**Dense prediction of label noise for learning building extraction from aerial drone imagery**

Nahian Ahmed[a]*, Rashedur M. Rahman[a], Mohammed Sarfaraz Gani Adnan[b], Bayes Ahmed[c]

[a]*Department of Electrical and Computer Engineering, School of Engineering and Physical Sciences, North South University, Bashundhara, Dhaka 1229, Bangladesh;*

[b]*Department of Urban and Regional Planning, Chittagong University of Engineering and Technology (CUET), Chittagong 4349, Bangladesh;*

[c]*Institute for Risk and Disaster Reduction (IRDR), University College London (UCL), Gower Street, London WC1E 6BT, UK;*

Email: nahian.ahmed@northsouth.edu

## 14  **Dense prediction of label noise for learning building extraction from**
## 15  **aerial drone imagery**

16  Label noise is a commonly encountered problem in learning building extraction
17  tasks; its presence can reduce performance and increase learning complexity.
18  This is especially true for cases where high resolution aerial drone imagery is
19  used, as the labels may not perfectly correspond/align with the actual objects in
20  the imagery. In general machine learning and computer vision context, labels
21  refer to the associated class of data, and in remote sensing-based building
22  extraction refer to pixel-level classes. Dense label noise in building extraction
23  tasks has rarely been formalized and assessed. We formulate a taxonomy of label
24  noise models for building extraction tasks, which incorporates both pixel-wise
25  and dense models. While learning dense prediction under label noise, the
26  differences between the ground truth clean label and observed noisy label can be
27  encoded by error matrices indicating locations and type of noisy pixel-level
28  labels. In this work, we explicitly learn to approximate error matrices for
29  improving building extraction performance; essentially, learning dense prediction
30  of label noise as a subtask of a larger building extraction task. We propose two
31  new model frameworks for learning building extraction under dense real-world
32  label noise, and consequently two new network architectures, which approximate
33  the error matrices as intermediate predictions. The first model learns the general
34  error matrix as an intermediate step and the second model learns the false positive
35  and false negative error matrices independently, as intermediate steps.
36  Approximating intermediate error matrices can generate label noise saliency
37  maps, for identifying labels having higher chances of being mis-labeled. We have
38  used ultra-high-resolution aerial images, noisy observed labels from
39  OpenStreetMap, and clean labels obtained after careful annotation by the authors.
40  When compared to the baseline model trained and tested using clean labels, our
41  intermediate false positive-false negative error matrix model provides
42  Intersection-Over-Union gain of 2.74% and F1-score gain of 1.75% on the
43  independent test set. Furthermore, our proposed models provide much higher
44  recall than currently used deep learning models for building extraction, while
45  providing comparable precision. We show that intermediate false positive-false
46  negative error matrix approximation can improve performance under label noise.

**Introduction**
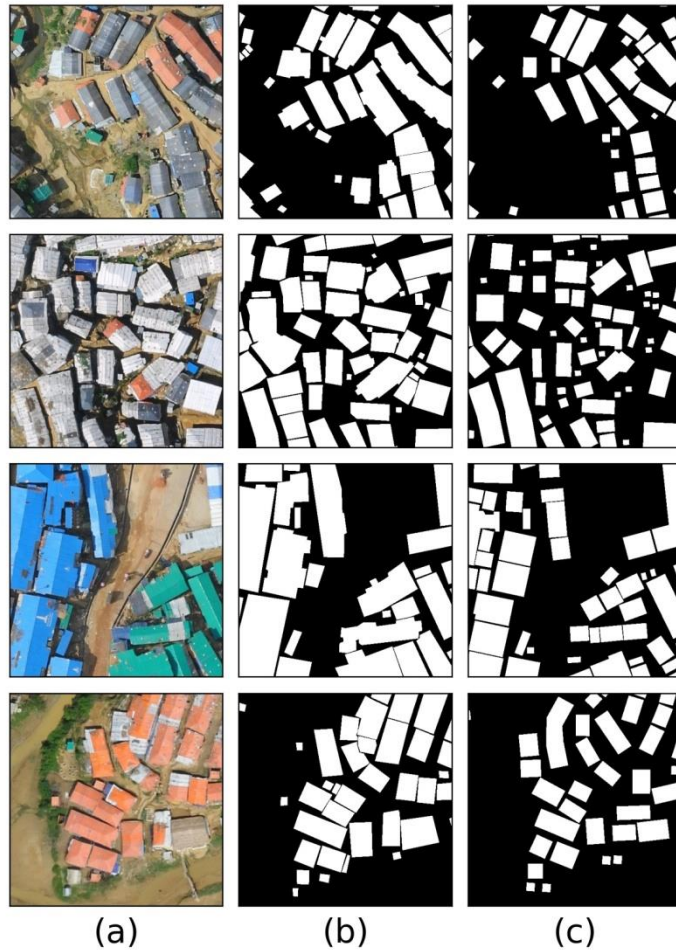
Building extraction involves learning mappings between remotely sensed aerial or satellite images and building labels from freely available vector data. The most commonly used source of labels, OpenStreetMap, though accurate to a large degree, contain various types of label noise (Mnih and Hinton, 2012; Ahmed et al., 2020; Zhang et al., 2020). Pixel-level predictions of building/non-building labels are performed, which is a binary dense prediction task. Label noise occurs when the observed label does not agree with the true label (Frénay and Verleysen, 2013; Frénay and Kabán, 2014) (Fig. 1). Presence of label noise in training data can reduce performance, while noise in testing data can lead to underestimation of model performance (Ahmed et al., 2020). However, most of the existing studies on deep learning-based building extraction do not acknowledge the presence of label noise. In general, complexity of the learning task is also increased under label noise (Garcia et al., 2015; Pelletier et al., 2017). Research on robust method of building extraction considering label noise requires formalization of the sources, processes and effects of noise on large scale freely available labels. Currently, the types of dense label noise processes have not been formalized in a comprehensively and inclusively in research. When building polygons are rasterized, the buildings are represented as superpixels in the prepared dense binary labels. Individual building polygon i.e. superpixel based errors are commonly considered as sources of noisy labels.

72      Coming from traditional remote sensing terminology, the most common are

73  registration errors, where building polygons are present but not aligned, annotated or

74  registered properly, and omission errors where buildings are left unlabeled (Mnih and

75  Hinton, 2012; Ahmed et al., 2020; Zhang et al., 2020). However, alternative

76  nomenclature has been proposed as well. Pixel-based nomenclature can be used to

77  express label noise processes in multiple scales, and therefore provides a more

78  generalized viewpoint. Even superpixel-based label noise processes are modeled using a

79  composite of pixel-based processes (Mnih and Hinton, 2012; Zhang et al., 2020). This

80  approach assumes that each pixel undergoing label noise is independent of and identical

81  to label noise processes in other (even neighboring) pixels. This scenario is analogous to

82  the use of label noise robust pixel-based building extraction methods such as logistic

83  regression (Maas et al., 2016), random forests (Maas et al., 2019), compared to the use

84  of deep learning-based label noise robust building extraction methods such as fully

85  convolutional networks and U-Nets (Zhang et al., 2020). The primary difference

86  between non-deep learning and deep learning-based building extraction is that the

87  former usually uses features from only the pixel being classified, whereas the latter

88  leverages context to predict dense labels for the entire image at once. Feature

89  representation is an important part of deep learning based remote sensing image

90  processing (Jing et al., 2021; He et al., 2021). Modeling of superpixel based label noise

91  process has been conducted for the general computer vision task of semantic

92  segmentation (Lu et al., 2016), but has largely been left unexplored for remote sensing

93  applications. If building extraction can be modeled using a dense prediction approach,

94  we argue that pixel-based label noise robustness approaches can also be extended to

95  dense prediction-based label noise robustness approaches.

96

Figure 1. Some examples of large image tiles from our dataset. (a) Image (b) True clean dense labels (c) Observed dense labels from OpenStreetMap with real world noise

There are various aspects of viewing the label noise generation process. Labeling tools used by human annotators also play a role in determining the label noise processes for dense prediction tasks (Frank et al., 2017). Simulated noise is common in label noise robust image classification scenarios (Ghosh et al., 2017; Rolnick et al., 2017; Patrini et al., 2017) and can be extended to dense prediction-based building extraction as well, however, we have access to data with real-world dense label noise. It is also important to acknowledge the limitations of simulated noise when compared to real-world noise (Jiang et al., 2020). Label noise processes can broadly be categorized by their randomness (Frénay, B., & Verleysen, 2013). For example, if certain building superpixels are being omitted in the observed labels, the question arises, are these

111 buildings being selected totally at random, or are certain types of buildings, perhaps

112 newly constructed buildings, being omitted. Randomness characterizes label noise

113 processes. Identifying this randomness is crucial for modeling label noise robust

114 learning systems. Randomness is unique to each dataset and is estimated prior to

115 modeling solutions.

116

117       We have quantified the effects of label noise on evaluation regimes for this

118 dataset and found that deep neural networks for semantic segmentation are intrinsically

119 robust to real world random label noise, specially aided if data augmentation and

120 regularization are introduced (Ahmed et al., 2020). However, robustness to label noise

121 is achieved as a by-product of overfitting-reduction schemes, and therefore the

122 modelling of label noise is implicit. In this work, we explicitly model dense label noise

123 as a subtask of building extraction, and show improved performance on independent test

124 set.

125

126       The primary objective of this study is to analyze label noise robustness of deep

127 semantic segmentation networks using our proposed evaluation regime. State-of-the-art

128 methods for deep learning-based building extraction from remotely sensed imagery

129 usually perform model evaluation using noisy labels as ground truth, we test the effects

130 of performing model evaluation against noisy labels and clean labels. Our contributions

131 are as follows. We outline approaches for modeling dense label noise and formalize a

132 multi-view and multi-scale taxonomy of label noise. We propose two new model

133 frameworks for building extraction from aerial drone imagery under dense label noise,

134 and consequently two new network architectures. Our network architectures

135 approximate the dense label noise characterizing error matrices as an intermediate step

136     to improve performance. Approximating intermediate error matrices can generate label

137     noise saliency/heat maps. We have made our dataset and method implementations

138     publicly available

139     (https://drive.google.com/uc?id=1UUGeewOaNzv_8kMGXOgEzR8_QKPlPsr8)

140     (https://github.com/nahian-ahmed/dense-label-noise).


141     **Dense label noise models**


142     *Preliminaries and definitions*

143     Formulations on label noise in non-dense approaches are well defined and studied

144     (Frénay, B., & Verleysen, 2013; Frénay and Kabán, 2014). Label noise processes are

145     defined based on the nature of the randomness of the process in question. The three
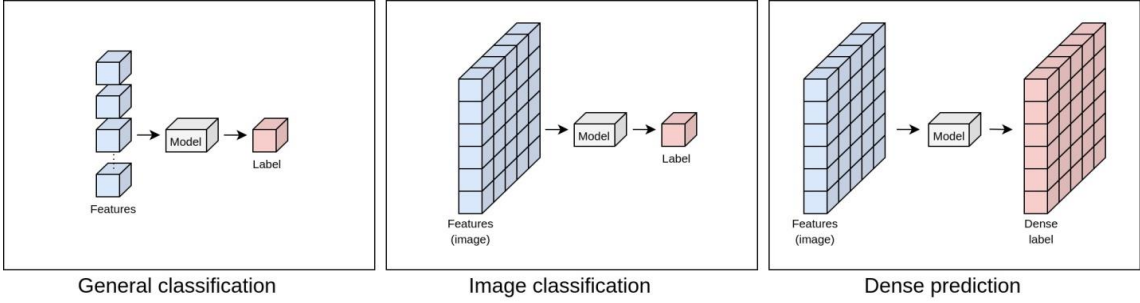
146     types of noisy labels are -

147

148     (1)     Noisy completely at random (NCAR) labels, where labels are flipped completely

149          independent of features and class label,

150     (2)     Noisy at random (NAR) labels, where labels are flipped independent of features

151          but dependent on class label,

152     (3)     Noisy not at random (NNAR) labels, where labels are flipped depending on

153          features and class label.

154

155          These label noise models are equally highly apt at expressing label noise

156     processes for classification on tabular data and image data. In image classification, each

157     image is assigned a single label; though the feature is more complex, the target is still a

158     single label and therefore the non-dense label noise models are sufficient in describing

159     the noise processes. However, for dense prediction, tasks the notation and process

160     models for label noise need extension. We have formulated label noise models for our

161     image segmentation task by extending the label noise models presented by Frénay, B.,

162     & Verleysen, (2013) and design according to pixel-wise and dense dependencies. Dense

163     label noise models can represent complex non-linear and fully-connected statistical

164     dependencies between the image tensors and label tensors. Fig. 2 shows the conceptual

165     differences between the label generation process for the general classification, image

166     classification, and dense prediction.

167



168           General classification        Image classification        Dense prediction

169 Figure 2. Differences among general classification, image classification and dense
170 prediction

171

172         Given an observed noisy dense label $\widetilde{Y} \in \{0,1\}^{n_h \times n_w}$ and its corresponding true

173     clean dense label $Y \in \{0,1\}^{n_h \times n_w}$, where height and width of image tile is $n_h$ and $n_w$

174     respectively. Indexing $n_h$ by $i$ and indexing $n_w$ by $j$ , $Y_{i,j}$ represents the pixel in $i$-th

175     row and $j$-th column of a label tile, $\widetilde{Y}_{i,j}$ is considered to be noisy if $\widetilde{Y}_{i,j} \neq Y_{i,j}$. We

176     extend the binary variable random in Frénay and Verleysen (2013) indicating presence

177     of label noise, to dense prediction settings. We define the *error matrix* $E \in \{0,1\}^{n_h \times n_w}$

178     as the matrix indicating positions of pixels with label noise. Thus, $E_{i,j} = 1$ when $\widetilde{Y}_{i,j} \neq$

179     $Y_{i,j}$ and $E_{i,j} = 0$ if $\widetilde{Y}_{i,j} = Y_{i,j}$. For binary labels, if the current observed pixel label $\widetilde{Y}_{i,j}$

180     and its labeling error presence $E_{i,j}$ is known, the true label $Y_{i,j}$ can directly be computed

181     by flipping the observed label when the pixel label in question is deemed to be noisy.

182  Each element $E_{i,j}$ is a binary random variable indicating if $Y_{i,j}$ is to be noised or not.

183  The relationship among $Y$, $\widetilde{Y}$ and $E$ in matrix form can be defined as

184

$$Y = \left|\widetilde{Y} - E\right| \tag{1}$$

185

186      All operations in Eq. (1) are element-wise matrix operations. Table 1 confirms

187  Eq. (1) and shows the different cases that may arise from combinations of $Y_{i,j}$ and $\widetilde{Y}_{i,j}$.

188  When the true label and observed label are the same (row no. 1 and 2 in Table 1), label

189  noise is absent; when the true label and observed label are not equal (row no. 3 and 4 in

190  Table 1), label noise is present. Given knowledge on the observed noisy label and error

191  matrix, the clean label can directly be computed using Eq. (1).

192

193  **Table 1.** The four possible cases arising from combinations of $Y_{i,j}$ and $\widetilde{Y}_{i,j}$

| No | Case | Label noise | $Y_{i,j}$ | $\widetilde{Y}_{i,j}$ | $E_{i,j}^{+}$ | $E_{i,j}^{-}$ | $E_{i,j}$ | $\left|\widetilde{Y} - E\right|$ |
|----|------|-------------|-----------|----------------------|---------------|---------------|-----------|----------------------------------|
| 1 | True negative observed pixel label | No | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | True positive observed pixel label | No | 1 | 1 | 0 | 0 | 0 | 1 |
| 3 | False positive observed pixel label | Yes | 0 | 1 | 1 | 0 | 1 | 0 |
| 4 | False negative observed pixel label | Yes | 1 | 0 | 0 | 1 | 1 | 1 |

194

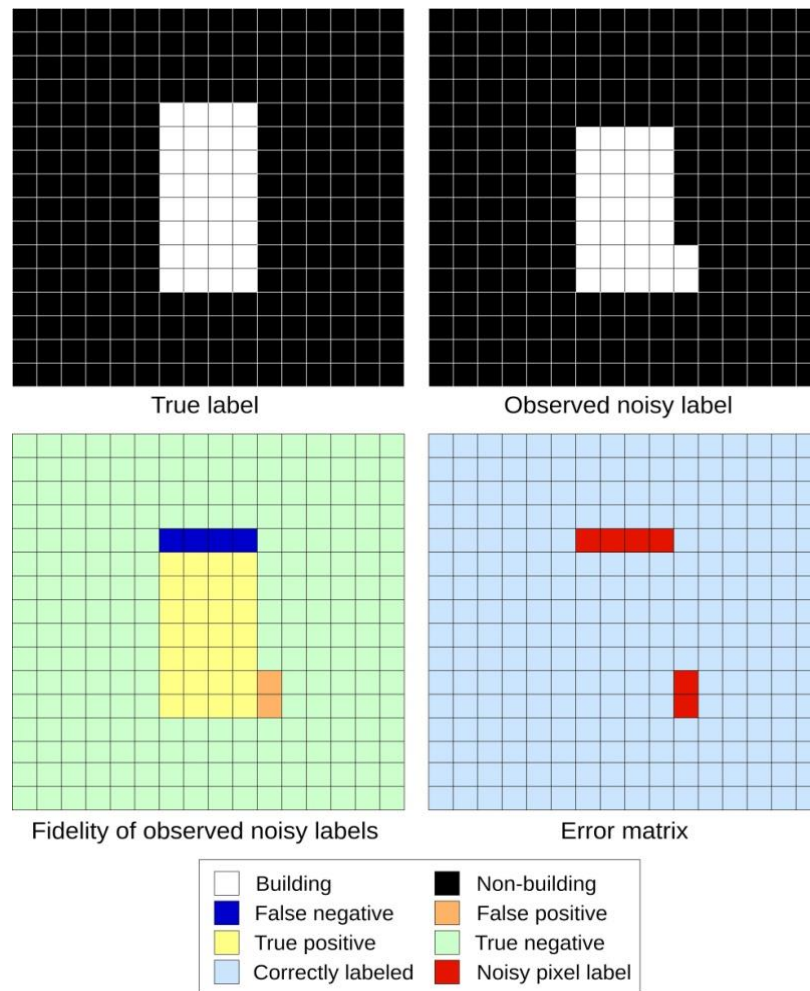195      The error matrix is the absolute difference between the true and observed labels

196

$$E = \left|Y - \widetilde{Y}\right| = \left|\widetilde{Y} - Y\right| \tag{2}$$

197

198        Let, the error matrix denoting *false positive* observed labels be $E^+ \in \{0,1\}^{n_h \times n_w}$

199    and the error matrix denoting *false negative* observed be $E^- \in \{0,1\}^{n_h \times n_w}$. Thus, $E$ is

200    the element-wise logical 'or' (expressed as summation) of $E^+$ and $E^-$ in matrix form,

201

$$E = E^+ + E^- \tag{3}$$

202

203        Fig. 3 shows an example of how label noise arises from disagreements between

204    the true label and observed label, displaying that a few positive pixel labels were missed

205    and a few true negative pixel labels were labeled as positives.

206



True label            Observed noisy label

Fidelity of observed noisy labels        Error matrix

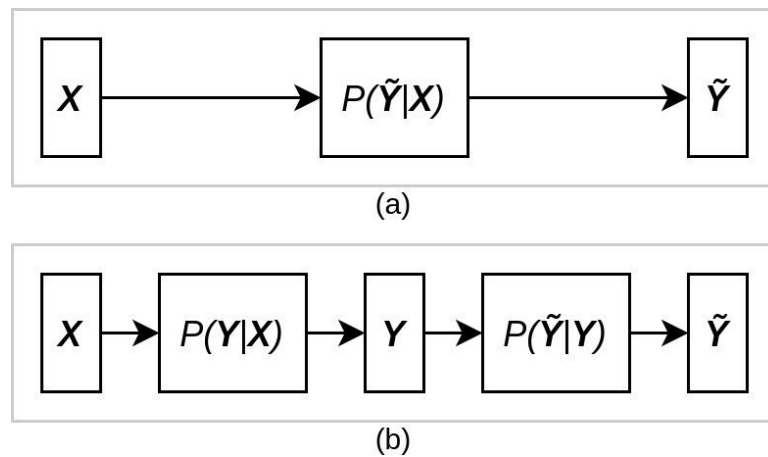| | Building | | Non-building |
| --- | --- | --- | --- |
| | False negative | | False positive |
| | True positive | | True negative |
| | Correctly labeled | | Noisy pixel label |

207

208 Figure 3. Example of how observed noisy dense labels differ from their corresponding

209 true dense labels. A 16 x 16 pixel image is used for demonstration. The error matrix $E$ is

210 shown in the bottom right subfigure, indicating positions of noisy pixel labels.

211

212      The label noise process involves the *corruption* of clean labels (Fig. 4). In

213 general learning schemes for building extraction, it is assumed that the observed labels

214 are clean and are directly used for learning/evaluation (Fig. 4(a)). However,

215 acknowledgement of label noise assumes the intermediary distribution of clean labels

216 over the images to be the clean labels and models the label noise process as the

217 distribution of observed noisy labels over the true clean labels (Fig. 4(b)), which means

218 that when label noise is present, the ground truth clean labels are unobserved

219

220



(a)

(b)

221

222 Figure 4. Observed label generation processes (a) Modeled without noise-free labels (b)

223 Modeled through noise-free labels

224

225 Having defined the important concepts i.e. $Y$, $\widetilde{Y}$ and $E$, for modeling dense label noise

226 processes, we move on to define the statistical dependencies for learning dense prediction

227 (Fig. 5). There are two main models -

228

229     ● *Pixel-wise models:* perform pixel classification using features from only the

230       corresponding input pixels (Fig. 5). Therefore, changing tile sizes does not have

231       significant effects if the same pixels are provided for training and testing

232       because only pixel-wise mappings are learned; features from neighboring pixels

233       are not considered. Without context, the rooftop of a building and a road may

234       appear identical to the model. However, learning pixel-wise mapping is common

235       in non-deep learning approaches to building extraction. Given, the input tensor

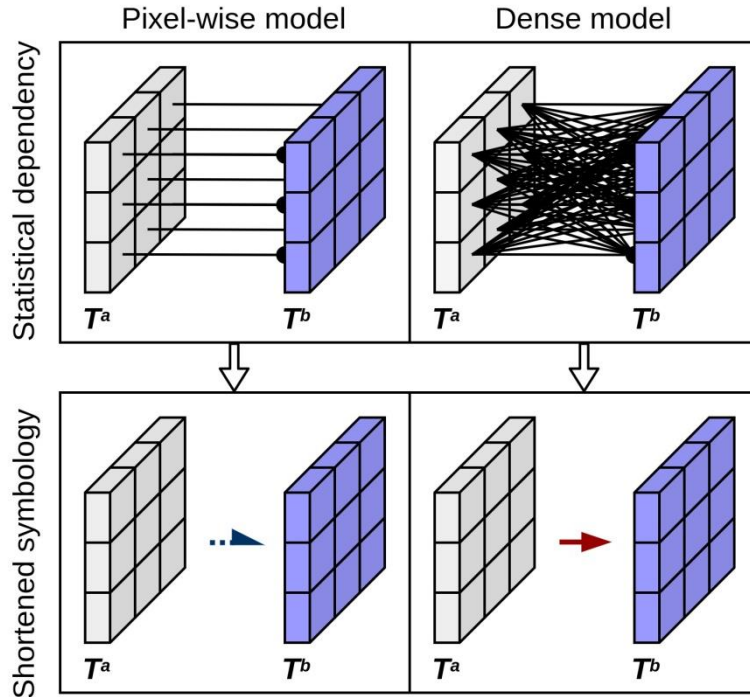236       $T^a$ and its dependent output tensor $T^b$, the pixel wise models learn,

$$P(T_{i,j}^b | T_{i,j}^a) \qquad (4)$$

237     ● *Dense models:* generates labels for pixels using features from all pixels of the

238       input tensor (Fig. 5). The model estimates each $P(T_{i,j}^b | T^a)$ and then uses the

239       product chain rule to learn $P(T^b | T^a)$,

$$P(T^b | T^a) = \prod_{i=1}^{n_h} \prod_{j=1}^{n_w} P(T_{i,j}^b | T^a) \qquad (5)$$

240     As Fig. 5 shows, we represent fully connected dense mappings using a red full

241 red arrow with continuous line and pixel wise mappings using a blue half-arrow with

242 dotted line.

243

244

Figure 5. Shortened symbology of statistical dependencies considered in pixel wise

models and dense models. In the pixel-wise model, each $T^b_{i,j}$ is only dependent on $T^a_{i,j}$.

In the dense model, each $T^b_{i,j}$ is dependent on the the entire matrix $T^a$ indicating fully

connectedness.

249

### *Taxonomy of dense label noise models*

The three types of label noise in Frénay and Verleysen (2013) are categorized according

to randomness. We refer to this approach as taxonomy characterized by randomness.

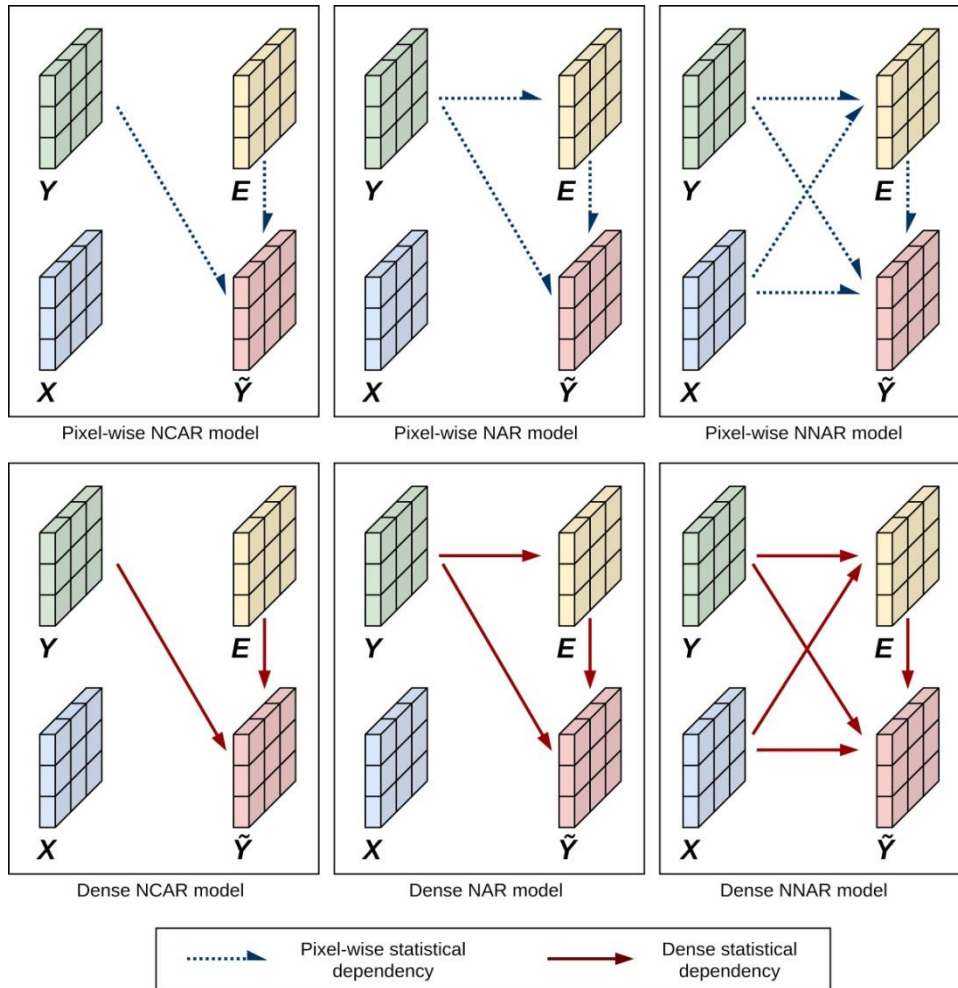However, in the context of dense prediction, structure (spatial information) in dense

labels also plays a role in label noise processes. We define the taxonomy of dense label

noise models. Given the two types of mapping models (pixel-wise and dense) and the

three types of stochasticity defined label noise processes (NCAR, NAR and NNAR),

there are six possible models (Fig. 6).

258

Figure 6. Statistical dependencies of different types of pixel based and dense label noise models. The dependency between *X* and *Y* are not shown for brevity.

(1) *Pixel-wise NCAR model*: NCAR models are class independent, therefore the only noise parameters for a pixel-wise NCAR model would be the *probability of error* $p_e = P(\widetilde{Y}_{i,j} \neq Y_{i,j})$. It is important to note that $p_e$ is constant for all pixels, and therefore NCAR models cannot model non-uniform label noise. All $E_{i,j}$ would have the same values because the probability of a pixel being noisy is constant and not dependent on any variables. The error matrix $E$ is completely independent (pixel-wise NCAR model in Fig. 6). For binary classification (which is our case for the pixel-wise models) having $p_e = 1/2$ would render the labels useless and inadequate to learn from (Angluin

271 and Laird, 1988). Furthermore, since NCAR models are class independent, asymmetric

272 noise cannot be modeled as well. NCAR models assume that labels of all classes have

273 equal chances of being observed as noisy labels. In real world settings, this is rarely the

274 case. For example, in building extraction tasks, the positive class is much more prone to

275 label noise. Furthermore, the positive class is also the minority class in most imbalanced

276 building extraction datasets.

277

278     (2) *Pixel-wise NAR model:* NAR models are able to model asymmetric and non-

279 uniform label noise processes. Each $\boldsymbol{E}_{i,j}$ is dependent on each $\boldsymbol{Y}_{i,j}$, which in turn affects

280 each $\widetilde{\boldsymbol{Y}}_{i,j}$ (pixel-wise NCAR model in Fig. 6). The probability of a specific label being

281 observed as another label is modelled using the transition matrix (Lawrence and

282 Schölkopf, 2001; Pérez et al., 2007). We define the *transition matrix* for noisy dense

283 binary labels as

284

$$\boldsymbol{\gamma} = \begin{bmatrix} \gamma_{0,0} & \gamma_{0,1} \\ \gamma_{1,0} & \gamma_{1,1} \end{bmatrix}$$

285

$$= \begin{bmatrix} P(\widetilde{\boldsymbol{Y}}_{i,j} = 0 | \boldsymbol{Y}_{i,j} = 0) & P(\widetilde{\boldsymbol{Y}}_{i,j} = 0 | \boldsymbol{Y}_{i,j} = 1) \\ P(\widetilde{\boldsymbol{Y}}_{i,j} = 1 | \boldsymbol{Y}_{i,j} = 0) & P(\widetilde{\boldsymbol{Y}}_{i,j} = 1 | \boldsymbol{Y}_{i,j} = 1) \end{bmatrix} \tag{6}$$

286

287     The conditional probabilities in Eq. (6) can be estimated from the observed and

288 corresponding clean labels. It is important to note that, the transition matrix is the same

289 for all $\widetilde{\boldsymbol{Y}}_{i,j}$ (and hence for all $\boldsymbol{Y}_{i,j}$). For uniform noise in dense binary labels, the

290 transition matrix becomes

291

$$\gamma = \begin{bmatrix} 1 - p_e & p_e \\ p_e & 1 - p_e \end{bmatrix} \tag{7}$$

292

293    (3) *Pixel-wise NNAR model*: In the case of NNAR models, the error matrix $\boldsymbol{E}$ is

294    dependent on the features as well (pixel-wise NNAR model in Fig. 6). The observed

295    pixel label $\widetilde{\boldsymbol{Y}}_{i,j}$ is dependent on $\boldsymbol{E}_{i,j}$ and $\boldsymbol{Y}_{i,j}$; if $\boldsymbol{E}_{i,j} = 1$, $\boldsymbol{Y}_{i,j}$ is flipped to get $\widetilde{\boldsymbol{Y}}_{i,j}$,

296    otherwise $\widetilde{\boldsymbol{Y}}_{i,j} = \boldsymbol{Y}_{i,j}$. The probability of error is a function of the pixel-wise feature

297    and pixel-wise true label,

298

$$p_e(\boldsymbol{X}_{i,j}, \boldsymbol{Y}_{i,j}) = P(\boldsymbol{E}_{i,j} = 1 | \boldsymbol{X}_{i,j} = x, \boldsymbol{Y}_{i,j} = y) \tag{8}$$

299

300    (4) *Dense NCAR model:* In the dense NCAR model, every $\widetilde{\boldsymbol{Y}}_{i,j}$ is affected by the

301    entire error matrix $\boldsymbol{E}$, and not just $\boldsymbol{E}_{i,j}$ (which is the case for the pixel-wise NCAR

302    model). Spatial information about label noise in terms of context (as opposed to pixel-

303    based information) can be modeled. Every $\boldsymbol{E}_{i,j}$ need not be constant; however, they are

304    still completely independent (of each other and of any other random variable) and thus

305    completely random (dense NCAR model in Fig. 4).

306

307    (5) *Dense NAR model*: The dense NAR model allows modeling asymmetric

308    dense label noise, which is not possible using the dense NCAR model. Unlike the pixel-

309    wise NAR model, the transition matrix for each $\widetilde{\boldsymbol{Y}}_{i,j}$ can be distinct and independent of

310    each other. The transition matrix for $\widetilde{\boldsymbol{Y}}_{i,j}$ in a dense NAR model can be defined as

$$\gamma^{(i,j)} = \begin{bmatrix} \gamma_{0,0}^{(i,j)} & \gamma_{0,1}^{(i,j)} \\ \gamma_{1,0}^{(i,j)} & \gamma_{1,1}^{(i,j)} \end{bmatrix} \tag{9}$$

311        The error matrix $E$ is directly dependent on the true dense label $Y$ (dense NAR

312    model in Fig. 6), but independent of the dense features $X$.

313

314        (6) *Dense NNAR model:* In the dense NNAR model all pixels from the image

315    affect the probabilities of label noise in certain observed pixels (dense NNAR model in

316    Fig. 6). Every $E_{i,j}$ is affected by the entire image tensor $X$, and every $\widetilde{Y}_{i,j}$ is affected by

317    the entire error matrix $E$. The error matrix can be estimated based on the observed dense

318    label $\widetilde{Y}$ and dense feature tensor $X$. We essentially model the conditional distribution of

319    the error matrix $E$, given the feature tensor $X$ and the observed dense label $\widetilde{Y}$ (Eq. (10)).

320    This estimated error matrix can then be used for generating the true labels using Eq. (1).

321

$$P(E \mid \widetilde{Y}, X) = \prod_{i=1}^{n_h} \prod_{j=1}^{n_w} P(E_{i,j} \mid \widetilde{Y}, X) \tag{10}$$

322

323    **Materials and methods**
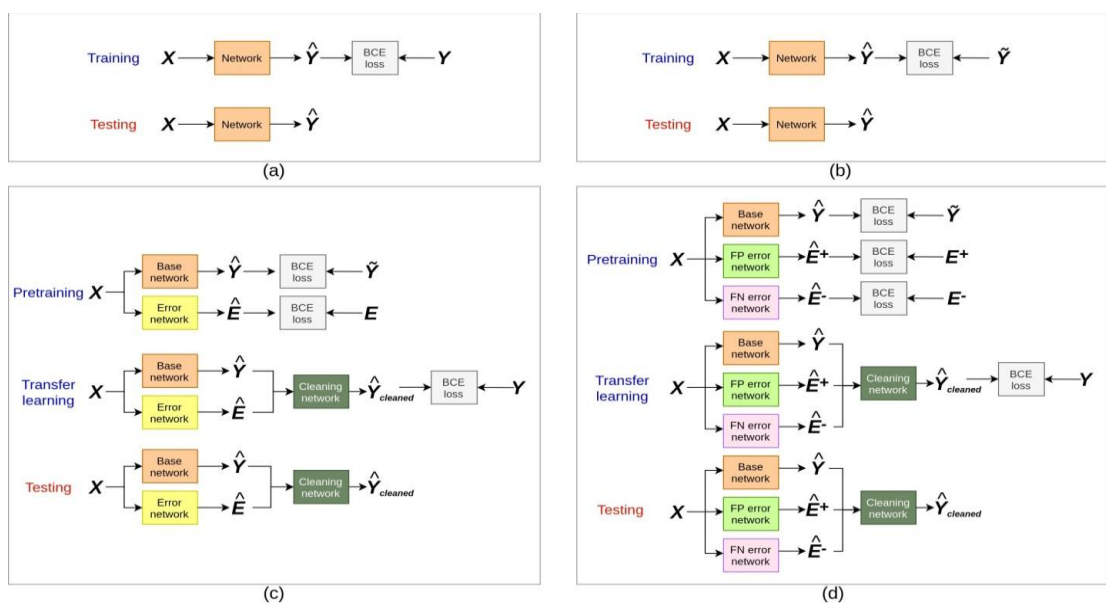
324    *Data*

325    The dataset consists of 258 large 512x512 ultra-high-resolution aerial image tiles over

326    the Kutupalong mega camp collected by the United Nations International Organization

327    for Migration on September 17, 2018. Kutupalong is the largest of the camps,

328    comprised of several sub-camps, situated in the south-eastern border region of

329    Bangladesh which acted as the corridor for the Rohingya refugees migrating from

330     Myanmar. For our case, $n_h = n_w = 512$. The observed noisy labels are collected from

331     OpenStreetMap. The true clean labels are obtained by relabeling performed by the

332     authors. The dataset is randomly split in half for denoting training and testing data.

333     Images have three channels/bands — Red, Green and Blue — with a spatial resolution

334     of 10 cm. These images have very high data quality i.e. without cloud or shadow cover

335     being collected by low flying unmanned aerial vehicles (UAVs) and capture fine-

336     grained details of the physical environments where the buildings are located. The

337     general error matrices are computed using Eq. (2), whereas the FP and FN error

338     matrices are computed without taking the absolute value, rather using the signed/un-

339     signedness of the difference matrix. Our dataset is relatively smaller than most

340     commonly used datasets for building extraction (such as Massachusetts, Potsdam and

341     Vaihingen datasets), this is because we have had to re-label all of our training and test

342     data by hand for obtaining the noise-free true clean labels, which is very time-

343     consuming. Moreover, datasets for semantic segmentation/dense prediction with the

344     corresponding observed labels (with real-world label noise) and counterpart clean labels

345     are virtually non-existent. Our dataset is unique in that aspect, since, having access to

346     the observed noisy labels and clean labels is crucial for obtaining ground truth error

347     matrices (Eq. (1)). It is important to note that the error matrices are only required for

348     pretraining the dense label noise prediction models, during testing/evaluation the

349     models directly output building maps corrected by error matrices.

350     ***Model frameworks***

351     The true clean dense label is solely dependent on the feature tensor in all six noise

352     models (caption of Fig. 6). The features (from satellite/aerial images), used for

353     approximating true labels, can be compared to the observed noisy label to obtain the

354     error matrix; the features have an important role in determining the observed label.

355    Therefore, the dense NNAR model is most suitable for expressing commonly observed

356    registration errors. Currently, deep learning is the state-of-the-art system for automated

357    building extraction (Vakalopoulou et al., 2015; Huang et al., 2016; Chen et al., 2017;

358    Yuan, 2017; Yang et al., 2018; Ji et al., 2018; Xu et al., 2018; Shrestha and Vanneschi,

359    2018; Boonpook et al., 2021; Sun et al., 2021). Fig 7(a) and 7(b) show the generally

360    used learning systems for deep learning-based building extraction i.e. with clean labels

361    (Fig. 7(a)) and with noisy labels (Fig. 7(b)). We propose two new models for automated

362    building extraction, and consequently, two novel network architectures, where error

363    matrices are approximated as an intermediate step (Fig 7(c) and Fig. 7(d)). As discussed

364    later, we draw from the dense NNAR model in modelling our learning frameworks. The

365    formulated dense noise models ultimately determine the architecture of the neural

366    networks. The base network in Fig. 7(a) represents the statistical dependency between

367    the feature and label tensors in the dense NNAR model (Fig. 6). Similarly, the error

368    matrix network in Fig. 7(a) represents the statistical dependency between the feature

369    and error matrix tensors in the dense NNAR model (Fig. 6). We elaborate on the model

370    frameworks, network architectures, learning and evaluation approaches.



371

372    Figure 7. Training and testing approaches (a) With clean labels - control, CL model (b)

373    With noisy labels - NL model (c) With intermediate error matrix approximation - I-EM

374    model (d)  With intermediate FP and FN error matrices approximation - I-FPFN-EM

375    model; BCE - binary cross entropy; FP - false positive, FN - false negative

376    *Intermediate error matrix (I-EM) model*

377    The first proposed intermediate error matrix (I-EM) model approximates error matrices

378    as an intermediate step of approximating building/non-building predictions. The noisy

379    observed labels are learned by the base network in Fig. 7(c) approximated as the mean

380    of the distribution in Eq. (11). The noisy observed labels are learned by the error

381    network in Fig. 7(c) approximated as the mean of the distribution in Eq. (12). Finally,

382    the outputs from the error matrix (EM) model and the observed label model are used

383    together by the cleaning network in Fig. 7(c) to learn noise free label approximation in

384    Eq. (13). Viewing the model framework from an end-to-end fashion in terms of testing

385    indicates (Testing in Fig. 7(c)) in Eq. (14).

386

387

$$P(\widetilde{Y}|X) = \prod_{i=1}^{n_h} \prod_{j=1}^{n_w} P(\widetilde{Y}_{i,j}|X) \tag{11}$$

$$P(E|X) = \prod_{i=1}^{n_h} \prod_{j=1}^{n_w} P(E_{i,j}|X) \tag{12}$$

$$P(\boldsymbol{Y}|\widetilde{\boldsymbol{Y}}, \boldsymbol{E}) = \prod_{i=1}^{n_h} \prod_{j=1}^{n_w} P(\boldsymbol{Y}_{i,j}|\widetilde{\boldsymbol{Y}}, \boldsymbol{E}) \qquad (13)$$

$$P(\boldsymbol{Y}|\boldsymbol{X}, \boldsymbol{E}) = \prod_{i=1}^{n_h} \prod_{j=1}^{n_w} P(\boldsymbol{Y}_{i,j}|\boldsymbol{X}, \boldsymbol{E}) \qquad (14)$$

388  *Intermediate FP and FN error matrix (I-FPFN-EM) model*

389  The second proposed intermediate FP and FN error matrix (I-FPFN-EM) model

390  approximates the FP and FN error matrices separately as an intermediate step of

391  approximating building/non-building predictions. The noisy observed labels are learned

392  by the base network in Fig. 7(d) approximated as the mean of the distribution in Eq.

393  (11). The FP (false positive) error matrix is learned by the FP error network in Fig. 7(d)

394  approximated as the mean of the distribution in Eq. (15). The FN (false negative) error

395  matrix is learned by the FNM error network in Fig. 7(d) approximated as the mean of

396  the distribution in Eq. (16). Finally, the outputs from the FP and FN error matrix

397  models, and the observed label model are used together by the cleaning network in Fig.

398  7(d) to learn noise free label approximation in Eq. (17). We refer to the FP error matrix

399  model as the FP-EM model and the FN error matrix model as the FN-EM model.

400  Viewing the model framework from an end-to-end fashion in terms of testing indicates

401  (Testing in Fig. 7(d)) in Eq. (18).

$$P(\boldsymbol{E}^+|\boldsymbol{X}) = \prod_{i=1}^{n_h} \prod_{j=1}^{n_w} P(\boldsymbol{E}_{i,j}^+|\boldsymbol{X}) \qquad (15)$$

$$P(\boldsymbol{E}^-|\boldsymbol{X}) = \prod_{i=1}^{n_h} \prod_{j=1}^{n_w} P(\boldsymbol{E}_{i,j}^-|\boldsymbol{X}) \tag{16}$$

$$P(\boldsymbol{Y}|\widetilde{\boldsymbol{Y}}, \boldsymbol{E}^+, \boldsymbol{E}^-) = \prod_{i=1}^{n_h} \prod_{j=1}^{n_w} P(\boldsymbol{Y}_{i,j}|\widetilde{\boldsymbol{Y}}, \boldsymbol{E}^+, \boldsymbol{E}^-) \tag{17}$$

$$P(\boldsymbol{Y}|\boldsymbol{X}, \boldsymbol{E}^+, \boldsymbol{E}^-) = \prod_{i=1}^{n_h} \prod_{j=1}^{n_w} P(\boldsymbol{Y}_{i,j}|\boldsymbol{X}, \boldsymbol{E}^+, \boldsymbol{E}^-) \tag{18}$$

402 *Network architectures*

403 Each intermediate network has four downsampling blocks and four upsampling blocks.

404 We use vanilla U-Nets with approximately 0.5 million parameters for intermediate

405 learning steps. The U-Net/autoencoder architecture is common for building extraction

406 tasks (Wang et al., 2020; Guo et al., 2020). The use of step-wise concatenation of

407 models has been employed for building extraction (Shao et al., 2020). Each

408 downsampling block has two convolutional layers punctuated by a single dropout layer,

409 which is then downsampled to half the output row and column size using max pooling.

410 Each upsampling block also has two convolutional layers punctuated by a single

411 dropout layer, which is then upsampled to double the output row and column size using

412 interpolation. We use the binary cross entropy loss function as it is commonly used for

413 most binary building extraction tasks (Ahmed et al., 2020). For the I-EM model (Fig.

414 8(a)) the outputs of the base network and error network are concatenated and fed to the

415 cleaning network. For the I-FPFN-EM network the outputs of the base network, FP

416 error matrix network and FN error matrix network are all fed into the cleaning network.

417 Please note that intermediate predictions of observed labels and error matrices (general,

418    FP and FN) are in the form of soft pixel level labels i.e. they are not converted to hard

419    labels based on threshold values. The I-EM model and I-FPFN-EM models have

420    approximately 1.5 million and 2 million parameters respectively. S1 details the network

421    architecture for NL, CL, EM, FP-EM and FN-EM models, Fig. S2 and Fig. S3 in

422    supplementary material contains the detailed network architectures of the I-EM and I-

423    FPFN-EM model respectively.



424

425    Figure 8. Proposed network architectures for building extraction under label noise (a) I-

426    EM model (b) I-FPFN-EM model

427    *Learning*

428    The I-EM model and I-FPFN-EM model are trained in two steps.

429     • *Step 1 - Pre-training*: For learning the parameters of the base and error

430          networks. Individual auto-encoders with skip connections are trained. For the I-

431          EM model, the base network is trained using the images $X$ as features and $\widetilde{Y}$ as

432          targets, the error network is trained using the images $X$ as features and $E$ as

433          targets. For the I-FPFN-EM model, the base network is also trained using the

434          images $X$ as features and $\widetilde{Y}$ as targets, the FP error network is trained using the

435          images $X$ as features and $E^+$ as targets, and the FN error network is trained

436          using the images $X$ as features and $E^-$ as targets.

437     • *Step 2 - Transfer learning*: After the base networks and error networks (general

438          for I-EM; FP and FN for I-FPFN-EM) are trained, their outputs are concatenated

439          and fed into the cleaning networks. In order to train the cleaning network, the

440          layers in the base and error networks are frozen i.e. they are set as non-trainable.

441          In this second step of training, the entire network is trained in an end-to-end

442          fashion against clean labels.

443

444 The baseline CL model and NL model both have approximately 0.5 million parameters.

445 The I-EM model and I-FPFN-EM models have approximately 1.5 million and 1.5

446 million parameters respectively. This larger number of parameters are due to the error

447 matrix networks and the cleaning networks used in the I-EM model and the I-FPFN-EM

448 models. The general error matrix sub-model in the I-EM model, and each of the false

449 positive error matrix model and the false negative error matrix models all have

450 approximately 0.5 million parameters. The time complexity of the I-EM model and I-

451 FPFN-EM model are also increased proportional to the increase of number of

452 parameters with respect to the CL and NL models. The total time needed for training the

453     sub-models of the I-EM model is triple that of the CL or NL models, and the total time

454     needed for training the I-FNFN-EM models is quadruple that of the CL or NL models.


455     *Method comparison*

456     In order to assess the qualitative and quantitative advantages/disadvantages of our two

457     proposed models, we also compare against generally used model frameworks for

458     automated building extraction. We compare four different deep learning-based building

459     segmentation models,

460

461         (1) Noisy label (*NL*) model (Ahmed et al., 2020): Dense building extraction with

462             noisy labels.

463         (2) Clean label (*CL*) model (Ahmed et al., 2020): Dense building extraction with

464             clean labels (control).

465         (3) *I-EM model*: The first proposed model described above.

466         (4) *I-FPFN-EM* model: The second proposed model described above.

467

468             Other than the CL and NL models in Ahmed et al., (2020), no other study

469     presents dataset/methods for dense prediction of label noise using clean and noisy labels

470     with real world noise. The threshold value determines the boundary value and

471     consequently the binary class label of each pixel. We vary the threshold for each model

472     with low (0.25), medium (0.5) and high (0.75) values to convert the soft labels (between

473     0 and 1 inclusive) to hard labels (0 or 1).


474     *Performance evaluation metrics*

475     We calculate the total number of true positives (TP), true negatives (TN), false positive

476     (FP) and false negative (FN) predictions on the approximately 33 million pixels of

477    testing data. Concurring to most building extraction scenarios, our dataset is also quite

478    imbalanced, being negative heavy. Therefore, we calculate the precision (Eq. (19),

479    recall (Eq. (20)), F1-score (Eq. (21) and Intersection-over-Union ($IoU$) (Eq. (22)).

480

$$Precision = \frac{TP}{TP + FP} \tag{19}$$

481

$$Recall = \frac{TP}{TP + FN} \tag{20}$$

482

$$F1 - score = \frac{2TP}{2TP + FP + FN} \tag{21}$$

483

$$IoU = \frac{TP}{TP + FP + FN} \tag{22}$$

484

485    **Results and discussion**

486    *Quantitative evaluation of performance*

487    The CL model provides the control/baseline against which we compare our two

488    proposed models since it represents the ideal scenario when the investigator has access

489    to both images and clean labels. Our I-FPFN-EM model at 0.5 medium threshold (row

490    no. 11 in Table 2) has the highest $IoU$ score (0.78514), which provides a gain of 2.74%

491    over the traditional CL model trained on clean labels (0.75768) and a gain of 25.65%

492    over the observed noisy labels with $IoU$ score of  0.52857. Similarly, our I-FPFN-EM

493    model at 0.5 threshold has the highest F1-score (0.87964), which provides a gain of

494    1.75% over the traditional model trained on clean labels with an F1-score of 0.86214,

495    and gain of 18.8% over the observed noisy labels with an F1-score of 0.69159.

496    Compared to the idealistic CL model, our I-FPFN-EM model has a better F1-score and

497    *IoU* score for high threshold value (0.75) as well, and has comparable/nearly identical

498    performance for low threshold value (0.25). At a threshold value of 0.75, the I-FPFN-

499    EM model (row no. 12 in Table 2) has an F1-score of 0.86009 which is 3.45% higher

500    than the F1-score of the CL model (0.8255) at a threshold value of 0.75. The I-FPFN-

501    EM model at a threshold value of 0.75, achieves an *IoU* score of 0.75453, providing a

502    gain of 5.16% over the CL model with an *IoU* score of 0.70285, at a threshold value of

503    0.75. Our I-FPFN-EM model provides better performance over traditional methods, for

504    the general threshold of 0.5 and the high threshold of 0.75.

505

506        The I-EM has slightly poorer/comparable performance to the CL model. This

507    indicates the importance of differentiating FP and FN error matrices as features, instead

508    of approximating an intermediate general error matrix, since that is the primary

509    conceptual difference between the I-EM model and I-FPFN-EM model. A lower

510    threshold means higher recall and lower precision. A higher threshold means higher

511    precision and lower recall. The threshold value determines the precision recall trade-off.

512    However, both the I-EM and I-FPFN-EM models have much higher recall and slightly

513    lower precision for corresponding threshold values when compared to the CL model. In

514    our case of highly imbalanced data, higher recall is preferred over higher precision.

515

516    Table 2. Performance of the four compared models for building extraction under label

517    noise and the fidelity of observed labels

| No. | Model | Threshold | Precision | Recall | F1-score | IoU |
|-----|-------|-----------|-----------|--------|----------|-----|

| | | | | | | |
|---|---|---|---|---|---|---|
| 1 | | 0.25 | 0.79584 | 0.82862 | 0.8119 | 0.68336 |
| 2 | NL | 0.50 | 0.91337 | 0.56184 | 0.69572 | 0.53342 |
| 3 | | 0.75 | 0.98292 | 0.0948 | 0.17291 | 0.09464 |
| 4 | | 0.25 | 0.79586 | 0.91111 | 0.84959 | 0.73851 |
| 5 | CL | 0.50 | 0.88502 | 0.84041 | 0.86214 | 0.75768 |
| 6 | | 0.75 | 0.93973 | 0.73603 | 0.8255 | 0.70285 |
| 7 | | 0.25 | 0.74541 | 0.93536 | 0.82965 | 0.70889 |
| 8 | I-EM | 0.50 | 0.84473 | 0.85928 | 0.85194 | 0.74207 |
| 9 | | 0.75 | 0.89968 | 0.76947 | 0.8295 | 0.70867 |
| 10 | | 0.25 | 0.76109 | 0.94634 | 0.84366 | 0.7296 |
| 11 | **I-FPFN-EM** | **0.50** | 0.86551 | 0.89424 | **0.87964** | **0.78514** |
| 12 | | 0.75 | 0.92819 | 0.80131 | 0.86009 | 0.75453 |
| 13 | OBSERVED | - | 0.82165 | 0.59708 | 0.69159 | 0.52857 |

518

519        Separated error matrices in the form of FP error matrix and FN error matrix is

520    crucial to surpassing the baseline CL model performance, as our I-EM model has

521    significantly poorer quantitative performance compared to the I-FPFN-EM model.

522    Comparing the I-EM model and the I-FPFN-EM model performances at the three

523    threshold values, the I-FPFN-EM model provides an F1-score increase of 1.4%

524    (0.84366 compared to 0.82965) and $IoU$ score increase of 2.071% (0.7296 compared to

525    0.70889) at a threshold value of 0.25, F1-score increase of 2.77% (0.87964 compared to

526    0.85194) and $IoU$ score increase of 4.307% (0.78514 compared to 0.74207) at a

527    threshold value of 0.5 and F1-score increase of 3.059% (0.86009 compared to 0.8295)

528  and *IoU* score increase of 4.586% (0.75453 compared to 0.70867) at a threshold value
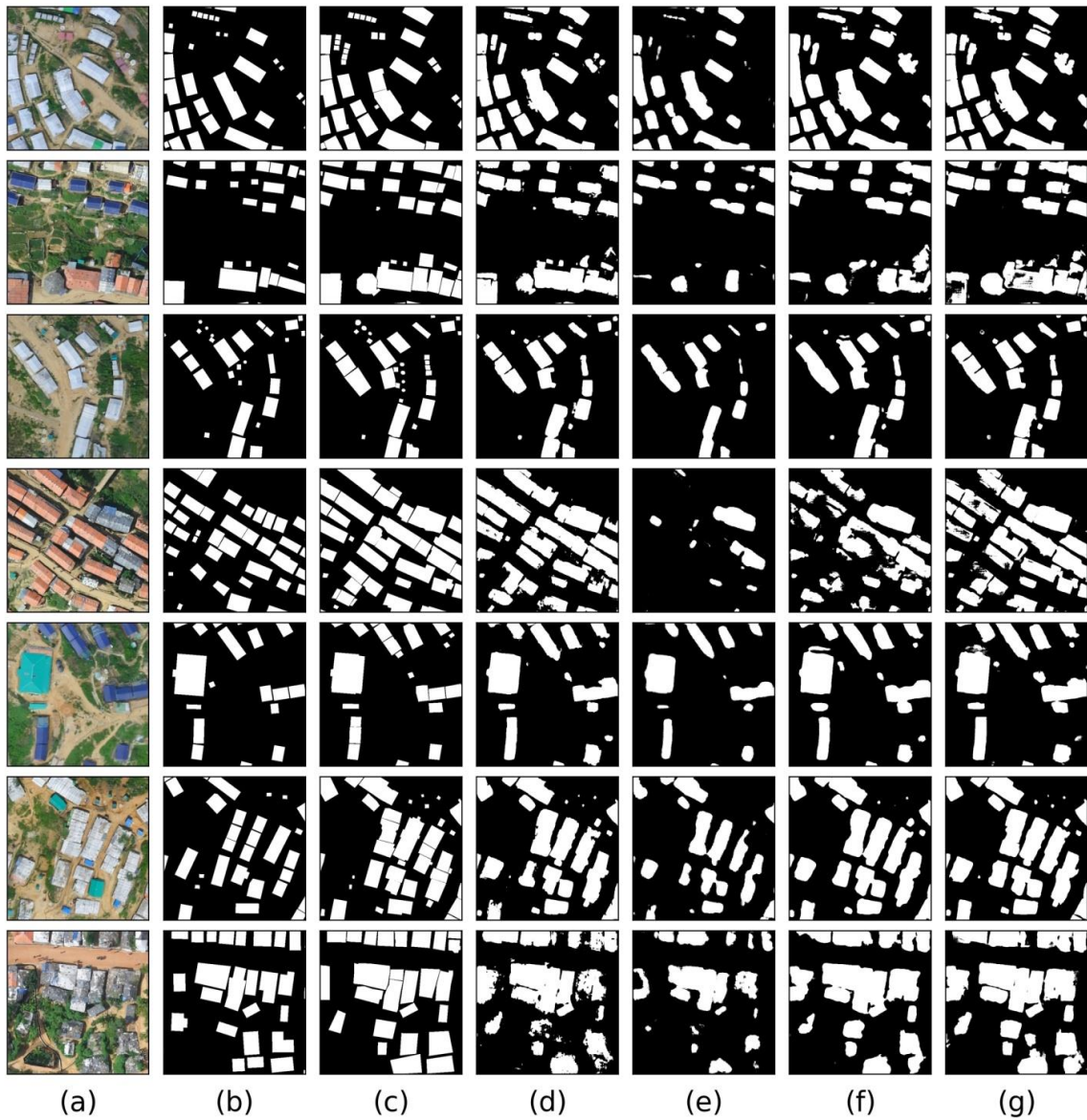
529  of 0.75.

530

531      The traditional model trained against noisy labels (NL model), quite obviously

532  has the poorest performance of the four tested models (row no. 1-3 in Table 2). At high

533  threshold values (0.75) the NL model (row no. 3 in Table 2) predictions become

534  practically useless, yielding an F1-score of 0.17291 and *IoU* score of 0.09464, whereas

535  the CL, I-EM and I-FPFN-EM model have much better performance at a high threshold

536  value of 0.75. The fidelity of noisy labels is also evaluated against the true clean labels

537  (row no. 13 in Table 2). Though the NL model has the poorest performance among four

538  tested models, predictions from the NL model have higher fidelity than the observed

539  labels with real world noise. This is commonly observed for building extraction under

540  real-world noisy conditions (Ahmed et al., 2020).

541  *Qualitative evaluation*

542  From a qualitative viewpoint, the predictions from the four models seem quite similar

543  prior to intensive inspection and photo-interpretation. We show some examples of

544  predictions on image tiles from the test set (Fig. 9). The CL model predictions (Fig.

545  9(d)) have the best qualitative properties, followed by the I-FPFN-EM model

546  predictions (Fig. 9(g)) which sometimes suffers from salt and pepper noise (all

547  predictions in Fig. 9 were made at a threshold value of 0.5 and can be remedied using

548  lower threshold values). Particularly, the I-FPFN-EM model predictions and I-EM

549  model predictions (Fig. 9(f)) for buildings with rare colored roofs (orange painted

550  corrugated metal roofs) contain salt and peppering. Rare colored building rooftops can

551  be challenging to learn due to the comparatively small number of examples in the

552  training set. The NL model predictions completely miss out on entire buildings with

553  orange-colored rooftops (Fig. 9(e)). The last row in Fig. 9 shows the issues of one-

554  storied building rooftops being obstructed partly or completely by vegetation. Building

555  rooftops obstructed by trees and vegetation are not easily detected, as the vegetation

556  over the rooftop is easily confused as non-building regions by the models (last row in

557  Fig. (9)). However, for buildings with vegetation on the rooftops, the I-FPFN-EM

558  model provides less peppering and errors compared to even the CL model (last row in

559  Fig. (9))

560



561        (a)          (b)          (c)          (d)          (e)          (f)          (g)

Figure 9. Examples of building predictions made by different models (a) Image (b) Noisy label (c) Clean label (d) Predictions from CL model (e) Predictions from NL model (f) Predictions from I-EM model (g) Predictions from I-FPFN-EM model
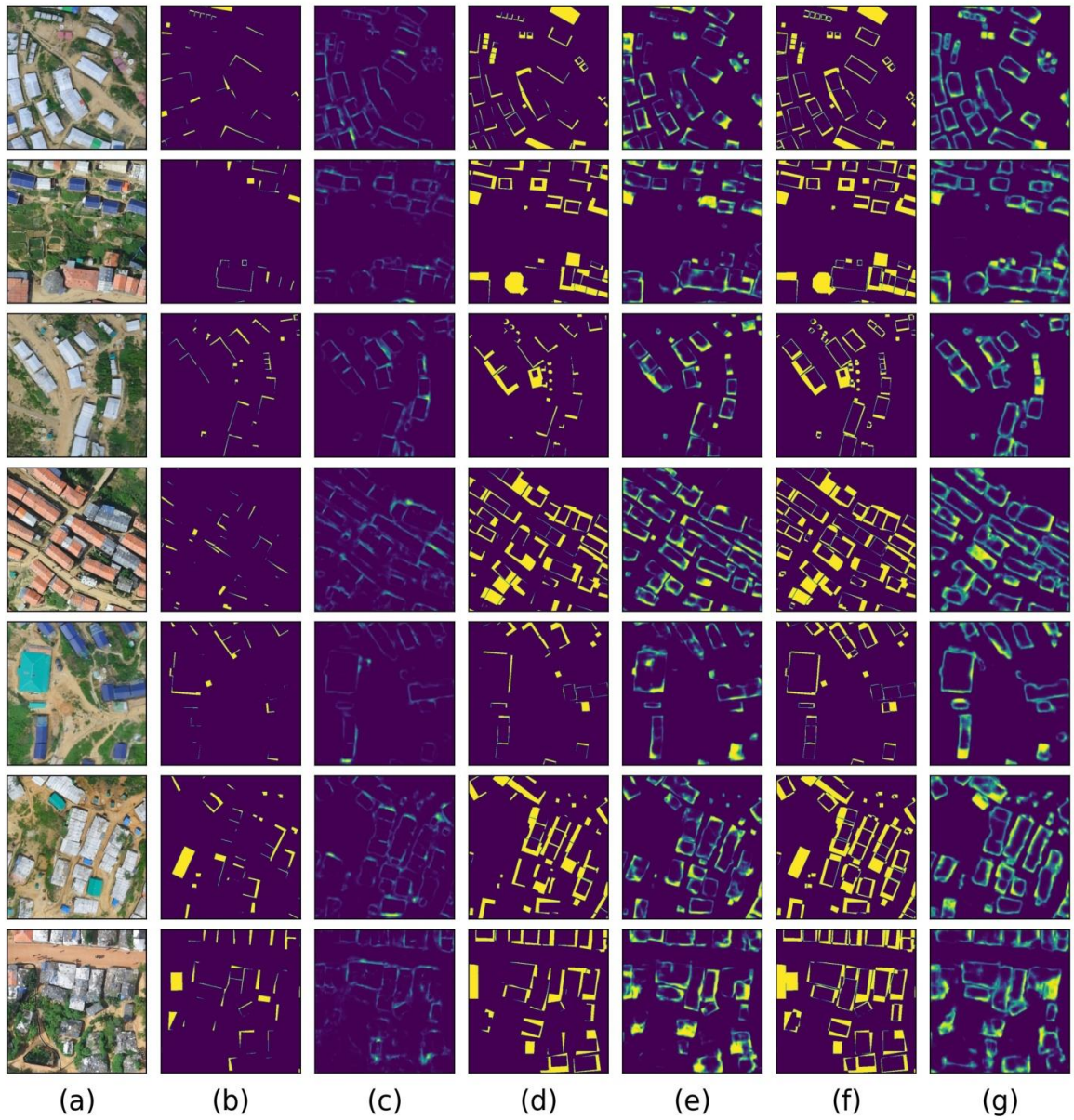
Some examples of error matrices predicted during the intermediate step are shown in Fig. 10. The error matrices are sparse, and weakly correlated to the images as the real world label noise can be random at times. However, they can provide insights about location having higher probabilities of being mislabeled. The ground truth FP error matrix is shown in Fig. 10(b) and the predicted FP error matrix is shown in Fig. 10(c). FP pixels are usually pixels adjacent to the clean building label boundary, but falling outside the boundary; this intuition is captured by the FP error matrix model as indicated by the predictions in Fig. 10(c). i.e. the regions adjacent to actual/clean boundaries have higher activations than other regions in the images, and thus have a higher probability of being an observed FP pixel. The predicted FP error matrix (non-thresholded) provides a heat map indicating the probability of each observed positive pixel label actually being true negative pixels.

FN pixels are less sparse than FP pixels since a major source of label noise in building extraction datasets comes from omitted/missed out buildings and shrunk label polygons. Fig. 10(d) shows the actual FN error matrix and Fig. 10(e) shows the predicted by the FN error matrix model. FN pixels are pixels within the clean building boundaries which are observed as non-building in the noisy labels, therefore regions in close proximity to the clean building boundaries but on the inner side have the highest probability of being observed as FN pixels, this is shown in Fig. 10(e). It is interesting to note that all pixels with significantly high FP error matrix activations lie outside and adjacent to the clean building boundaries whereas all pixels with significantly high FN error matrix activations lie inside the clean building boundaries; the modeling intuition

587    is expressed in the qualitative results.

588

589



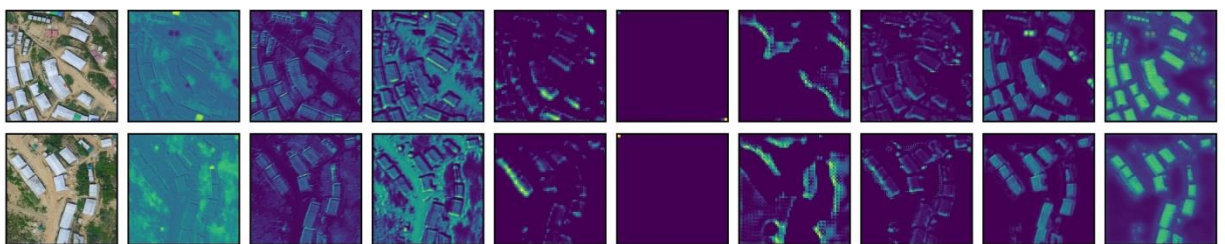(a)    (b)    (c)    (d)    (e)    (f)    (g)

590

591    **Figure 10.** Examples of error matrix predictions (a) Image (b) FP error matrix (c)

592    Predicted FP error matrix (d) FN error matrix (e) Predicted FN error matrix (f) General

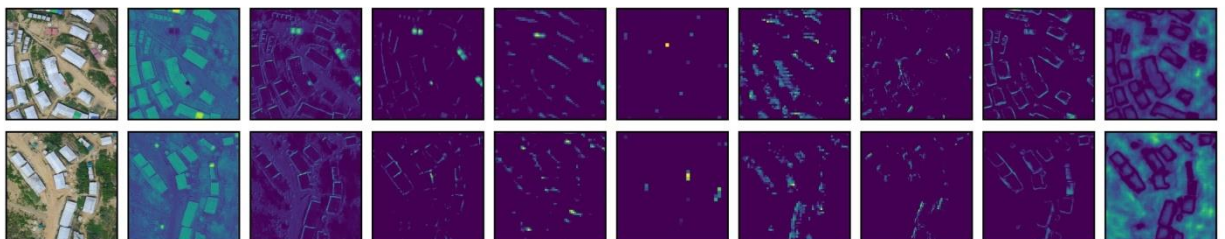593    error matrix (g) Predicted general error matrix

594    The general error matrix predictions are shown in Fig. 10(f) and the predicted

595    general error matrix predictions are shown in Fig. 10(g). Among the three types of error

596    matrices (general, FP and FN) the general error matrices are least sparse, since they are

597    the element wise addition of the FP and FN error matrices. The extra information

598    provided by separated FP and FN matrices are crucial to approximating useful noise

599    features. Experimental results on our dataset confirm this statement. The I-EM model

600    results are poorer than the CL model (albeit providing higher recall values at all

601    thresholds) qualitatively and quantitatively (in terms of F1-score and IoU score on the

602    independent test set). The predicted intermediate observed label also affects the

603    predicted true label. The outputs of hidden blocks of different models are shown in Fig.

604    11, feature maps for learning error matrices (Fig. 11(b), 11(c), 11(d) are quite different

605    from feature maps for learning base level building extraction (Fig. 11(a)). The

606    activation maps in Fig. 11 are outputs of the blocks for each model architecture. The

607    first feature map for each output is shown. The block outputs in Fig. 11 ($U_1$-$U_4$. the

608    bottleneck and $D_1$-$D_4$) show discriminative properties of the learned mappings in terms

609    of resolution and separability.
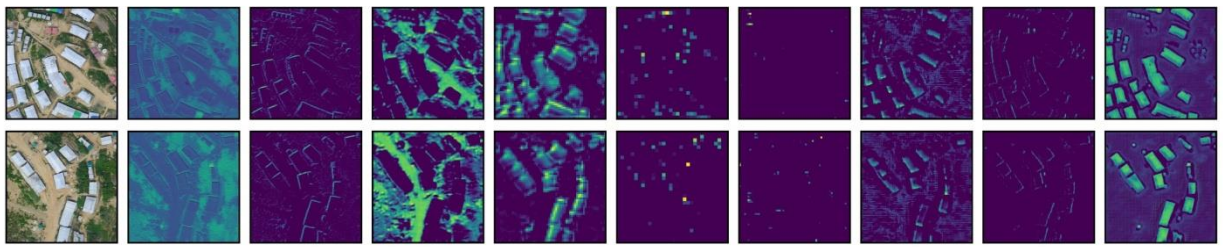
610



611
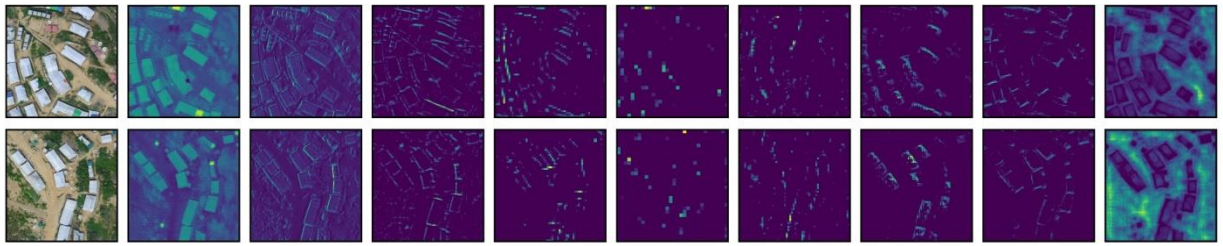612                                         (a)



613

614

(b)



615
616

(c)



| Image | $D_1$ | $D_2$ | $D_3$ | $D_4$ | Bottleneck | $U_1$ | $U_2$ | $U_3$ | $U_4$ |

617
618

(d)

619 Figure 11. Outputs learned by hidden convolutional layers for building extraction and

620 for error matrix approximation. (a) Clean label network (b) Error network (c) FP error

621 network (d) FN error matrix network

## Conclusion

623 In this work, we have provided a comprehensive taxonomy of label noise, in which the

624 six formulated label noise models can be used to express any kind of label noise in

625 building extraction tasks. Dense models are more apt than pixel-wise models for

626 building extraction. We propose two new model frameworks for dense prediction based

627 building extraction under label noise. The first model approximates the general error

628 matrix as an intermediate step, but has poor performance improvements compared to the

629 clean model. However, approximating the FP error matrix and the FN error matrix

630 separately greatly improves performance over the idealistic scenario presented in the

631 form of the CL model. Therefore, it is important to model the false positives and false

632 negatives independently rather than using a general model for both types of pixel-level

633 observed labels. Label noise in most building extraction cases is asymmetric, as also

634    observed for our case; there is a massive imbalance in the pixel-level label noise i.e.

635    there are much more false negatives than false positives. Therefore, a general model is

636    not sufficient in modeling the FP and FN noise processes to a degree that can aid the

637    larger task of noise-free building extraction. Qualitative results show that the error

638    matrix models (FP, FN and general) all capture the intuition behind the model

639    framework. The FP error matrix dense model has higher activations for regions right

640    outside and adjacent to the actual clean building boundaries. Similarly, FN error matrix

641    dense model has higher activations for regions inside and adjacent to the actual clean

642    building boundaries. Clean labels and corresponding observed labels with real-world

643    label noise are rarely available in conjunction with each other, which are essential for

644    obtaining the error matrices outlined in our proposed methodologies, and thus limit the

645    applicability.

649    **Disclosure statement**

650    No potential conflict of interest was reported by the authors.

651    **Data and Codes Availability Statement**

652    The data and codes that support the findings of this study are available at dedicated

653    GitHub repository (https://github.com/nahian-ahmed/dense-label-noise).

657    **References**

658    1.  Mnih, V., & Hinton, G. E. (2012). Learning to label aerial images from noisy
659        data. In Proceedings of the 29th International conference on machine learning
660        (ICML-12) (pp. 567-574).

661

662    2.  Ahmed, N., Mahbub, R. B., & Rahman, R. M. (2020). Learning to extract
663        buildings from ultra-high-resolution drone images and noisy labels. International
664        Journal of Remote Sensing, 1-22.

665    3.  Zhang, Z., Guo, W., Li, M., & Yu, W. (2020). GIS-supervised building
666        extraction with label noise-adaptive fully convolutional neural network. IEEE
667        Geoscience and Remote Sensing Letters.

668

669    4.  Frénay, B., & Kabán, A. (2014, April). A comprehensive introduction to label
670        noise. In ESANN.

671

672    5.  Frénay, B., & Verleysen, M. (2013). Classification in the presence of label
673        noise: a survey. IEEE transactions on neural networks and learning systems,
674        25(5), 845-869.

675

676    6.  Lu, Z., Fu, Z., Xiang, T., Han, P., Wang, L., & Gao, X. (2016). Learning from
677        weak and noisy labels for semantic segmentation. IEEE transactions on pattern
678        analysis and machine intelligence, 39(3), 486-500.

679

680    7.  Maas, A., Rottensteiner, F., & Heipke, C. (2016). Using label noise robust
681        logistic regression for automated updating of topographic geospatial databases.
682        In XXIII ISPRS Congress, Commission VII 3 (2016), Nr. 7 (Vol. 3, No. 7, pp.
683        133-140). Göttingen: Copernicus GmbH.

684

8. Maas, A. E., Rottensteiner, F., & Heipke, C. (2019). A label noise tolerant random forest for the classification of remote sensing data based on outdated maps for training. Computer Vision and Image Understanding, 188, 102782.

688

9. Ghosh, A., Kumar, H., & Sastry, P. S. (2017, February). Robust loss functions under label noise for deep neural networks. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 31, No. 1).

692

10. Rolnick, D., Veit, A., Belongie, S., & Shavit, N. (2017). Deep learning is robust to massive label noise. arXiv preprint arXiv:1705.10694.

695

11. Patrini, G., Rozza, A., Krishna Menon, A., Nock, R., & Qu, L. (2017). Making deep neural networks robust to label noise: A loss correction approach. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 1944-1952).

700

12. Jiang, L., Huang, D., Liu, M., & Yang, W. (2020, November). Beyond synthetic noise: Deep learning on controlled noisy labels. In International Conference on Machine Learning (pp. 4804-4815). PMLR.

704

13. Garcia, L. P., de Carvalho, A. C., & Lorena, A. C. (2015). Effect of label noise in the complexity of classification problems. Neurocomputing, 160, 108-119.

707

708 14. Pelletier, C., Valero, S., Inglada, J., Champion, N., Marais Sicre, C., & Dedieu,
709     G. (2017). Effect of training class label noise on classification performances for
710     land cover mapping with satellite image time series. Remote Sensing, 9(2), 173.

711

712 15. Frank, J., Rebbapragada, U., Bialas, J., Oommen, T., & Havens, T. C. (2017).
713     Effect of label noise on the machine-learned classification of earthquake
714     damage. Remote Sensing, 9(8), 803.

715

716 16. Angluin, D., & Laird, P. (1988). Learning from noisy examples. Machine
717     Learning, 2(4), 343-370.

718

719 17. Lawrence, N., & Schölkopf, B. (2001, July). Estimating a kernel fisher
720     discriminant in the presence of label noise. In 18th International Conference on
721     Machine Learning (ICML 2001) (pp. 306-306). Morgan Kaufmann.

722

723 18. Pérez, C. J., Girón, F. J., Martín, J., Ruiz, M., & Rojano, C. (2007).
724     Misclassified multinomial data: a Bayesian approach. RACSAM, 101(1), 71-80.

725

726 19. Xu, Y., Wu, L., Xie, Z., & Chen, Z. (2018). Building extraction in very high
727     resolution remote sensing imagery using deep learning and guided filters.
728     Remote Sensing, 10(1), 144.

729 20. Yuan, J. (2017). Learning building extraction in aerial scenes with convolutional
730     networks. IEEE transactions on pattern analysis and machine intelligence,
731     40(11), 2793-2798.

732

733    21. Chen, K., Fu, K., Gao, X., Yan, M., Sun, X., & Zhang, H. (2017, July). Building
734         extraction from remote sensing images with deep learning in a supervised
735         manner. In 2017 IEEE International Geoscience and Remote Sensing
736         Symposium (IGARSS) (pp. 1672-1675). IEEE.

737

738    22. Yang, H., Wu, P., Yao, X., Wu, Y., Wang, B., & Xu, Y. (2018). Building
739         extraction in very high resolution imagery by dense-attention networks. Remote
740         Sensing, 10(11), 1768.

741

742    23. Ji, S., Wei, S., & Lu, M. (2018). Fully convolutional networks for multisource
743         building extraction from an open aerial and satellite imagery data set. IEEE
744         Transactions on Geoscience and Remote Sensing, 57(1), 574-586.

745

746    24. Vakalopoulou, M., Karantzalos, K., Komodakis, N., & Paragios, N. (2015, July).
747         Building detection in very high resolution multispectral data with deep learning
748         features. In 2015 IEEE International Geoscience and Remote Sensing
749         Symposium (IGARSS) (pp. 1873-1876). IEEE.

750

751    25. Huang, Z., Cheng, G., Wang, H., Li, H., Shi, L., & Pan, C. (2016, July).
752         Building extraction from multi-source remote sensing images via deep
753         deconvolution neural networks. In 2016 IEEE International Geoscience and
754         Remote Sensing Symposium (IGARSS) (pp. 1835-1838). IEEE.

755

756    26. Shrestha, S., & Vanneschi, L. (2018). Improved fully convolutional network
757         with conditional random fields for building extraction. Remote Sensing, 10(7),
758         1135.

759

760   27. Boonpook, W., Tan, Y., & Xu, B. (2021). Deep learning-based multi-feature
761        semantic segmentation in building extraction from images of UAV
762        photogrammetry. International Journal of Remote Sensing, 42(1), 1-19.

763

764   28. Sun, S., Mu, L., Wang, L., Liu, P., Liu, X., & Zhang, Y. (2021). Semantic
765        Segmentation for Buildings of Large Intra-Class Variation in Remote Sensing
766        Images with O-GAN. Remote Sensing, 13(3), 475.

767

768   29. Wang, S., Hou, X., & Zhao, X. (2020). Automatic building extraction from
769        high-resolution aerial imagery via fully convolutional encoder-decoder network
770        with non-local block. IEEE Access, 8, 7313-7322.

771

772   30. Guo, M., Liu, H., Xu, Y., & Huang, Y. (2020). Building extraction based on U-
773        Net with an attention block and multiple losses. Remote Sensing, 12(9), 1400.

774

775   31. Shao, Z., Tang, P., Wang, Z., Saleem, N., Yam, S., & Sommai, C. (2020).
776        BRRNet: A fully convolutional neural network for automatic building extraction
777        from high-resolution remote sensing images. Remote Sensing, 12(6), 1050.
778

779   32. Jing, H., Sun, X., Wang, Z., Chen, K., Diao, W., & Fu, K. (2021). Fine Building
780        Segmentation in High-Resolution SAR Images via Selective Pyramid Dilated
781        Network. IEEE Journal of Selected Topics in Applied Earth Observations and
782        Remote Sensing.

783

784   33. He, Q., Sun, X., Yan, Z., & Fu, K. (2021). DABNet: Deformable contextual and
785        boundary-weighted network for cloud detection in remote sensing images. IEEE
786        Transactions on Geoscience and Remote Sensing.

787

788

789