of the training data-set.



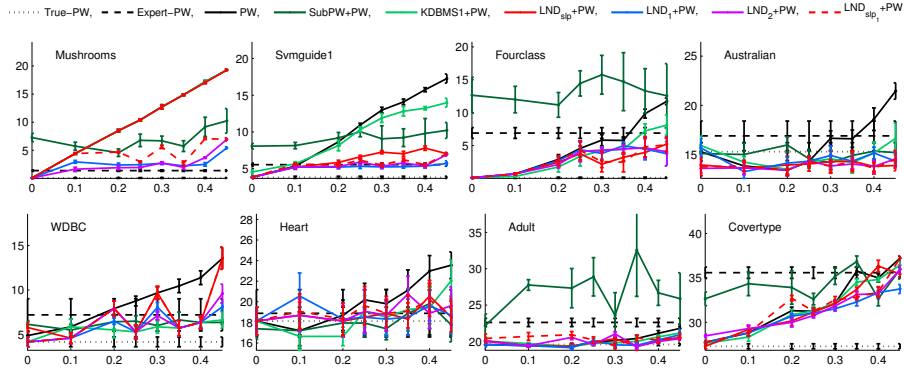
(a) Biased Annotator Noise



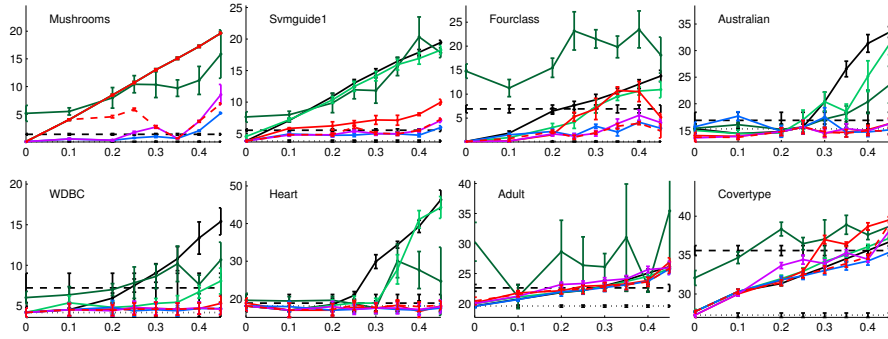
(b) Adversarial Noise

Figure 2: Comparison of different label noise correction methods and robust SVMs in terms of classification error of the **SVM** classifiers trained with corrected labels. In each plot the y-axis shows the test error (in %) and the x-axis shows the fraction of noisy labels present in both classes of the training data-set.

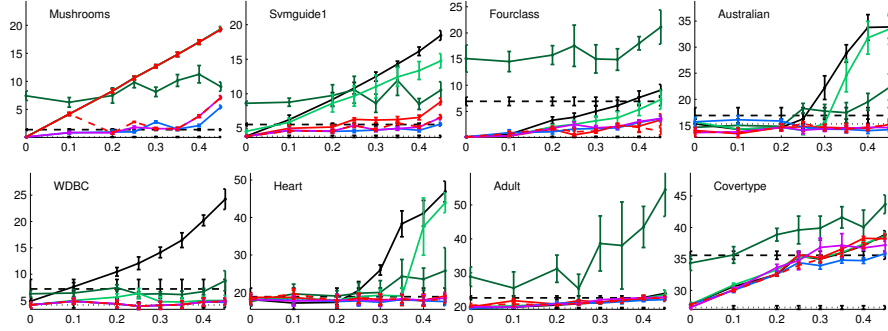
332 the test error of the classifier trained with true labels. **KDBMS1** performs well
 333 when the number of noisy data points is small but by increasing the number
 334 of noisy labels it deteriorates. The reason is that with a high regularization
 335 parameter it can only detect a few data points as noisy on the other hand by
 336 decreasing the value of the regularization parameters after a certain value, the
 337 effect of regularization become negligible and it starts detecting noisy data-
 338 points only from one class (as discussed in Section 2). The performance of
 339 **SubSVM** and **SubPW** are not consistent and vary for different data-sets and



(a) Boundary Noise

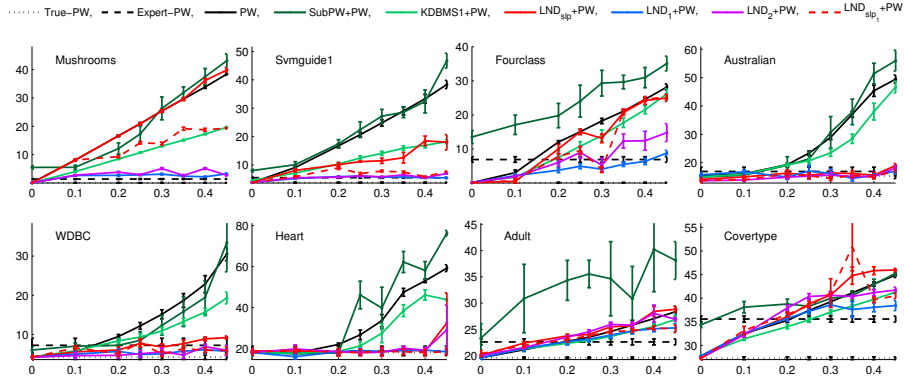


(b) Biased Annotator Noise

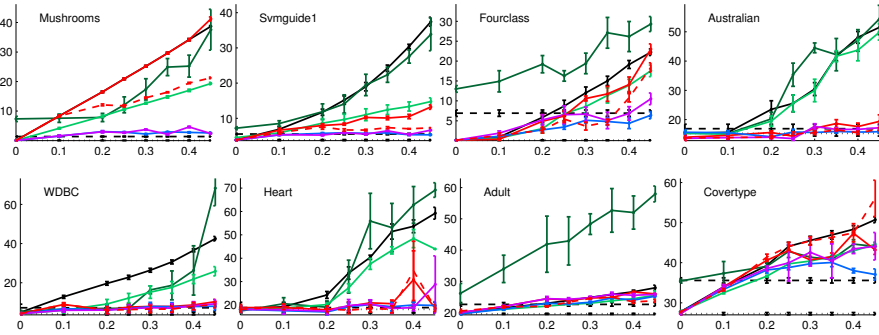


(c) Adversarial Noise

Figure 3: Comparison of different label noise correction methods in terms of classification error of **PW** classifiers trained with corrected labels. In a plot the y-axis shows the test error (in %) and the x-axis shows the fraction of noisy labels present in the majority class of the training data-set.



(a) Biased Annotator Noise



(b) Adversarial Noise

Figure 4: Comparison of different label noise correction methods in terms of classification error of **PW** classifiers trained with corrected labels. In a plot the y-axis shows the test error (in %) and the x-axis shows the fraction of noisy labels present in both classes of the training data-set.

340 also have large variance.

341 Figures 1,2 show that in general, for the biased annotator noise and adversarial noise, **SVM** classifiers trained with the labels corrected by **LND** also
 342 outperforms all the other robust counter parts (**CSVM** and **RSVM**). Only
 343 in the case of boundary noise, **CSVM** performs best but our method is close.
 344 But for other kinds of noise the performance of **CSVM** deteriorates. When the
 345 noise is symmetric for both the classes, the performance of **CSVM** is not better
 346 than the performance of **SVM** trained with noisy data. As the margin noise in
 347 both the classes does not affect **SVM** much, we have not studied that kind of
 348

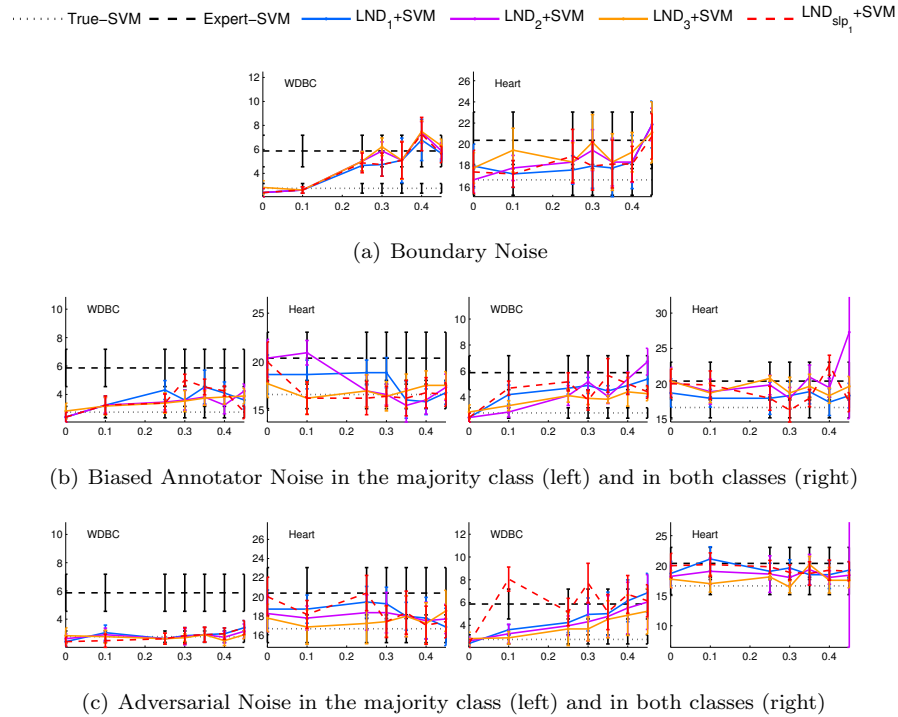


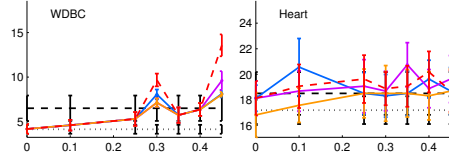
Figure 5: Effect of approximation quality for **LND** in terms of classification error of the **SVM** classifiers trained with labels corrected by **LND**₁, **LND**₂ and **LND**₃. In the plots the y-axis shows the test error (in %) and the x-axis shows the fraction of noisy labels present in the majority class (left) and in both classes (right) of the training data-set.

349 noisy data.

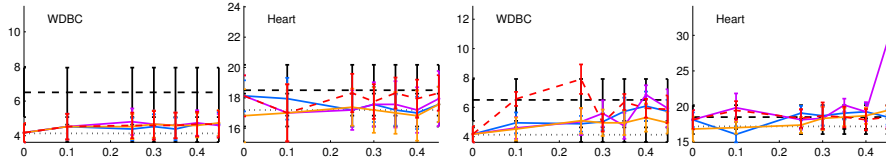
350 We also study performance of rank-3 approximation of **LND**(**LND**₃) for
 351 the smaller data-sets: WDBC and Heart. Figure 5, 6 show that for all kinds of
 352 noise in the case of WDBC and for biased annotator and boundary noise in the
 353 case of Heart the performance of rank-1 approximation of **LND**(**LND**₁) is also
 354 very close (sometimes better) to the performance of higher order (rank-2 and
 355 rank-3) approximations. Figure 1, 2, 3, 4 also show that one can get a good
 356 amount of improvement in terms of classification accuracy over **SVM** and **PW**
 357 classifiers by pre-processing the noisy data-sets using (**LND**₁) with a running
 358 time of $O(n^2)$ which is also feasible for large data-sets.

359 We compare classification accuracy when availability of the expert labels

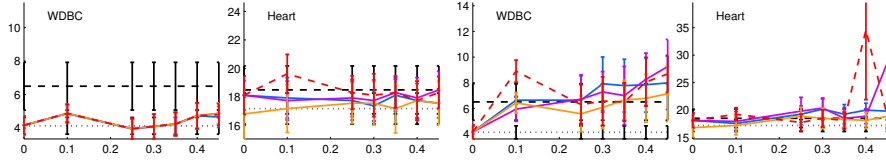
..... True-PW - - - Expert-PW — LND₁+PW — LND₂+PW — LND₃+PW - - - LND_{slp₁}+PW



(a) Boundary Noise



(b) Biased Annotator Noise in the majority class (left) and in both classes (right)



(c) Adversarial Noise in the majority class (left) and in both classes (right)

Figure 6: Effect of approximation quality for **LND** in terms of classification error of the **PW** classifiers trained with labels corrected by **LND**₁, **LND**₂ and **LND**₃. In the plots the y-axis shows the test error (in %) and the x-axis shows the fraction of noisy labels present in the majority class (left) and in both classes (right) of the training data-set.

360 varies. Figure 7, 8 show how the performance of **Expert-SVM (Expert-PW)**,
 361 nominal **SVM (PW)** and **LND**₁ differ when the percentage of expert labels in
 362 the training data-set varies among 1%, 2%, 5%, 10% and 20%. Important to note
 363 that, even with 1% of expert labels, the proposed **LND** improves classification
 364 accuracy sometimes more than 10% of that achieved by nominal **SVM** and
 365 **PW**. **LND**₁ also beats **Expert-SVM** and **Expert-PW** with high margin for
 366 small noise level and the difference is more when small number of expert labels
 367 are available. Please note that, we are using expert label information to tune
 368 ρ_- and ρ_+ in cross-validation. Hence, we are not able to tune these parameters
 369 when there are no expert label available. That is why, we do not include the

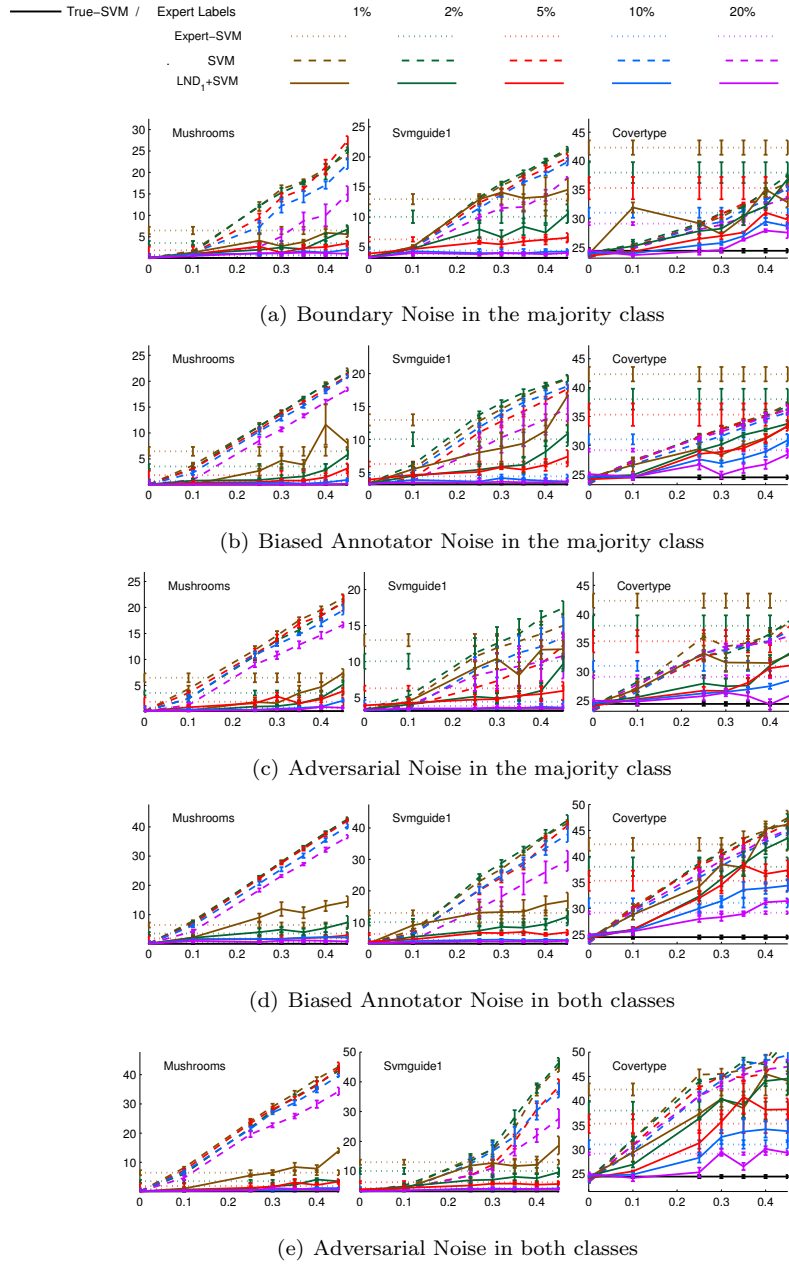


Figure 7: Plots show the effect of different fraction of expert labels for **LND** in terms of classification error of **SVM** trained with corrected label. In the plots the y-axis shows the test error (in %) and the x-axis shows the fraction of noisy labels present in the training data-set.

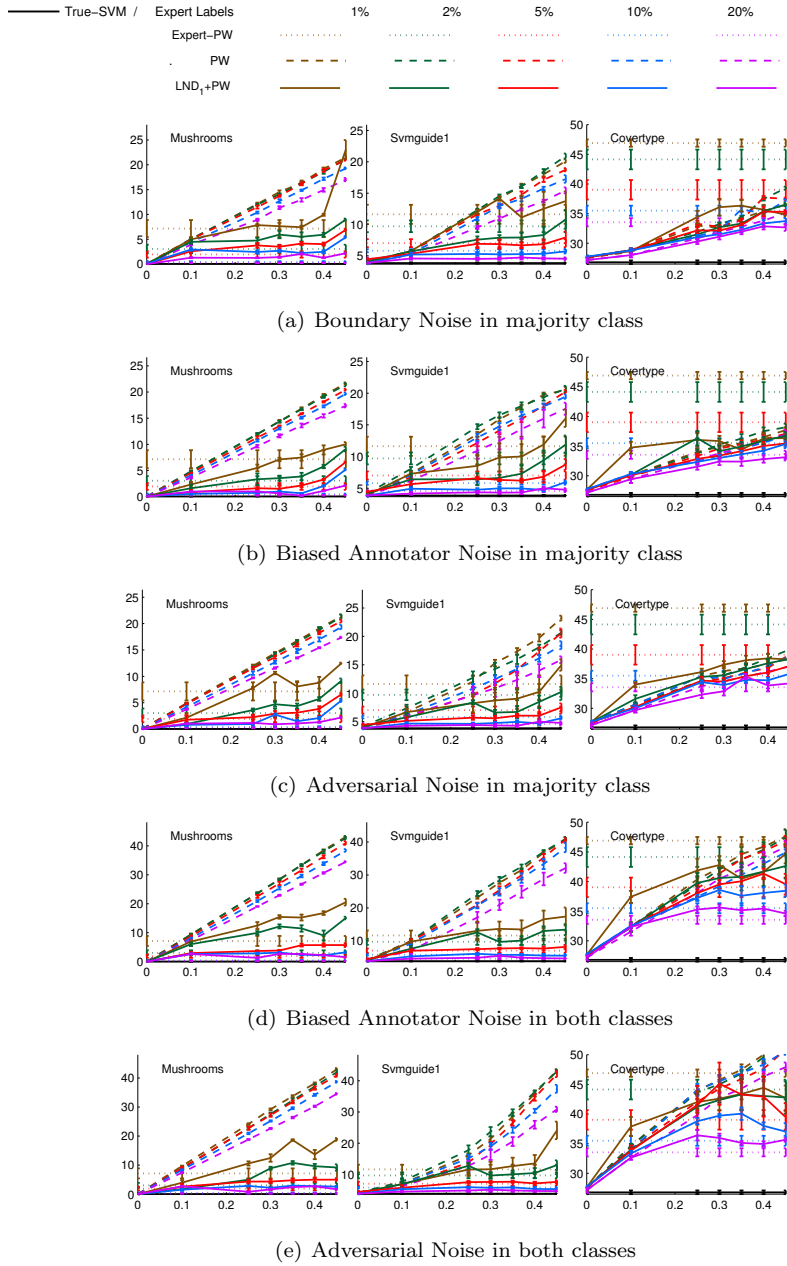
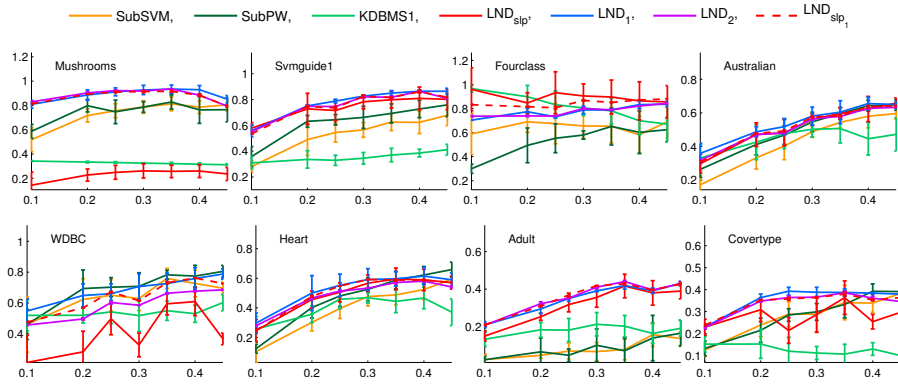
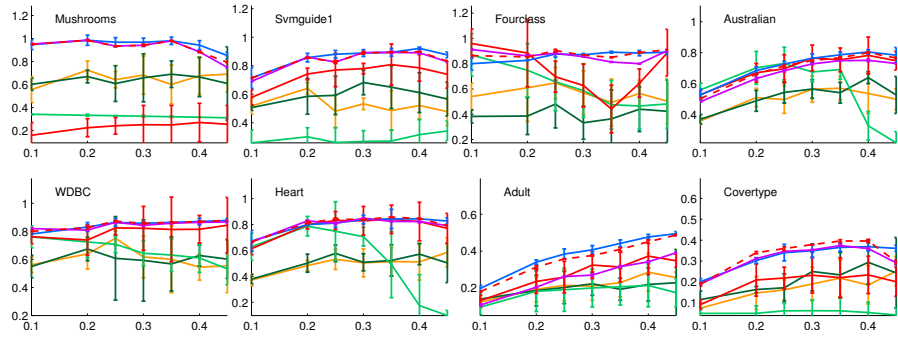


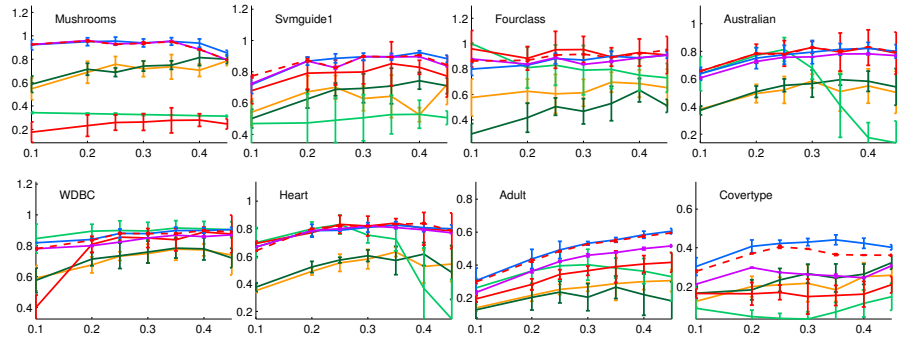
Figure 8: Plots show the effect of different fraction of expert labels for **LND** in terms of classification error of **PW** trained with corrected label. In the plots the y-axis shows the test label error (in %) and the x-axis shows the fraction of noisy labels present in the training data-set.



(a) Boundary Noise

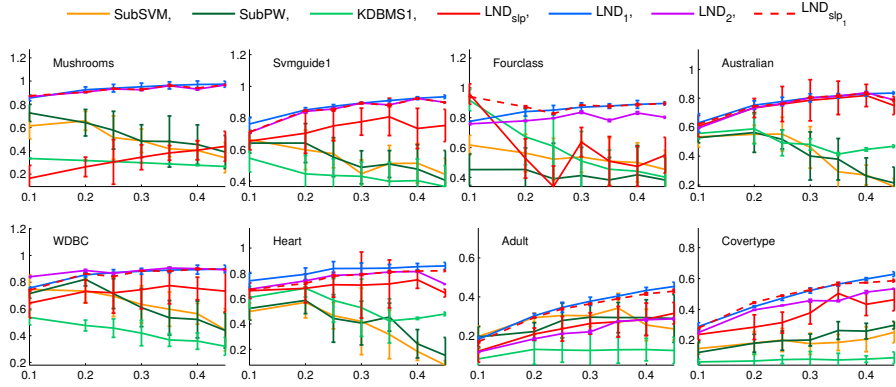


(b) Biased Annotator Noise

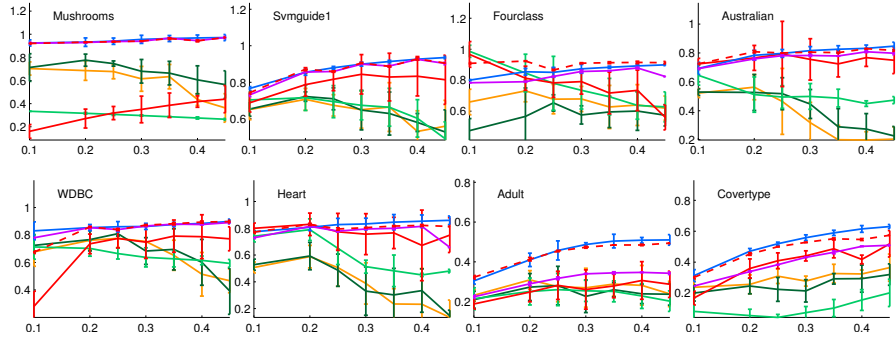


(c) Adversarial Noise

Figure 9: The label noise detecting ability of various methods in term of the F1-score. In the plot the y-axis shows the F1-score of noise detection and the x-axis shows the fraction of noisy labels present in the majority class of the training data-set.



(a) Biased Annotator Noise



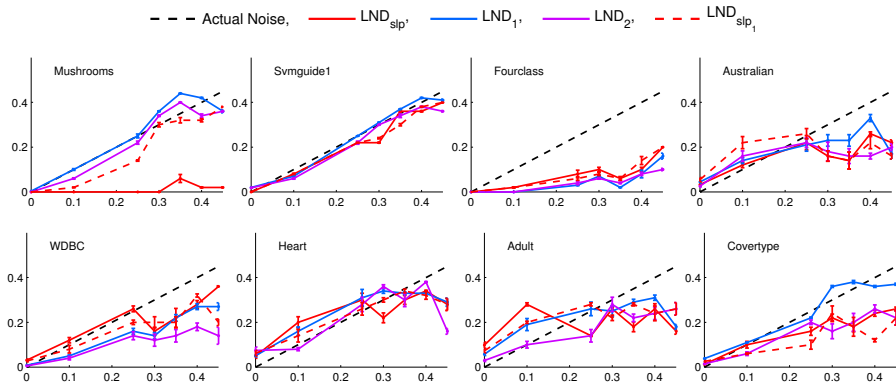
(b) Adversarial Noise

Figure 10: The label noise detecting ability of various methods in term of the F1-score. In the plot the y-axis shows the F1-score of noise detection and the x-axis shows the fraction of noisy labels present in both classes of the training data-set.

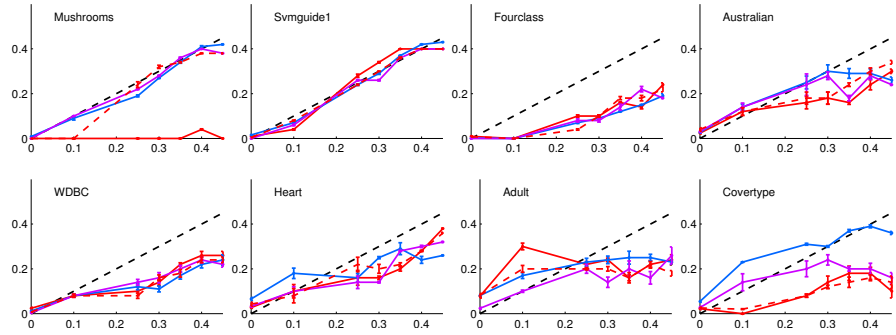
370 plot for the worst case scenario, i.e., the plots for no expert labels.

371

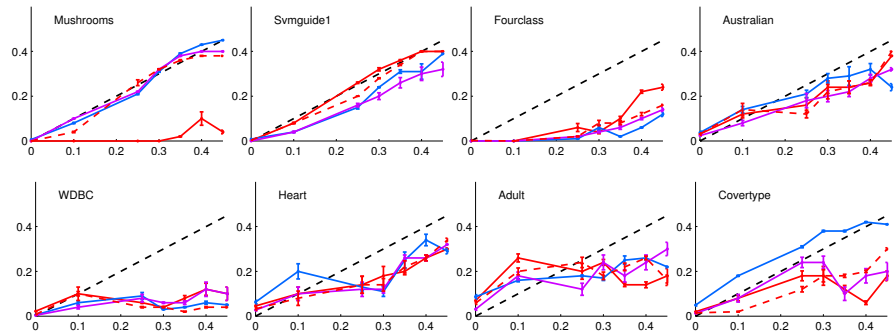
372 **Noisy label detection ability:** We compare in Figures 9, 10 the average F1
 373 scores of detecting noisy labels achieved by different methods with proper pa-
 374 rameter tuning through cross-validation. For this set of experiments we show
 375 the results only up-to rank-2 approximation of **LND** as for WDBC and Heart,
 376 the F1 score of rank-3 approximation is very close to that of the rank-2 approx-
 377 imation. In case of adversarial noise and biased annotator noise our proposed
 378 method is able to detect noisy labels with F1 scores higher than 0.8 except



(a) Boundary Noise

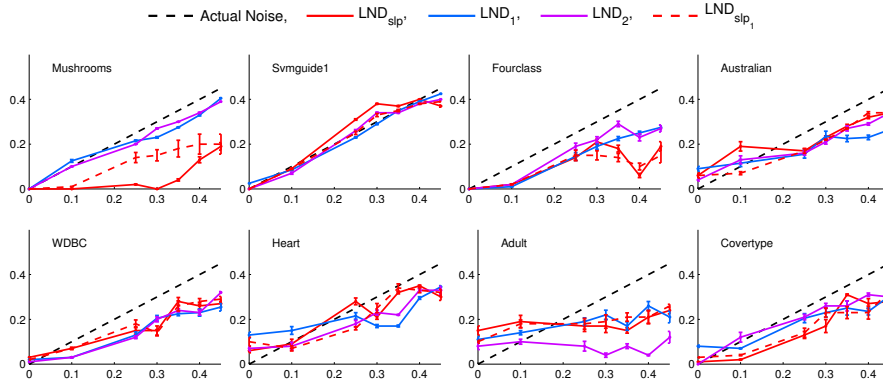


(b) Biased Annotator Noise

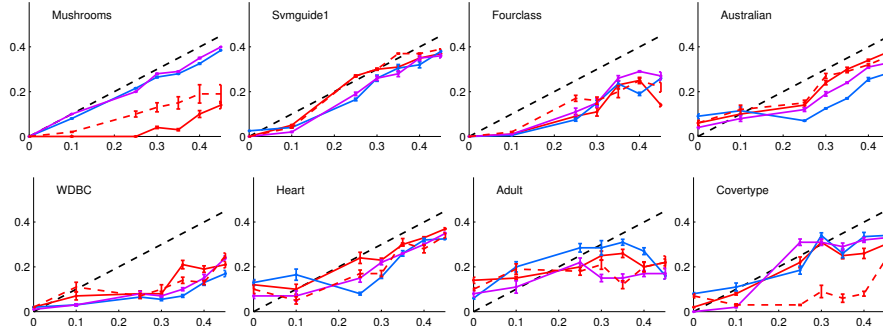


(c) Adversarial Noise

Figure 11: Accurate estimation of noise ratio. In each plot the y-axis shows the average estimation of noise ratio (ρ_+ or ρ_-) and the x-axis shows the fraction of noisy labels present in the majority class of the training data-set.



(a) Biased Annotator Noise



(b) Adversarial Noise

Figure 12: Accurate estimation of noise ratio. In each plot the y-axis shows the average estimation of noise ratio (ρ_+ or ρ_-) and the x-axis shows the fraction of noisy labels present in both classes of the training data-set.

379 for the Adult and the Coverttype data-set. For the Mushrooms the F1 score
 380 is even closer to 1 for all kind of noise. Whereas the F1 scores for the other
 381 algorithms are below 0.6 in most of the cases. For small number of noisy labels,
 382 **KBDMS1** sometimes performs better than **LND** but when more noisy data
 383 points are present, **LND₁**, **LND₂** and **LND_{slp₁}** outperform **KBDMS1** with a
 384 margin of more than 0.2. Performance of **LND_{slp}** is not consistent and varies
 385 a lot. This shows the superiority of Algorithm 4.

386

387 **Accurate estimation of the noise ratio:** We verify how close the estimated

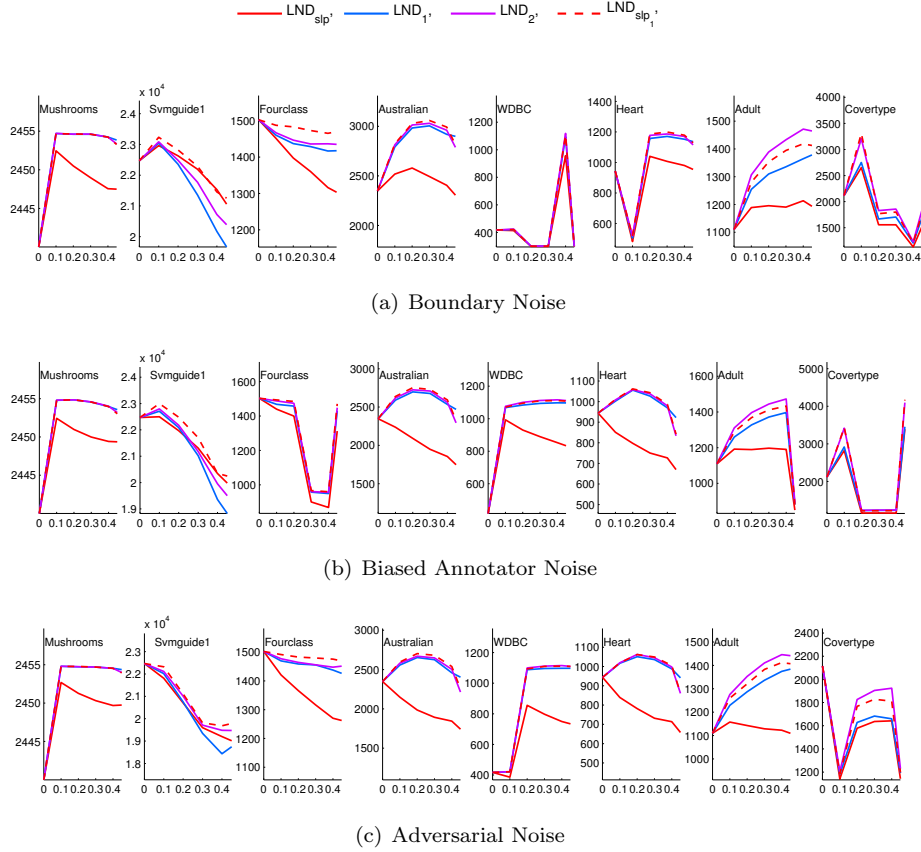
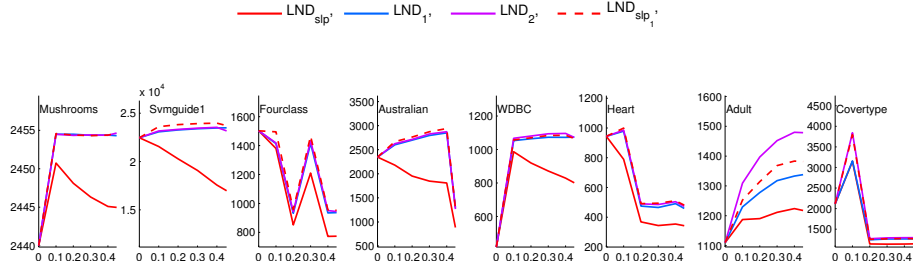


Figure 13: Plots show objective values achieved by various algorithms and approximations of LND , i.e., LND_{slp} , LND_1 , LND_2 and LND_{slp_1} . In each plot the y-axis shows the objective values of LND (Equation (6)) and the x-axis shows the fraction of noisy labels present in the majority class of the training data-set.

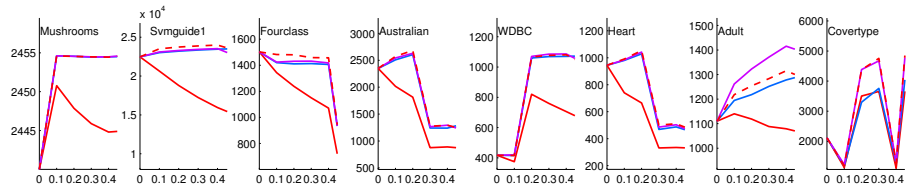
388 noise ratios (ρ_+ and ρ_-) are to the actual percentage of noisy labels. Figures
389 11, 12 show that for proposed methods (LND_1 , LND_2 and LND_{slp_1}) the noise
390 ratios are estimated accurately by cross-validation. For all data-sets except
391 Adult and Fourclass and all types of noise models the estimated values are very
392 close to the actual proportion of the noise. When the noise is larger (≥ 0.3) the
393 estimated values are sometimes smaller than the true values.

394

395 **Solution quality:** We compare the objective values achieved by various al-



(a) Biased Annotator Noise



(b) Adversarial Noise

Figure 14: Plots show objective values achieved by various algorithms and approximations of \mathbf{LND} , i.e., \mathbf{LND}_{slp} , \mathbf{LND}_1 , \mathbf{LND}_2 and \mathbf{LND}_{slp_1} . In each plots the y-axis shows the objective values of \mathbf{LND} (Equation (6)) and the x-axis shows the fraction of noisy labels present in both classes of the training data-set.

396 algorithms for solving (6). In these experiments, we use ρ_+ and ρ_- to be same
 397 as that used during the phase of insertion of noise. The experimental results
 398 (Figure 13,14 show that in almost all the cases the objective values achieved by
 399 \mathbf{LND}_1 and \mathbf{LND}_2 based on Algorithm 4 are much higher than that of \mathbf{LND}_{slp} .
 400 Whereas the objective values achieved by \mathbf{LND}_{slp_1} is higher than that of \mathbf{LND}_1 ,
 401 and also for few data-sets (Svmguide1, Fourclass and Australian) objective val-
 402 ues achieved by \mathbf{LND}_1 are even higher than \mathbf{LND}_2 . Performance of \mathbf{LND}_{slp} is
 403 not consistent as most of the time the objective values achieved by it are very
 404 far from the optimal one.

405 We plot (Figure 15,16) $\frac{\text{Bound of } OPT^*}{OPT^*}$, where the bound is calculated by
 406 (18). For the Mushrooms the bound is very tight with $\frac{\text{Bound of } OPT^*}{OPT^*} \leq 1.03$. In
 407 general, the approximation guarantees are very good and quite stable across all
 408 the data-sets.

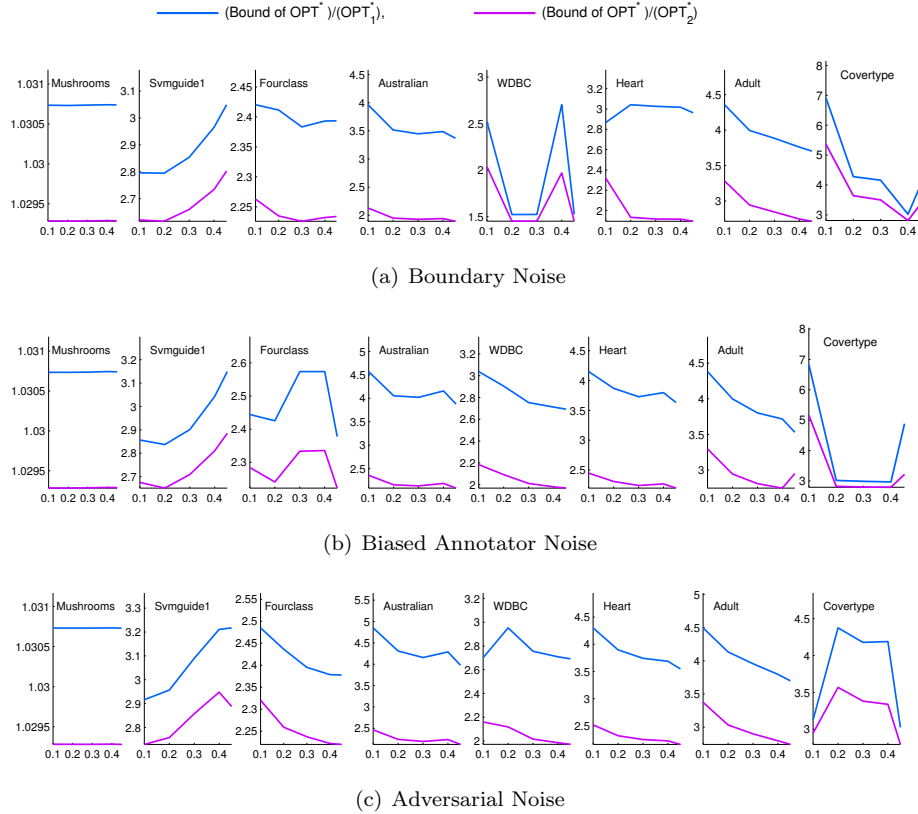
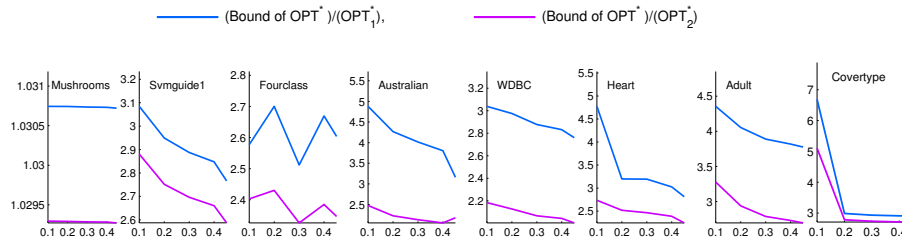


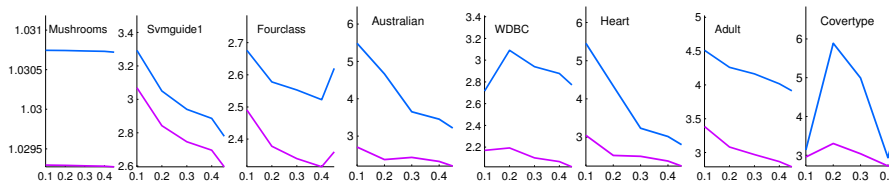
Figure 15: Plots show $\frac{\text{Bound of } OPT^*}{OPT_r^*}$ for $r = 1$ and 2. In each plot the y-axis shows $\frac{\text{Bound of } OPT^*}{OPT_r^*}$ and the x-axis shows the fraction of noisy labels present in the majority class of the training data-set.

409

410 **Comparison of time complexity for noise label detection:** Here we
 411 compare time complexities for various methods and algorithms. Each number
 412 in Table 2 reports average time required for each algorithm to detect noisy
 413 labels for a fixed set of parameter h , C , C_- , C_+ , ρ_- and ρ_+ . Time required
 414 for LND_1 is less than that for KBDMS1 , LND_{slp} and for smaller data sets
 415 it is comparable to SubPW and SubSVM . For moderately sized data-sets
 416 training time of LND_1 is more than that of SubPW and SubSVM , where
 417 timing for SubPW and SubSVM depends on the number of local classifiers



(a) Biased Annotator Noise



(b) Adversarial Noise

Figure 16: Plots show $\frac{\text{Bound of } OPT^*}{OPT_r^*}$ for $r = 1$ and 2. In the plot the y-axis shows $\frac{\text{Bound of } OPT^*}{OPT_r^*}$ and the x-axis shows the fraction of noisy labels present in both classes of the training data-set.

418 learned. Here we learn **SubSVM** and **SubPW** with only 50 local classifiers.
 419 On the other hand, classification accuracy achieved by **SubSVM** and **SubPW**
 420 is much worse compared to that of **LND**₁. As expected, the running time for
 421 **LND**₂ is high but for larger data-sets it is better than the running time of
 422 **KBDMS1**. The time required for rank-3 approximation of **LND** for Heart
 423 and WDBC are 954.67 ± 106.2 and 591.23 ± 41.62 respectively.

424 In Figure 17 we study how the time complexities of noise detection algo-
 425 rithms vary with the increase of training data-sets. We do this using the Cover-
 426 type data-set and we are not able to do this experiment beyond 3×10^4 for the
 427 proposed **LND** algorithms in our system.

428 7. Conclusion

429 In this paper we propose **LND** as a novel method for correcting labels of
 430 mislabeled data points. Although, the optimization problem underlying **LND**
 431 is NP-hard, we extend the Spannogram algorithm to obtain a solution with

Table 2: The time required for detecting noisy labels by various methods (in seconds)

Dataset	SubSVM	SubPW	KBDMS1	LND _{slp}	LND ₁	LND ₂	LND _{slp1}
Heart 150x13	0.035 ±0.019	0.007 ±0.002	0.213 ±0.056	1.060 ±0.245	0.012 ±0.003	1.369 ±0.571	0.1473 ±0.051
WDBC 290x30	0.034 ±0.012	0.006 ±0.004	0.291 ±0.060	2.249 ±0.752	0.013 ±0.002	2.826 ±1.359	0.0752 ±0.015
Australian 310x14	0.047 ±0.015	0.009 ±0.002	0.358 ±0.155	2.510 ±0.498	0.013 ±0.002	2.081 ±0.399	0.0850 ±0.015
Fourclass 430x2	0.054 ± 0.020	0.011 ±0.003	0.621 ±0.187	3.8754 ±0.834	0.019 ±0.002	3.412 ±0.282	0.155 ±0.072
Adult 1605x123	0.566 ±0.394	0.155 ±0.129	2.739 ±0.681	6.198 ±1.99	0.151 ±0.016	7.076 ±1.009	0.2538 ±0.543
Svmguide1 3089x4	0.329 ±0.263	0.297 ±0.133	16.893 ±1.362	28.994 ±5.405	0.323 ±0.021	10.486 ±1.380	0.850 ±0.385
Mushroom 4062x112	0.437 ±0.279	0.350 ±0.085	27.062 ±2.999	99.849 ±17.813	0.824 ±0.030	18.3610 ±1.885	2.8788 ±1.240
Coverttype 5810x54	0.502 ±0.277	0.286 ± 0.053	41.386 ± 3.845	156.517 ±26.02	1.394 ±0.0074	35.266 ±2.111	3.702 ±0.708

432 a provable approximation guarantee. Experimental results using a variety of
 433 data-sets and different noise models demonstrate that the proposed approach
 434 **LND** outperforms all existing methods for the detection of noisy labels by large
 435 margin. **SVM** trained with data pre-processed by **LND** also outperforms all
 436 other existing version of SVMs which are supposed to be robust to the label
 437 noise.

438 8. Acknowledgment

439 Sahely Bhadra is funded by the Indo-German Max Planck Center for Com-
 440 puter Science (IMPECS).

441 References

- 442 [1] J. Howe, Crowdsourcing: How the power of the crowd is driving the future
 443 of business, Random House, 2008.
- 444 [2] C. E. Brodley, M. A. Friedl, Identifying mislabeled training data, Journal
 445 of Artificial Intelligence Research 11 (1999) 131–167.

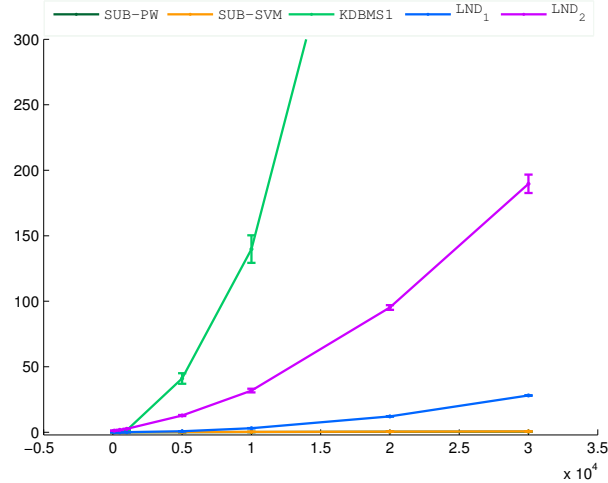


Figure 17: This plot compares the training time of various noise detection algorithms with various size of training data-sets. In the plot the y-axis shows the time required for detecting noisy labels for a fixed set of parameters (in *sec*) and the x-axis shows the number of training data points

- 446 [3] V. Guruswami, P. Raghavendra, Hardness of learning halfspaces with noise,
 447 in: FOCS, 2006, pp. 543–552.
- 448 [4] C. Caramanis, S. Mannor, Learning in the limit with adversarial distur-
 449 bances, in: COLT, 2008, pp. 467–478.
- 450 [5] R. Collobert, F. Sinz, J. Weston, L. Bottou, Trading convexity for scala-
 451 bility, in: ICML, 2006, pp. 201–208.
- 452 [6] L. Xu, K. Crammer, D. Schuurmans, Robust support vector machine train-
 453 ing via convex outlier ablation, in: AAAI, 2006, pp. 536–542.
- 454 [7] G. Stempfel, L. Ralaivola, Learning svms from sloppily labeled data, in:
 455 ICANN (1), 2009, pp. 884–893.
- 456 [8] B. Biggio, B. Nelson, P. Laskov, Support vector machines under adversarial
 457 label noise, in: ACML, Vol. 20, 2011, pp. 97–112.

- 458 [9] N. Natarajan, I. Dhillon, P. Ravikumar, A. Tewari, Learning with noisy
459 labels, in: NIPS, 2013, pp. 1196–1204.
- 460 [10] H. Valizadegan, P. N. Tan, Kernel based detection of mislabeled training
461 examples., in: SDM, 2007.
- 462 [11] J. Cao, S. Kwong, R. Wang, A noise-detection based adaboost algorithm
463 for mislabeled data, Pattern Recognition 45 (12) (2012) 4451–4465.
- 464 [12] S. Laxman, S. Mittal, R. Venkatesan, Error correction in learning using
465 svms, CoRR abs/1301.2012.
- 466 [13] D. S. Papailiopoulos, A. G. Dimakis, S. Korokythakis, Sparse pca through
467 low-rank approximations, in: ICML, Vol. 3, 2013, pp. 747–755.
- 468 [14] D. S. Papailiopoulos, I. Mitliagkas, A. Dimakis, C. Caramanis, Finding
469 dense subgraphs through low-rank approximations, preprint.
470 URL https://webpace.utexas.edu/dp26726/papers/DkS_long.pdf
- 471 [15] X. Zhu, X. Wu, Q. Chen, Eliminating class noise in large datasets, in:
472 ICML, 2003, pp. 920–927.
- 473 [16] R. Duda, P. Hart, D. Stork, Pattern Classification, Wiley, 2001.
- 474 [17] D. Wied, R. Weißbach, Consistency of the kernel density estimator: a
475 survey, Statistical Papers 53 (1) (2012) 1–21.
- 476 [18] B. Moghaddam, Y. Weiss, S. Avidan, Spectral bounds for sparse pca: Ex-
477 act and greedy algorithms, in: Advances in neural information processing
478 systems, 2005, pp. 915–922.
- 479 [19] C. Chang, C. Lin, LIBSVM: A library for support vector machines,
480 ACM Transactions on Intelligent Systems and Technology 2 (2011)
481 27:1–27:27, dataset available at [http://www.csie.ntu.edu.tw/~cjlin/](http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets)
482 [libsvmtools/datasets](http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets).