

The Weak Supervision Landscape

Rafael Poyiadzi¹, Daniel Bacaicoa-Barber², Jesus Cid-Sueiro²,
Miquel Perello-Nieto¹, Peter Flach¹, Raul Santos-Rodriguez¹

¹*Intelligent Systems Lab, University of Bristol, Bristol, UK*

²*Signal Theory and Communications Dept., Universidad Carlos III of Madrid, Spain*

Contact Email: rp13102@bristol.ac.uk

Abstract—Many ways of annotating a dataset for machine learning classification tasks that go beyond the usual class labels exist in practice. These are of interest as they can simplify or facilitate the collection of annotations, while not greatly affecting the resulting machine learning model. Many of these fall under the umbrella term of weak labels or annotations. However, it is not always clear how different alternatives are related. In this paper we propose a framework for categorising weak supervision settings with the aim of: (1) helping the dataset owner or annotator navigate through the available options within weak supervision when prescribing an annotation process, and (2) describing existing annotations for a dataset to machine learning practitioners so that we allow them to understand the implications for the learning process. To this end, we identify the key elements that characterise weak supervision and devise a series of dimensions that categorise most of the existing approaches. We show how common settings in the literature fit within the framework and discuss its possible uses in practice.

Index Terms—weak supervision, weak labels, annotation process

I. INTRODUCTION

A machine learning classification task requires having a dataset of pairs of instances and labels. Obtaining labels of good quality can be expensive, time-consuming, or difficult in general. These constraints have led to the development and study of several flexible settings where annotations are assumed to not be perfect, but still suitable for the learning process. These are usually referred to as weak supervision (WS). For example, specialised products like Amazon’s Mechanical Turk provide access to pools of (non-expert) annotators whose labels are sometimes ambiguous and noisy. Interestingly, weak labels can take several forms and not only be the result of the annotation process (e.g., the data not annotated by an expert, or automatically extracted from the web) but they can come from the choices made by the dataset owner when deciding on the annotation process (e.g., allow annotators to provide more than one candidate classes when uncertain).

In machine learning, weakly supervised classification refers the task of obtaining a classifier from a given *weak* dataset, such that it has a low generalisation error with respect to the true data distribution. Even though weakly supervised learning encompasses settings which are widely applicable and studied, it is yet to become a standard machine learning setting

such as traditional supervised classification or clustering [1]. They have been previously referred to as indirect supervision [2], distant supervision [3], inaccurate, incomplete, or inexact supervision [4], learning from measurements [5], learning imprecise and fuzzy observations [6] and many more. In this work we refer to all the approaches where the observed label is not perfect as *weak*. Where applicable we will also refer to the observed label as *weak*, as opposed to the unobserved *clean* label. We will also be referring to the process by which a clean label is transformed to a weak label as the *weakening process*.

In this paper we introduce a framework for categorising weak supervision settings by identifying the key set of dimensions that should be used to describe the different instantiations that exist in the literature and in practice. This framework can then be used by both dataset owners / annotators and machine learning practitioners / researchers. For the former, it will be a tool to navigate the landscape of options when designing the data collection or annotation process or to describe an existing dataset. In both cases it will aid communication and can help future users understand implications of the type of annotations present in the dataset. For the later, it can be used to find the right place for a new or existing technique to learn from weak labels and identify (open) research problems in the field and have a clearer understanding of the generalisation of their contributions.

The remainder of the paper is structured as follows: in Section II we introduce the framework and in Section III we present several well studied WSL settings and discuss how they fit within our framework. Lastly in Section IV we discuss the main implications and limitations of the work and directions for future research.

II. DIMENSIONS OF WEAK SUPERVISION

In this section we present the dimensions of weak supervision, which are the building blocks of our framework. We separate the dimensions into three groups based on whether they refer to the true label space, the weak label space or the weakening process. Table I condenses all this information together with questions to illustrate the meaning of each dimension and the options that they offer. This table can be used as a standalone tool to plan or understand weak annotations.

Partially supported by project TEC2017-83838-R, funded by FEDER/Ministerio de Ciencia, Innovación y Universidades – AEI; by the SPHERE Next Steps Project funded by EPSRC [grant EP/R005273/1]; RP and RSR are funded by the UKRI Turing AI Fellowship EP/V024817/1.

TABLE I

FRAMEWORK FOR CATEGORISING WEAK SUPERVISION SETTINGS CONSTITUTED BY CATEGORIES AND DIMENSIONS. THE THIRD COLUMN LISTS OPTIONS FOR EACH. INDIVIDUAL SETTINGS CAN BE IDENTIFIED BY SELECTING ONE OPTION PER DIMENSION. WE NOTE THAT CERTAIN COMBINATIONS OF OPTIONS ARE NOT COMPATIBLE. THE FINAL COLUMN SUMMARIZES THE MEANING OF EACH OF THE DIMENSIONS FOR USE IN THE PRESCRIPTION OR DESCRIPTION OF THE CORRESPONDING DIMENSION.

Category	Dimension	Options	Question
True Label Space	Number of classes	Binary / Multi-class	How many classes does the task involve?
	Multi-label	Yes / No	Can instances belong to more than one class?
Weak Label space	Unsupervised	Yes / No	Are annotators allowed to not annotate certain samples?
	Soft labels	Yes / No	Are annotators allowed to use soft or probabilistic labels?
	Number of Annotators	1 / >1	How many annotators will annotate the data?
	Number of candidate classes	1 / >1	Are annotators allowed to provide annotations covering more than one class?
Weakening Process	Aggregation	Yes / No	Are samples annotated individually or as a group?
	Class dependent	Yes / No	Are classes equally prone to annotation errors?
	Instance dependent	Yes / No	Are samples equally prone to annotation errors?

A. True Label Space

Here we present dimensions that derive directly from the nature and description of the task. Classification tasks can take different forms, from binary classification, where we might be interested in identifying an image as a dog or a cat, to multi-class, where we might aim to recognize images of dogs according to their breed. This is summarized in the first dimension – **number of classes**. Additionally, we accept that it is sometimes the case that more than one categories can be assigned to a single instance as in **multi-label** settings, which forms the second dimension. The types of true/clean label that we consider are then as follows.

Number of classes:

$$\mathcal{Y}_k = \{\mathbf{y} \mid \mathbf{y} \in \{0, 1\}^k, \mathbf{1}^\top \mathbf{y} = 1\} \quad (1)$$

Multi-label:

$$\mathcal{Y}_{m,k} = \{\mathbf{y} \mid \mathbf{y} \in \{0, 1\}^k, 1 \leq \mathbf{1}^\top \mathbf{y} \leq m \leq k\} \quad (2)$$

Although this can be extended to structured data, we do not explore it here for simplicity. We also assume that the dataset only contains what is defined by the task, e.g., in the case of the task being classifying dogs vs cats, there would not be images of other animals in the dataset. Even though we do not discuss the learning stage in this paper, the framework is constructed keeping in mind that part of the task would be to obtain a classifier: $f: \mathcal{X} \rightarrow \mathcal{Y}_{\text{Clean}}$ and in minimizing $\mathbb{E}_{X,Y} \ell(Y, f(X))$, where $\ell(\cdot, \cdot)$ is a loss function and the expectation is over the clean label data distribution.

B. Weak Label Space

In this section we focus on dimensions concerned with the weak label space, $\mathcal{Y}_{\text{Weak}}$ (See Eq. 5) and the possible forms that it might take. These characterise the degrees of freedom of the annotator.

Access to unlabelled data. In certain situations, besides the potentially weakly labelled dataset, we also have access to

a separate unlabelled dataset. This could be either because annotators are allowed to return an empty annotation, or because this set of data was just not chosen for annotation. We refer to this dimension as **unsupervised**.

Access to multiple annotators. We usually assume a dataset is annotated by one annotator, but it might be the case that a **number of annotators** provide annotations for the same dataset. Annotations on the same instance will not necessarily agree, potentially creating ambiguity.

Restriction on number of assigned candidate classes. In binary and multi-class classification, the classes are mutually exclusive and exhaustive, which means that every instance is associated with one true class only. In the weak label space, we could allow an annotator to provide a set of candidate classes, instead of just one (**number of candidate classes**).

An example of a weak label and how this set of dimensions can affect what we get to observe is shown in Figure 1. It shows how these three dimensions can increase the complexity of a weak label, as we can have several annotators, allow them to provide no annotation, and also allow them to assign more than one class to each sample.

Soft labels. Also known as probabilistic labels, the **soft labels** relaxation allows for more flexibility in the annotation process by letting an annotator express a degree of belief. For example, instead of resorting to ‘dog’ for an image, they could say “70% confident this image contains a dog” [7].

C. Weakening Process

The final set of dimensions captures the different aspects of the weakening process, i.e., the (usually unknown) transformation that maps true/clean to weak labels. Interestingly, the weakening process depends on the annotator and how they make suboptimal annotation decisions, but it can also depend on choices made by the dataset owner or on the task itself.

Aggregation. A key dimension is whether the labels provided correspond to a single instance, or whether they depend on,

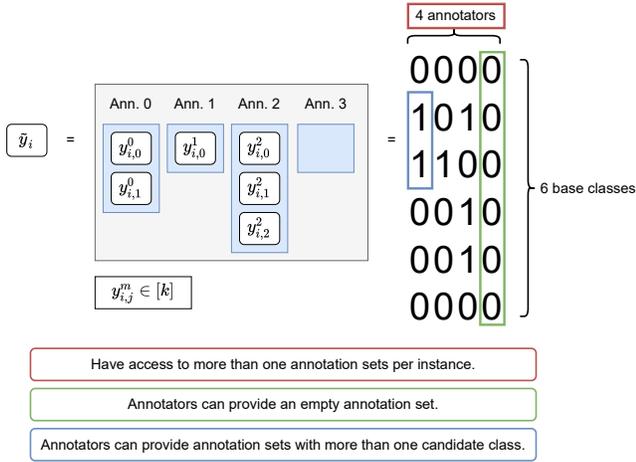


Fig. 1. An example of a non-aggregate weak label. For this case we consider six base classes and access to four annotators. This example showcases three dimensions: (1) having access to more than one annotators, (2) annotators can provide an empty annotation set (see annotator 3), and (3) annotators can provide an annotation set with more than one candidate class (see annotators 1 and 2).

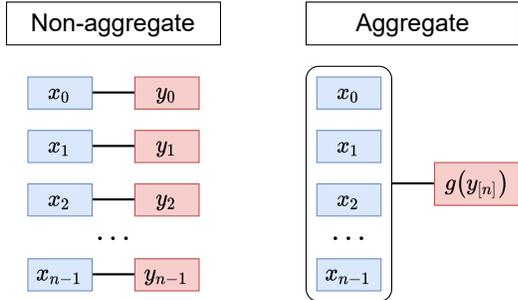


Fig. 2. Depiction of the difference between non-aggregate and aggregate settings. In the case of non-aggregate settings, there is a one-to-one relationship between instances and labels, while in the case of aggregate noise settings we observe one label per group. We note that a label could be one class, multiple or none, as we describe below.

and correspond to, a set of instances through an **aggregation** mechanism. For the aggregate cases, the labels provided correspond to a set of instances which we will refer to as a ‘bag’. Instead of being provided with a set of instance-label pairs $\{(x_i, y_i)\}_{i=1}^n$, we are provided with a set of $\{(x_i, b_i)\}_{i=1}^n$ where b_i ’s are bag indicators and $\{(b_j, t_j)\}_{j=1}^m$, where t_j ’s are the corresponding labels. The labels in this case are of the form:

$$t_j = g\left(\{y_i \mid i \in b_j\}\right) \quad (3)$$

where we use $i \in b_j$ to imply that the i^{th} instance belongs to the j^{th} bag, and g is the label aggregation function. It should be noted that t_j is the label for *all* instances x_i , $i \in b_j$. The distinction between the two is shown pictorially in Fig. 2.

For non-aggregate settings, we will refer to $\mathbb{P}(\tilde{Y} \mid X, Y)$, where \tilde{Y} denotes the random variable for the weak label, as the *weakening function*, which upon conditioning has the form,

$$\tau(y; x) = \mathbb{P}(\tilde{Y} \mid X = x, Y = y) \quad (4)$$

with

$$\tau(y; x) : \mathcal{Y}_{\text{Clean}} \rightarrow \Delta_{\mathcal{Y}_{\text{Weak}}} \quad (5)$$

where $\Delta_{\mathcal{Y}_{\text{Weak}}}$ is the probability simplex over $\mathcal{Y}_{\text{Weak}}$. Eqs. 4 & 5 involve three objects:

- the clean label space $\mathcal{Y}_{\text{Clean}}$ (covered in Sec. II-A),
- the weak label space $\mathcal{Y}_{\text{Weak}}$ (covered in Sec. II-B) and
- the weakening function $\tau(y; x)$.

The weakening function acts between discrete sets and can therefore be described as a *Categorical* (Cat) distribution. This distribution is parametrized by a column of a *mixing matrix* or transition matrix, \mathbf{T} , which is a non-negative column-wise stochastic matrix. The weakening process is modelled as,

$$\tilde{Y} \mid Y = y, X = x \sim \text{Cat}(\mathbf{T}\mathbf{y}) \quad (6)$$

With regards to dimensions arising of the weakening function we will consider two, the dependence on the instance and on the true class.

Dependence on instance. We consider two cases, **instance dependence** or instance independence [8].

Instance-independent noise (IIN) :

$$\mathbb{P}(\tilde{Y} = \tilde{y} \mid X = x, Y = y) = \mathbb{P}(\tilde{Y} = \tilde{y} \mid Y = y),$$

Instance-dependent noise (IDN) :

$$\mathbb{P}(\tilde{Y} = \tilde{y} \mid X = x, Y = y) \neq \mathbb{P}(\tilde{Y} = \tilde{y} \mid Y = y).$$

Dependence on true class. With regards to the **class dependence** we consider *symmetric* (uniform) or *asymmetric* (class-conditional) with respect to the original classes [9]. In the case of multi-class classification, symmetric noise would imply:

$$\mathbb{P}(\tilde{Y} = e_u \mid X = x, Y = e_i) = \mathbb{P}(\tilde{Y} = e_v \mid X = x, Y = e_i) \quad (7)$$

$$\forall u, v, i \in [k], u \neq v \neq i.$$

We have asymmetric label noise when this does not hold.

III. THE DIMENSIONS IN PRACTICE

We now present several well-studied settings within the field of weak supervision and discuss how they fit in the proposed framework. We first start by exploring examples that belong to the non-aggregate category and then move to aggregate, as this separation simplifies the formulation of the settings.

A. Non-aggregate WS Settings

For non-aggregate WS settings, the selected examples are summarized in Table II to illustrate their constraints in the true label space, weak label space and weakening process.

a) *Noisy labels (Flipping noise)*: In this setting, instances switch labels with a certain probability. This type of noise could be introduced when the data is labelled by a non-expert annotator, and nicely extends to having multiple such annotators. In this setting the clean and noisy label spaces are the same:

$$\begin{aligned} \mathcal{Y}_k &\rightarrow \mathcal{Y}_k \\ \mathbf{T} &\in \mathbb{R}^{k \times k} \end{aligned}$$

TABLE II

NON-AGGREGATE WSL SETTINGS ACCORDING TO THE LABEL-SPACE OF TRUE LABELS AND WEAK LABELS, AND THE MIXING MATRIX USED TO MODEL THE WEAKENING PROCESS.

Name	True Label Space	Weak Label Space	Mixing matrix
Noisy Labels	\mathcal{Y}_k	\mathcal{Y}_k	$\mathbf{T} \in \mathbb{R}^{k \times k}$
Partial Labels	\mathcal{Y}_k	$\mathcal{Y}_{m,k}$	$\mathbf{T} \in \mathbb{R}^{(2^k-2) \times k}$
Superset Learning	\mathcal{Y}_k	$\mathcal{Y}_{m,k}(z)$	$\mathbf{T} \in \mathbb{R}^{(2^k-2) \times k}$
Semi-supervised Learning	\mathcal{Y}_k	\mathcal{Y}_{k+1}	$\mathbf{T} \in \mathbb{R}^{(k+1) \times k}$
Positive-Unlabelled	\mathcal{Y}_2	\mathcal{Y}_3	$\mathbf{T} \in \mathbb{R}^{3 \times 2}$
Multiple Annotators	\mathcal{Y}_k	\mathcal{Y}_k^n	$\{\mathbf{T} \in \mathbb{R}^{k \times k}\}_{i=1}^n$

An example of a mixing matrix for binary classification:

$$+ \quad - \\ + \begin{pmatrix} 1 - \gamma_0 & \gamma_1 \\ \gamma_0 & 1 - \gamma_1 \end{pmatrix} \\ -$$

where if $\gamma_0 = \gamma_1$ it implies symmetric label noise, and asymmetric label noise otherwise. In the case of instance-independent noise, γ_0 and γ_1 are constants, i.e., all instances have their label flipped with same probability. On the other hand, with instance-dependent noise, γ_0 and γ_1 could be functions of x , where we could have rates differ $\gamma_0(x_i) \neq \gamma_1(x_j)$ for $x_i \neq x_j$. A comprehensive review of classification in the presence of label noise is given by [10].

b) Superset Learning: In this setting a weak label could be any combination of class labels, subject to the constraint that this combination contains the true label. For example, in the case of multi-class classification with three classes, for a true label [001] (class 3 in one-hot encoding), we could observe [011] or [101], but not [110] (represented as the binary OR operator on the true one-hot encoding). Superset learning has been studied in the literature under different names such as ‘learning with partial-labels’ (see below), ‘learning with ambiguous labels’ and ‘learning from complementary labels’ [11]. Some of the earlier works on the topic include [12] and [13] where the setting is referred to as ‘partial-labels’ and ‘multiple labels’ respectively. An example of a mixing matrix:

$$\begin{matrix} & 001 & 010 & 100 \\ \begin{matrix} 001 \\ 010 \\ 100 \\ 110 \\ 101 \\ 011 \end{matrix} & \begin{pmatrix} \alpha_0 & 0 & 0 \\ 0 & \beta_0 & 0 \\ 0 & 0 & \gamma_0 \\ 0 & \beta_1 & \gamma_1 \\ \alpha_1 & 0 & 1 - \gamma_0 - \gamma_1 \\ 1 - \alpha_0 - \alpha_1 & 1 - \beta_0 - \beta_1 & 0 \end{pmatrix} \end{matrix}$$

$$\mathcal{Y}_k \rightarrow \mathcal{Y}_{m,k}(z) \\ \mathbf{T} \in \mathbb{R}^{(2^k-2) \times k}$$

where,

$$\mathcal{Y}_{m,k}(z) = \{\mathbf{y} | \mathbf{y} \in \{0, 1\}^k, 1 \leq \mathbf{1}^\top \mathbf{y} \leq m \leq k, z \in \mathbf{y}\} \quad (8)$$

We abuse notation and use \mathbf{y} both as a vector and as a set and with $z \in \mathbf{y}$ we imply that \mathbf{y} covers z , i.e. has a non-zero entry wherever z has a non-zero entry.

c) Partial Labels (PLL): PLL is sometimes used to refer to superset learning. This setting is similar to that of superset learning but where there is no restriction that the observed weak label includes the true label [14, 15].

$$\mathcal{Y}_k \rightarrow \mathcal{Y}_{m,k} \\ \mathbf{T} \in \mathbb{R}^{(2^k-2) \times k}$$

where,

$$\mathcal{Y}_{m,k} = \{\mathbf{y} | \mathbf{y} \in \{0, 1\}^k, 1 \leq \mathbf{1}^\top \mathbf{y} \leq m \leq k\} \quad (9)$$

In both settings, the set of potential observations extends from k to $2^k - 2$. While standard presentations of the settings consider $2^k - 1$, we choose to exclude the all inclusive potential observation and instead consider it as an extra dimension.

d) Semi-supervised Learning (SSL): In semi-supervised learning [16], on top of the usual supervised dataset, we are also provided with an unlabelled dataset.

$$\mathcal{Y}_k \rightarrow \mathcal{Y}_{k+} \\ \mathbf{T} \in \mathbb{R}^{(k+1) \times k}$$

where,

$$\mathcal{Y}_{k+} = \{\mathbf{y} | \mathbf{y} \in \{0, 1\}^k, \mathbf{1}^\top \mathbf{y} \in \{0, 1\}\} \quad (10)$$

and is an extension of the multi-class label space (Eq. 1) that allows for no classes to be provided. An example of a mixing matrix for binary classification:

$$+ \quad - \\ + \begin{pmatrix} 1 - \gamma_0 & 0 \\ 0 & 1 - \gamma_1 \end{pmatrix} \\ - \quad \text{na} \begin{pmatrix} \gamma_0 & \gamma_1 \end{pmatrix} \quad (11)$$

e) Positive-Unlabelled (PU) Learning: Learning with positive and unlabelled instances [17] is the setting of binary classification where the training dataset only consists of positive and unlabelled instances. Situations where PU learning arises include medical records where only known previous diseases are listed and personalised advertising where visited pages and clicks are the positive cases [18].

$$\mathcal{Y}_2 \rightarrow \mathcal{Y}_{2+} \\ \mathbf{T} \in \mathbb{R}^{3 \times 2}$$

An example of a mixing matrix for PU learning:

$$+ \quad - \\ + \begin{pmatrix} 1 - \gamma_0 & 0 \\ 0 & 0 \end{pmatrix} \\ - \quad \text{na} \begin{pmatrix} \gamma_0 & 1 \end{pmatrix}$$

PU learning can be seen as a special case of semi-supervised learning where $\gamma_1 = 1$ (Eq. 11). It has been extended to the multi-class case, under the name Multi-Positive and Unlabeled learning, where labeled data from multiple positive classes are provided for training along with unlabeled data from a mixture of the positive classes and a negative class [19].

f) *Multiple annotators*: Assuming we have m annotators, we also have potentially m distinct mixing matrices [20, 21, 22]. Therefore, $\tilde{Y}^n = \{\tilde{Y} \sim \text{Cat}(\mathbf{T}_j \mathbf{y})\}_{j=1}^n$, where \mathbf{T}_j denotes the j^{th} annotator's mixing matrix.

$$\mathcal{Y}_k \rightarrow \mathcal{Y}_k^n \\ \left\{ \mathbf{T} \in \mathbb{R}^{k \times k} \right\}_{j=1}^n$$

where,

$$\mathcal{Y}_k^n = \{ \mathbf{y} \mid \mathbf{y} \in \{0, 1\}^{k \times n}, \mathbf{y}^\top \mathbf{1} = \mathbf{1} \} \quad (12)$$

where the equality should be understood as elementwise.

B. Aggregate WSL Settings

Aggregate settings have comparatively received much less attention in the literature. Here, we present multiple instance learning and learning from label proportions as the main representatives.

a) *Multiple Instance Learning*: Multiple instance learning (MIL) is usually considered in the binary classification case. The label provided for a bag of samples is an indicator of the presence of the positive class. In other words, is there at least one positive instance in the set? It was first introduced in [23] with the motivation of drug activity prediction.

$$g(y_1, \dots, y_n) = \max(y_{[1..n]})$$

In [24] the authors extend MIL, from requiring at least one positive instance in a bag, to requiring r .

$$g(y_1, \dots, y_m) = \mathbf{1} \left\{ \sum_{i=1}^m y_i \geq r \right\}$$

b) *Learning from Label Proportions (LLP)*: In LLP, the aggregation function is the count function (or proportion) for each of the classes present in a bag of samples. LLP was introduced in [25] with the motivation of learning with mass spectrometry data.

$$g(y_1, \dots, y_m) = \left[\sum_{i=1}^m y_{i,0}, \dots, \sum_{i=1}^m y_{i,c} \right]$$

Interestingly, PU Learning was presented as being a non-aggregate setting, but under certain conditions it can also be seen as a case of learning from label proportions. In PU Learning, we are provided with a set of positive data and a set of unlabelled data. If we know the portion of positives and negatives in the unlabelled set, we can view this as an instance of LLP with two bags.

C. Towards a Unified Formulation

For non-aggregate settings, we can now see how to find a common formulation for weakly supervised settings. Starting with the label space corresponding to multiple annotators in Eq. 12, we can extend it to allow annotators to provide an empty set for a sample (e.g., when they find an instance difficult to annotate):

$$\mathcal{Y}_{k+}^n = \{ \mathbf{y} \mid \mathbf{y} \in \{0, 1\}^{k \times n}, \mathbf{y}^\top \mathbf{1} \in \{0, 1\} \} \quad (13)$$

and then it could be extended to allow annotators to provide an annotation set with more than one candidate,

$$\mathcal{Y}_{m,k+}^n = \{ \mathbf{y} \mid \mathbf{y} \in \{0, 1\}^{k \times n}, \mathbf{0} \leq \mathbf{y}^\top \mathbf{1} \leq \mathbf{m} \leq \mathbf{k} \} \quad (14)$$

With regards to describing the annotation process for a dataset, we could model the weakening process through its dependence on the instance and the true class.

Aggregate settings can also be extended to the case of having access to unsupervised data or multiple annotators. Instance-dependence has a different meaning in this case though. While previously it had to do with whether the weakening function was uniform across all instances, in this case it has to do with bag creation and intra-bag similarities [26, 27, 28]

IV. CONCLUSION & FUTURE WORK

We have presented a framework with dimensions that can help in navigating the weak supervision field, but that can also help in understanding exploiting the flexibility of the annotation process.

However, this is nothing but a first step towards categorizing weak supervision. In future work we wish to extend this work in three directions. First, we aim to complement our framework with corresponding algorithms where a practitioner can turn to after identifying the characteristics of their problem. This would also allow for understanding the implications that certain choices on the annotation process have on available algorithms, their theoretical guarantees and practical considerations. Second, we want to expand this work into a transparent process for documenting the annotation of a dataset. This very much aligns with the proposal of accompanying a dataset with a datasheet that documents its motivation, composition, collection process and recommended uses in *Datasheets for Datasets* [29]. Third, we want to strengthen the framework itself. Non-aggregate weak supervision settings have been more widely studied and hence the difference in weight they have received in this paper. As seen, they can be unified through Eq. 6 and the mixing matrix which can be used to reverse the noise process and make learning unbiased [30]. An important aspect in these settings is whether the mixing matrix is known [31] or whether it has to be estimated [22]. Also, the unified formulation is not only a matter of aesthetics, but more importantly allows for the transferability of methods and theory. In the case of aggregate WSL settings, while we have Eq. 3 showing how aggregation is abstracted away, it does not improve our theoretical understanding, or allow for algorithms to be applied across settings yet.

REFERENCES

- [1] G. Patrini, “Weakly supervised learning via statistical sufficiency (phd thesis),” Ph.D. dissertation, The Australian National University, 2016.
- [2] A. Raghunathan, R. Frostig, J. Duchi, and P. Liang, “Estimation from indirect supervision with linear moments,” in *International conference on machine learning*. PMLR, 2016, pp. 2568–2577.
- [3] M. Surdeanu, J. Tibshirani, R. Nallapati, and C. D. Manning, “Multi-instance multi-label learning for relation extraction,” in *joint conference on empirical methods in natural language processing and computational natural language learning*, 2012, pp. 455–465.
- [4] Z.-H. Zhou, “A brief introduction to weakly supervised learning,” *National science review*, vol. 5, no. 1, pp. 44–53, 2018.
- [5] P. Liang, M. I. Jordan, and D. Klein, “Learning from measurements in exponential families,” in *International Conference on Machine Learning*, 2009, pp. 641–648.
- [6] E. Hüllermeier, “Learning from imprecise and fuzzy observations: Data disambiguation through generalized loss minimization,” *International Journal of Approximate Reasoning*, vol. 55, no. 7, pp. 1519–1534, 2014.
- [7] P. Peng, R. C.-W. Wong, and P. S. Yu, “Learning on probabilistic labels,” in *Proceedings of the SIAM International Conference on Data Mining*, 2014, pp. 307–315.
- [8] A. K. Menon, B. Van Rooyen, and N. Natarajan, “Learning from binary labels with instance-dependent noise,” *Machine Learning*, vol. 107, no. 8, pp. 1561–1595, 2018.
- [9] —, “Learning from binary labels with instance-dependent corruption,” *arXiv preprint arXiv:1605.00751*, 2016.
- [10] B. Frénay and M. Verleysen, “Classification in the presence of label noise: a survey,” *IEEE transactions on neural networks and learning systems*, vol. 25, no. 5, pp. 845–869, 2013.
- [11] T. Ishida, G. Niu, W. Hu, and M. Sugiyama, “Learning from complementary labels,” *arXiv preprint arXiv:1705.07541*, 2017.
- [12] Y. Grandvalet, “Logistic regression for partial labels,” in *Proc. IPMU*, 2002.
- [13] R. Jin and Z. Ghahramani, “Learning with multiple labels,” in *NIPS*, vol. 2. Citeseer, 2002, pp. 897–904.
- [14] J. Cid-Sueiro, “Proper losses for learning from partial labels,” in *Advances in neural information processing systems*. Citeseer, 2012, pp. 1565–1573.
- [15] J. Cid-Sueiro, D. Garcia-Garcia, and R. Santos-Rodríguez, “Consistency of losses for learning from weak labels,” in *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases*, 2014, pp. 197–210.
- [16] O. Chapelle, B. Scholkopf, and A. Zien, “Semi-supervised learning (chapelle, o. et al., eds.; 2006)[book reviews],” *IEEE Transactions on Neural Networks*, vol. 20, no. 3, pp. 542–542, 2009.
- [17] B. Liu, Y. Dai, X. Li, W. S. Lee, and P. S. Yu, “Building text classifiers using positive and unlabeled examples,” in *Third IEEE International Conference on Data Mining*, 2003, pp. 179–186.
- [18] J. Bekker and J. Davis, “Learning from positive and unlabeled data: A survey,” *Machine Learning*, vol. 109, no. 4, pp. 719–760, 2020.
- [19] Y. Xu, C. Xu, C. Xu, and D. Tao, “Multi-positive and unlabeled learning,” in *IJCAI*, 2017, pp. 3182–3188.
- [20] V. C. Raykar, S. Yu, L. H. Zhao, G. H. Valadez, C. Florin, L. Bogoni, and L. Moy, “Learning from crowds,” *Journal of Machine Learning Research*, vol. 11, pp. 1297–1322, Aug. 2010.
- [21] Y. Ni, M. McVicar, R. Santos-Rodríguez, and T. De Bie, “Understanding effects of subjectivity in measuring chord estimation accuracy,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 12, pp. 2607–2615, 2013.
- [22] M. Perello-Nieto, R. Santos-Rodríguez, D. Garcia-Garcia, and J. Cid-Sueiro, “Recycling weak labels for multiclass classification,” *Neurocomputing*, vol. 400, pp. 206–215, 2020.
- [23] T. G. Dietterich, H. R. Lathrop, and T. Lozano-Perez, “Solving the multiple instance problem with axis-parallel rectangles,” *Artificial intelligence*, vol. 89, no. 1-2, pp. 31–71, 1997.
- [24] S. Scott, J. Zhang, and J. Brown, “On generalized multiple-instance learning,” *International Journal of Computational Intelligence and Applications*, vol. 5, no. 01, pp. 21–35, 2005.
- [25] D. R. Musicant, J. M. Christensen, and J. F. Olson, “Supervised learning by training on aggregate outputs,” in *Seventh IEEE International Conference on Data Mining (ICDM 2007)*. IEEE, 2007, pp. 252–261.
- [26] Z.-H. Zhou, Y.-Y. Sun, and Y.-F. Li, “Multi-instance learning by treating instances as non-iid samples,” in *International Conference on Machine Learning*, 2009.
- [27] M.-A. Carbonneau, V. Cheplygina, E. Granger, and G. Gagnon, “Multiple instance learning: A survey of problem characteristics and applications,” *Pattern Recognition*, vol. 77, pp. 329–353, 2018.
- [28] C. Scott and J. Zhang, “Learning from label proportions: A mutual contamination framework,” *arXiv preprint arXiv:2006.07330*, 2020.
- [29] T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. D. Iii, and K. Crawford, “Datasheets for datasets,” *Communications of the ACM*, vol. 64, no. 12, pp. 86–92, 2021.
- [30] B. Van Rooyen, A. K. Menon, and R. C. Williamson, “Learning with symmetric label noise: The importance of being unhinged,” *arXiv preprint arXiv:1505.07634*, 2015.
- [31] D. Bacaicoa-Barber, M. Perelló-Nieto, R. Santos-Rodríguez, and J. Cid-Sueiro, “On the selection of loss functions under known weak label models,” in *International Conference on Artificial Neural Networks*, vol. 12892, 2021, pp. 332–343.