

Data Collection and Quality Challenges in Deep Learning: A Data-Centric AI Perspective

Steven Euijong Whang · Yuji Roh · Hwanjun Song · Jae-Gil Lee

Received: date / Accepted: date

Abstract Software 2.0 is a fundamental shift in software engineering where machine learning becomes the new software, powered by big data and computing infrastructure. As a result, software engineering needs to be re-thought where data becomes a first-class citizen on par with code. One striking observation is that 80–90% of the machine learning process is spent on data preparation. Without good data, even the best machine learning algorithms cannot perform well. As a result, data-centric AI practices are now becoming mainstream. Unfortunately, many datasets in the real world are small, dirty, biased, and even poisoned. In this survey, we study the research landscape for data collection and data quality primarily for deep learning applications. Data collection is important because there is lesser need for feature engineering for recent deep learning approaches, but instead more need for large amounts of data. For data quality, we study data validation and data cleaning techniques. Even if the data cannot be fully cleaned, we can still cope with imperfect data during model training where using robust model training techniques. In addition, while bias and fairness

have been less studied in traditional data management research, these issues become essential topics in modern machine learning applications. We thus study fairness measures and unfairness mitigation techniques that can be applied before, during, or after model training. We believe that the data management community is well poised to solve problems in these directions.

Keywords Data Collection · Data Quality · Deep Learning · Data-Centric AI

1 Overview

Deep learning is widely used to glean knowledge from massive amounts of data. There is a wide range of applications including natural language understanding, healthcare, self-driving cars, and more. Deep learning has become so prevalent thanks to its excellent performance with the availability of big data and powerful computing infrastructure. According to the IDC [3], the amount of data worldwide is projected to grow exponentially to 175 zettabytes (ZB) by 2025. In addition, powerful GPUs and TPUs enable software to have superhuman performances in various tasks.

Software 2.0 [10] is a fundamental paradigm shift in software engineering where machine learning becomes the new software. Conventional software engineering involves designing, implementing, and debugging code. In comparison, machine learning starts with data and trains a function on the data. It is known that data preparation is an expensive step in end-to-end machine learning. In particular, collecting data, cleaning it, and making it suitable for machine learning training takes 80–90% of the entire time [125]. In addition, the code on a machine learning platform (e.g., PyTorch [95]) is high level and thus requires significantly fewer lines

This article extends tutorials the authors delivered at the VLDB 2020 [134] and KDD 2021 [82] conferences.

S. E. Whang
KAIST
E-mail: swhang@kaist.ac.kr

Y. Roh
KAIST
E-mail: yuji.roh@kaist.ac.kr

H. Song
Naver AI Lab
E-mail: hwanjun.song@navercorp.com

J. Lee
KAIST
E-mail: jaegil@kaist.ac.kr

compared to conventional software. Finally, the trained model may need to be continuously improved with hyperparameter tuning. This entire process from data preparation to model deployment is widely viewed as a new software engineering paradigm, and companies have been actively developing open source [19, 34] and proprietary Software 2.0 systems [27, 13, 14, 60, 63, 114]. Solving data issues is increasingly becoming critical in machine learning research.

While data collection and quality issues are important, machine learning research has mainly focused on training algorithms instead of the data. In the industry, a common complaint [125] is that research institutions spend 90% of their machine learning efforts on algorithms and 10% on data preparation, although based on the amounts of time spent, the numbers should be 10% and 90% the other way.

At the same time, many companies are declaring Responsible and data-centric AI practices. For example, Google [8] says that AI has a significant potential to help solve challenging problems, but it is important to develop responsibly. Microsoft [9] pledges to advance AI using ethical principles that put people first. Other companies make similar statements [11, 7]. More recently, data-centric AI [4] is becoming critical where the primary goal is not to improve the model training algorithm, but to improve the data pre-processing for better model accuracy.

These trends motivate us to investigate data collection and quality challenges for deep learning from a data-centric AI perspective. Figure 1 shows a simplified end-to-end process starting from data collection to model deployment. Deep learning systems are more complicated in practice [117, 101], and we only show the essential steps. The first topic we cover is *data collection*. In comparison to traditional machine learning, in deep learning feature engineering is less of a concern, but there is instead a need for large amounts of training data. Unfortunately, many industries do not adopt deep learning simply because of the lack of data and the lack of explainability of the trained models. The second topic is *data cleaning and validation*. While there is a vast literature on data cleaning, unfortunately not all the techniques directly benefit deep learning accuracy [84]. In addition, there are recent deep learning issues including data poisoning that needs to be addressed, especially by the data management community. Data poisoning has the intention of reducing the model accuracy, and there is a branch of research called data sanitization where the goal is to defend against such attacks. The third topic is *robust model training*. Even after we carefully validate and clean our data, the data quality may still be problematic because there

is no guarantee that we fixed all the data problems. Hence, we may still need to cope with dirty, missing, or even poisoned data in model training. Fortunately, there are various robust training techniques [123] available. The fourth topic is *fair model training*. Traditional research on data management has not primarily focused on bias and fairness issues. However, in addition to improving model accuracy, also showing fairness against biased data is becoming essential for responsible AI. Model fairness research [18, 131, 38, 87] largely consists of fairness measures and unfairness mitigation before, during, or after model training.

While the coverage of this survey is broad, we believe it is important to have a birds-eye view of data issues in the entire deep learning process in order to make informed decisions. Each subtopic is not only substantial, but studied by different communities. Data collection, cleaning, and validation have been traditionally studied in the data management community. Robust model training is a central topic in the machine learning and security communities while fair model training is a popular topic in the machine learning and fairness communities. Both fairness and robustness topics are increasingly being researched in the data management community as well because they are closely related to the input data. Our contribution is to connect these related topics together at a high level with a focus on recent and significant works. Table 1 shows a taxonomy of the techniques covered in this survey. Our work targets researchers and practitioners who need a starting point of understanding how data plays a key role in data-centric AI.

In summary, deep learning is becoming prevalent thanks to big data and fast computation. Software 2.0 is a new paradigm for software engineering. However, big data for deep learning has been relatively understudied, but is becoming critical in data-centric AI. In each of the following sections, we cover the topics in Table 1.

- Data collection techniques for machine learning (Section 2).
- Data validation and cleaning techniques for machine learning (Section 3).
- Robust training techniques for coping with noisy and poisoned data (Section 4).
- Fair training techniques for coping with biased data (Section 5).

2 Data Collection

Our coverage of data collection originates from a survey [109] by one of the authors, but is refined and rewritten based on a tutorial [134]. There are three

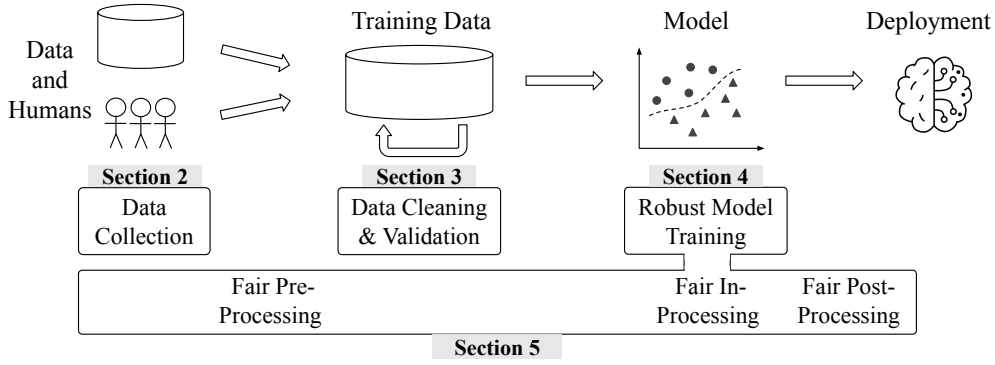


Fig. 1: Deep learning challenges from a data-centric AI perspective. Data collection and quality issues cannot be resolved in a single step, but throughout the entire machine learning process. This survey thus focuses on the breadth of available techniques.

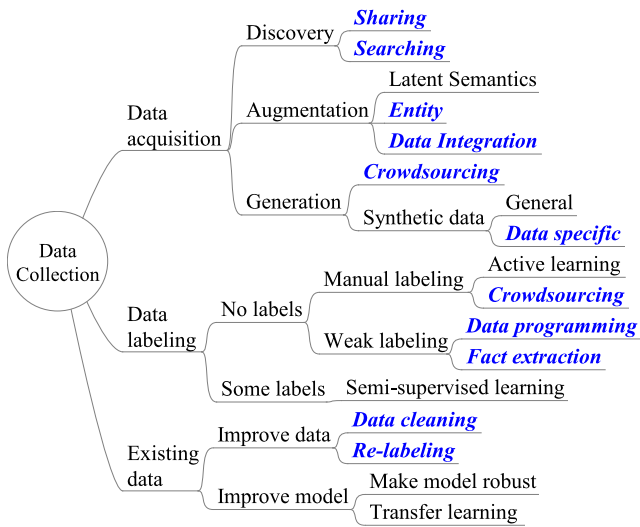


Fig. 2: Data collection for machine learning (ML) [109]. In the leaf nodes, the techniques highlighted in italic blue font are proposed by the data management community, and others are proposed by the ML community.

main approaches for data collection (Figure 2). First, *data acquisition* is the problem of discovering, augmenting, or generating new datasets. Second, *data labeling* is the problem of adding informative annotations to data so that a machine learning model can learn from them. Since labeling is expensive, there is a variety of techniques to use including semi-supervised learning, crowdsourcing, and weak supervision. Finally, if one already has data, *improving existing data and models* can be done instead of acquiring or labeling from scratch.

2.1 Data Acquisition

If there is not enough data, the first option is to perform data acquisition, which is the process of finding

datasets that are suitable for training machine learning models. In this survey we cover three approaches: data discovery, data augmentation, and data generation. *Data discovery* is the problem of indexing and searching datasets. *Data augmentation* is related, but focuses on complementing existing datasets by integrating them with external data. If there is not enough data around, the last resort is to take matters in one's own hands and create datasets using crowdsourcing or synthetic *data generation* techniques.

2.1.1 Data Discovery

Data discovery is the problem of indexing and searching datasets that exist either in corporate data lakes [128, 48] or the Web [30]. A representative example is the Goods system [56], which searches tens of billions of datasets in Google's data lake. A key challenge is that there are many data sources, ranging from file systems to Big tables and Spanners. It is infeasible to consolidate all the data in a single place, and one cannot simply ask the dataset owners to move their data either. Thus, Goods takes a post-hoc approach where it crawls the datasets and extracts metadata to maintain a central dataset catalog (Figure 3), which does not require any work from the dataset owners. Each entry in the catalog contains metadata about one dataset including its size, provenance on which job created it and which job read it, and schema information. Goods provides search, monitoring, and dataset annotations as well. More recently, a public version of Goods called Google Dataset Search [29], which facilitates search primarily for science datasets, has been released. In comparison to the internal Goods system, this public version assumes an open ecosystem where the dataset owners can publish semantically-enhanced metadata, including the purposes of using the dataset. Thus, once metadata is

Operation	Category	Section	Technique	Key References
Data Collection	Data Acquisition	2.1	Data discovery	[128, 48, 30, 56, 29, 147]
			Data augmentation	[52, 51, 77, 104, 62, 42, 145, 50]
			Data generation	[1, 98, 74]
	Data Labeling	2.2	Utilize existing labels	[129, 152, 138]
			Manual labeling	[2, 6, 118, 108]
			Automatic labeling	[103, 102, 17]
Data Validation and Cleaning	Improve Existing Data	2.3	Improve labels and model	[120, 93, 54]
	Data Validation	3.1	Visualization	[81, 130]
			False discovery control	[148, 49]
			Schema-based validation	[28, 19, 100]
			New functionalities	[55, 115, 105, 116, 88]
	Data Cleaning	3.2	Data-only cleaning	[65, 107]
			ML-aware cleaning	[85, 45, 76, 44]
	Data Sanitization	3.3	Data poisoning	[119, 151]
			Poisoning defenses	[41, 75, 64, 97]
Robust Model Training	Noisy Features	4.1	Adversarial learning	[106, 96, 32]
	Missing Features	4.2	Data imputation	[33]
	Noisy Labels	4.3	Robust learning	[123, 83, 86, 68, 57]
	Missing Labels	4.4	Semi-supervised learning	[127, 24]
Measuring Fairness	Statistical Fairness	5.2	Independence criteria ($\hat{Y} \perp Z$)	[47, 46]
			Separation criteria ($\hat{Y} \perp Z Y$)	[58, 140, 22]
			Sufficiency criteria ($Y \perp Z \hat{Y}$)	[37, 43, 22]
	Other Fairness	5.2	Individual fairness	[46]
			Causal fairness	[73, 78, 146, 91, 71]
Unfairness Mitigation	Pre-processing	5.3	Repair data	[47, 113, 69]
			Generate data	[36, 135]
			Acquire data	[126, 35, 16]
	In-processing	5.3	Fairness constraints	[141, 12, 70]
			Adversarial training	[142, 40, 110, 110]
			Adaptive reweighting	[111, 143, 67, 66]
	Post-processing	5.3	Fix predictions	[58, 39, 99]
	Convergence with Robustness	5.4	Fairness-oriented	[133, 80, 59, 79]
			Robustness-oriented	[144, 72, 136]
			Equal Mergers	[110, 112, 132, 122]

Table 1: Taxonomy of data collection and quality techniques for deep learning.

published, Google Dataset Search crawls that information automatically and builds an index out of it.

More recently, these data discovery tools have become more interactive. A representative system is Juneau [147], which is an interactive data search and management tool built on top of the Jupyter Notebook data science platform. In Figure 4, a notebook has a workflow of cells; the user can post a query search-

ing for related tables for a given table, and then the server search engine searches the data lake including other notebooks and returns a ranked list of related tables. Next, the user can choose to import some of the tables like T' . The server then inserts a new cell to generate the code to import table T' . Here the key technical challenge is finding the related tables. Juneau uses similarity measures based on row column overlap

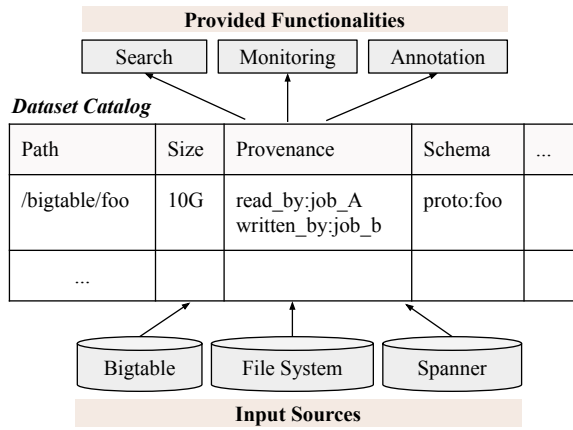


Fig. 3: The Goods system [56] maintains a catalog of dataset metadata.

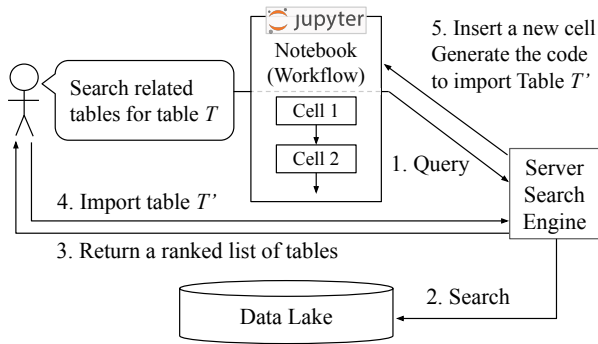


Fig. 4: The Juneau workflow [147] demonstrates interactive data discovery using Jupyter Notebooks.

and provenance information that intuitively captures the purpose of creating each data set.

2.1.2 Data Augmentation

For data augmentation, a popular method for generating data in the machine learning community is generative adversarial networks (GANs) [52, 51, 77]. As shown in Figure 5, we start from a training set that has real data, e.g., MNIST digit images. There are two components: a *generator* that generates fake data that is realistic using some random noise as an input and a *discriminator* that tries to distinguish the real data from the fake data of the generator. The generator and discriminator compete such that the generator tries to generate data that is realistic enough to trick the discriminator, while the discriminator tries to be good at identifying the generated data. If the training is done properly, the resulting generator produces realistic data.

One limitation of a GAN is that it cannot generate data that is completely different than the existing data. Using policies [104, 62] is a way to complement that limitation where one can apply various custom transforma-

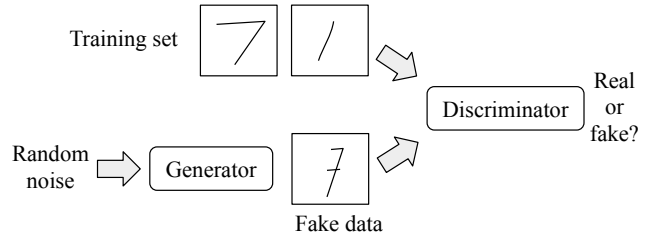


Fig. 5: A Generative Adversarial Network [52] competes two networks: a generator of fake data and a discriminator that distinguishes real from fake data.

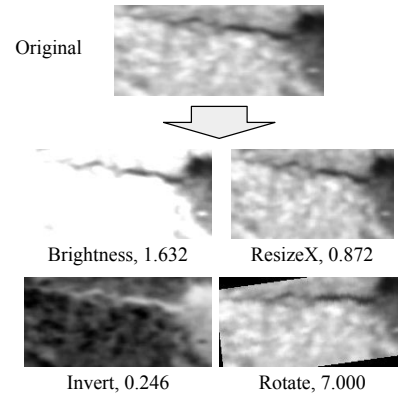


Fig. 6: Policies for data augmentation [62] can be used to transform existing data into other realistic data.

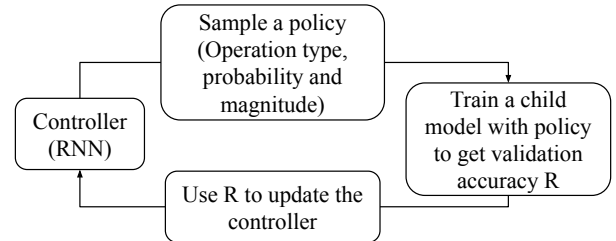


Fig. 7: Automatic policy searching [42] identifies the best data transformations for the purpose of data augmentation.

tions provided by domain experts as long as the data remains realistic. For example, starting from an original image, one can change the brightness, resize, invert, or rotate the image to create new data as long as the resulting images are also acceptable (Figure 6). AutoAugment [42] (Figure 7) automates this process where the idea is to have a controller that suggests a strategy for applying transformations with certain probabilities and magnitudes on the data. The system then trains a child model on this augmented data and measures the accuracy on a validation set. This result is then used to decide whether the strategy produces useful data that is within realistic bounds and should thus be used.

The data augmentation literature continues to grow rapidly. Mixup [145, 24, 139, 23, 61] has been proposed as a simple, but effective augmentation technique where the key idea is to mix pairs of data points of different classes. The additional data effectively regularizes the model to predict in-between training data points assuming linearity. Model patching [50] utilizes GANs to augment the data of specific subgroups of a class so that the model accuracy is similar across subgroups.

2.1.3 Data Generation

Another option for collecting or acquiring new data is to generate data. A popular option is to use crowdsourcing platforms like Amazon Mechanical Turk [1] where one can create tasks and pay human workers to create data. For example, a task may ask workers to find restaurant and review information. In addition, one can use a simulator or generator for specific domains, e.g., Hermoupolis [98] for mobility data and Crash to Not Crash [74] for driving data. We note that GANs also generate new data, but they require sufficient amounts of real data for model training.

2.2 Data Labeling

Once there are enough datasets, the next step is to label the examples. We cover data labeling techniques for utilizing existing labels and manually or automatically labeling from no labels.

2.2.1 Utilize Existing Labels

The traditional approach for labeling is semi-supervised learning [129, 152] where the idea is to use existing labels to predict the other labels as much as possible. The simplest form is *Self-training* [138] (Figure 8) where a model is trained on the available labeled data and then applied to the unlabeled data. Then the predictions with the highest confidence values are trusted and added to the training set. This approach assumes that we can trust the high confidence, but there are other techniques including Tri-training [150], Co-learning [149], and Co-training [26] that do not rely on this assumption.

2.2.2 Manual Labeling from No Labels

If there are no labels to start with, but one has funds to employ workers, a standard approach is to use crowdsourcing platforms like Amazon Mechanical Turk to perform labeling. Since labeling is such an important task, there are labeling-specialized services like Amazon

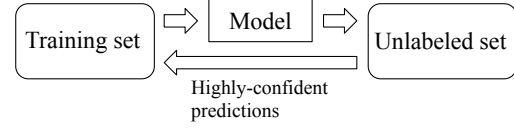


Fig. 8: Self training [138] for adding highly-confident predictions into the training set.

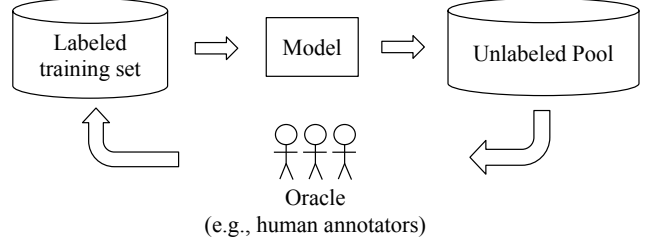


Fig. 9: Active learning workflow [118] where uncertain examples that are most likely to improve the model accuracy when labels are selected.

Sagemaker Ground truth [2] and Google Cloud Labeling [6]. When using Sagemaker, one can choose labeling tasks and recruit labelers who are assisted with a UI and tools to label the data. Sometimes, crowdsourcing may not be feasible because the workers do not have the right expertise. Hence, the last resort is to rely on domain experts, but this option can be expensive.

Active learning [118, 108] is an effective method to reduce the crowdsourcing cost. The idea is to ask human labelers to label uncertain examples that, when answered, are likely to improve model accuracy the most. Figure 9 show the overall process of starting from unlabeled data, asking an oracle or human annotators to label examples, and adding the labels back to the training set, which is then used to train better models. While a full coverage of active learning is out of scope, the example selection techniques can largely be categorized into identifying uncertain examples and using decision theoretic approaches to analyze the effect of a newly-labeled example on the model accuracy.

2.2.3 Automatic Labeling from No Labels

Recently, *weak supervision* is becoming popular where the idea is to (semi-)automatically generate labels that are not perfect (therefore called “weak” labels), but at scale where the larger volume may compensate for the lower label quality. Weak supervision is useful in applications where there are few or no labels to start with. Data programming is a specific weak supervision technique where multiple labeling functions are developed and combined to generate weak labels.

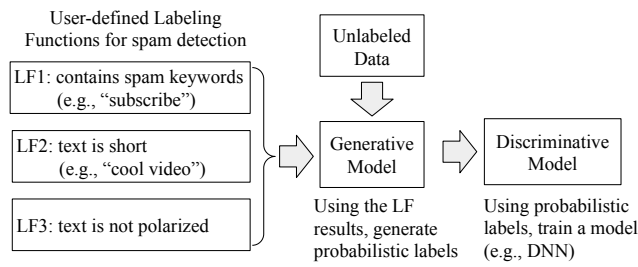


Fig. 10: Data programming workflow [103] of using multiple user-defined labeling functions to generate probabilistic labels, which are then used to train a final discriminative model.

Snorkel [103, 102] is the seminal system for data programming. To illustrate labeling functions, let us say that we want to perform spam detection as shown in Figure 10. The first labeling function looks for certain spam keywords like “subscribe” and, if so, it guesses that the text is spam. This logic by itself is not perfect, but it is useful enough. The second labeling function checks if the text is too short where the intuition is that short text is an indicator of spam. The third labeling function checks if the text is not polarized, which is yet another good signal for determining spam. There can be other labeling functions as well. Snorkel then combines all of the labeling functions in one generative model by intuitively taking a probabilistic consensus. Then given unlabeled data, Snorkel can generate probabilistic labels. The unlabeled data and the probabilistic labels are used to train a final discriminative model like a deep neural network. Another way to combine labeling functions is to use majority voting. Empirically, the number of labeling functions determines whether a generative model or majority voting is better.

Snorkel has been deployed on various workloads in Google and compared with various supervised learning systems [17]. As a result, Snorkel has the accuracy advantage over supervised learning approaches when there are few or no labels, and the supervised systems eventually catch up with more available labels.

2.3 Improving Existing Data

In addition to searching and labeling datasets, one can also improve the quality of existing data and models. This approach is useful in several scenarios. Suppose the target application is novel or non-trivial where there are no relevant datasets outside, or collecting more data no longer benefits the model’s accuracy due to its low quality. Here a better option may be to improve the existing data, which can be done in several ways. Alternatively,

suppose the target application is relevant to the previously studied applications. In this case, one can utilize transfer learning to directly improve the performance of the pre-trained models while avoiding the cost of re-generating new models.

An effective fix is to improve the labels through *re-labeling*. Sheng et al. [120] demonstrates the importance of improving labels by showing the model accuracy trends against more training examples for datasets with different qualities. As the data quality decreases, even if more data is used, the accuracy of the model does not increase from some point and plateaus. In this case the only way to improve the model accuracy is to improve the label quality, which can be done by re-labeling and taking majority votes on multiple labels per example. Another important approach is to clean the entire data, which we cover in Section 3.

In addition to improving labels and data, one can use *transfer learning* [93] to directly improve the model using pre-trained models. Nowadays, there are tools [54] for finding pre-trained models that perform similar tasks. The model training can then be performed starting from the pre-trained model and thus save labeling costs. For example, if the task is to perform sentiment analysis on product reviews where some are positive while others are negative, using a pre-trained model that also performs sentiment analysis on different products can boost the accuracy of model training without having to search for large numbers of reviews. Whenever training a model, it is thus a good idea to look for any related models to train on top of.

In summary, data collection techniques consist of data acquisition, data labeling, and improving existing data and models. Some of the techniques have been studied by the data management community while others by the machine learning community, but we are now seeing a convergence of techniques.

3 Data Validation and Cleaning

It is common for the training data to contain various errors. Machine learning platforms like TensorFlow Extended (TFX) [19] have data validation [100] components to detect such data errors in advance using data visualization and schema generation techniques. Data cleaning can be used to actually fix the data, and there is a heavy literature [65] on various integrity constraints. However, recent studies [84] show that only focusing on improving the data with well-defined errors does not necessarily benefit machine learning accuracy. In addition, in a machine learning point of view, we also need to address critical issues including data sanitization against malicious poisoning and model fairness,



Fig. 11: Facets [5] supports feature statistics (Overview) and data sample exploration (Dive).

which are actively studied in the security and machine learning communities, respectively. We cover data sanitization in this section and model fairness techniques in Section 5.

3.1 Data Validation

Data visualization is a widely-used and effective way to validate data for machine learning. A representative open-source tool is Facets [5], which shows various statistics of datasets that are relevant for machine learning (Figure 11). The visualization of statistics is useful for sanity checks on data to prevent larger errors downstream. There are two components: Overview and Dive. Overview shows for each feature statistics including mean, standard deviation, the portion of zero values, histograms, and so on. While this information is simple, it can help identify critical errors early on. For example, a histogram showing the number of positive and negative labels is useful to ensure that there is no class imbalance. The other component Dive enables a closer look at data samples and visually finding correlations among features. Any problems that are not fixed here may be harder to fix downstream because the model training will not necessarily crash, but may just perform slightly worse, making the problem go undetected. Debugging such errors may take hours or days to figure out the root cause.

In addition to manual visualization, there has also been research on *automatic generation* of new visual-

izations [81] that can be used for validation purposes. SeeDB [130] is a seminal framework that repeatedly generates visualizations of interest. To capture the notion of interestingness, SeeDB uses a deviation-based utility metric that gives a high value when groupings of the data result in different probability distributions. Suppose we compare the internet access usage of different groups where we compare the desktop usage and the mobile usage. When comparing female versus male groups, the distribution of desktop versus mobile may not be that different, so the visualization is considered uninteresting. On the other hand, if we compare emerging and mature markets, emerging markets may have relatively higher mobile users while mature markets still mainly use desktops for Internet access. The different distributions of usage are then considered interesting.

Automatically generating visualizations can run into the problem of false positives, so there is also a line of research that proposes *false discovery control* techniques. CUDE [148] is a technique for controlling false discovery in the context of multiple hypothesis testing for visual interactive data exploration. Here users can repeatedly generate visualizations and mark the ones that are significant. Based on this user feedback, we would like to automatically choose the visualizations that are significantly interesting in a statistical sense. Unfortunately, traditional methods for controlling false discoveries like Bonferroni correction and the Benjamini-Hochberg procedure [21] assume static hypotheses where there is a fixed number of visualizations to compare. However, these methods cannot be applied to an interactive data exploration setting where there is a continuous stream of visualizations. CUDE thus proposes a hypothesis testing method called alpha investing [49] where it controls a measure that captures the false discovery rate. Intuitively, we start with a budget α for hypothesis testing that can either run out or increase if interesting visualizations are discovered.

Schema-based validation [19, 100] is widely used in practice. Tensorflow Data Validation (TFDV) [28] assumes a continuous training setting where input data periodically streams in as shown in Figure 12. TFDV generates a data schema from previous data sets and uses the schema to validate future data sets and alert users on data anomalies. For each anomaly, TFDV provides concrete action items to possibly fix the root cause. A schema here is different from a traditional database schema where it contains a summary of data statistics of the features. In case a new dataset violates the current schema, either the data needs to be fixed, or the schema needs to be updated, and the user makes the decision. For example, a country code feature may have been previously known to only have certain coun-

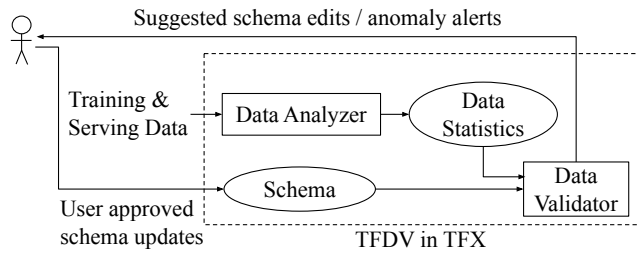


Fig. 12: TensorFlow Data Validation (TFDV) [28] uses a user-approved schema to validate the statistics of the training and serving data.

try codes according to the schema. However, if the data contains a new country code that the user considers valid, the country code can be added to the schema. TFDV has been used by hundreds of product teams within Google and validates several petabytes of production data per day. According to usage statistics in a 30-day period, the most common anomalies include the appearance of new features and out-of-bound values. These results demonstrate that most of the anomalies are simple, but essential to fix early on.

More recently, data validation systems are equipped with additional functionalities including declarative abstractions [55, 115], automatic identification of error types [105], testing the impact of errors on models [116], and ease of usage [88].

3.2 Data Cleaning

Data cleaning has a long history of removing various well-defined errors by satisfying integrity constraints including key constraints, domain constraints, referential integrity constraints, and functional dependencies. For an introduction see the book *Data Cleaning* [65].

We first introduce one of the state-of-the-art data cleaning techniques to give a sense of how sophisticated these techniques have become. HoloClean [107] repairs data using probabilistic inference using three main ingredients: satisfying various integrity constraints, using external dictionaries to check the validity of values, and using quantitative statistics. Suppose we have a database of restaurants where the schema contains the attributes name, address, city, state, and zip code. If there is a restaurant with an incorrect city name, a functional dependency from zip code to city can be used to fix the error by looking at other restaurants with correct city names and zip codes. If a zip code is incorrect, but there are no functional dependencies available, then external address listings can be used to fix the zip code values. In the case where a name looks like a statistical outlier compared to other names, but there is no

external information, the name can still be fixed using non-outliers. The main point here is that the data can have various data errors and there is no one technique that fixes them all. Instead, HoloClean applies all the fixes together probabilistically so that they are satisfied as much as possible.

Unfortunately, only focusing on fixing the data does not necessarily guarantee the best model accuracy. At first glance, it seems that perfectly cleaning the data would be most useful for the model training. However, the notion of clean data is not always clear cut, and removing all possible errors is not always feasible. CleanML [85] is a framework that evaluates various data cleaning techniques and seeing if they actually help model accuracy. The authors use 13 real data sets and considered five common error types including missing values, outliers, duplicates, inconsistencies, and missed labels. Seven machine learning models are evaluated: decision tree, random forest, Adaboost, XGBoost, and Naïve Bayes. As result, data cleaning does not necessarily improve downstream machine learning models. In fact, the cleaning may sometimes have a negative effect on the models. However, by selecting an appropriate machine learning model, one can eliminate the negative effects of data cleaning. Also there is no single cleaning algorithm that performs the best, and one must adaptively choose the algorithm depending on the type of noise to clean. Hence, data cleaning techniques that are not originally designed for machine learning must be used carefully.

Recently there are data cleaning techniques with the specific purpose of improving model accuracy [45]. ActiveClean [76] is a seminal framework that iteratively cleans samples of dirty data and updates the model. Figure 13 shows the workflow where there is a sampler that chooses an example that is likely to be dirty, and data quality rules can be used to identify such dirty samples. The reason for sampling data is that cleaning the entire data is presumed to be very expensive. Each sample can be cleaned by an oracle or domain expert. Then the updater updates the model to be more accurate and also chooses the next sample. ActiveClean has theoretical guarantees where, by repeatedly training a model on the clean sample plus previously-cleaned data, the model eventually obtains an accuracy as if it was trained on clean data only. ActiveClean assumes convex loss models like SVMs, and the data cleaning is assumed to be done perfectly.

In addition, another branch of research is to clean the labels for the purpose of improving model accuracy. TARS [44] is a system that predicts model accuracy out of noisy labels that are produced from crowdsourcing. TARS first chooses labels that are likely to be flipped

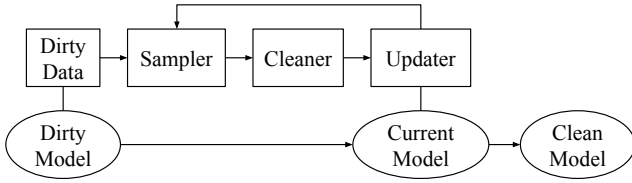


Fig. 13: ActiveClean [76] iteratively selects data that is likely to be dirty and cleans it.

because they were labeled by incompetent workers and thus have low confidence values. TARS then estimates how much the model will improve if the label is flipped after cleaning. The confidence values of labels can be computed by using confusion matrices of workers, which capture the history of how well they performed in past tasks. A confusion matrix thus contains the previous false positive, false negative, true positive, and true negative rates. Given the probability that a label is flipped, TARS estimates the resulting model accuracy and subtracts that by the estimated accuracy of the current model to determine whether the label is worth examining. Hence, TARS can selectively clean labels that are expected to benefit model accuracy the most.

3.3 Data Sanitization

Data poisoning has recently become a serious issue because changing a fraction of the training data, which may come from an untrusted source, may alter the model’s behavior. Compared to dirty data, there is a malicious intent to make the model fail. Data poisoning is a real problem because data is now easier to publish through dataset search engines. A dataset owner can simply post metadata to the public, which will be automatically crawled by the search engine. Then one can simply harvest that data using web crawlers without knowing that the data is poisoned. Data sanitization [41] is the problem of defending against such poisoning attacks and can be viewed as an extreme version of data cleaning.

The basic type of data poisoning is called *label flipping* where a label is flipped from one class to another. More recently the data poisoning techniques have become much more sophisticated and therefore harder to defend against [119, 151]. We illustrate a state-of-the-art data poisoning techniques for deep learning [151]. A major challenge when poisoning data for deep learning is that the victim’s model cannot be easily analyzed. Hence, transferable poisoning attacks have been proposed, which can still succeed without accessing the victim’s model. The idea is to train an ensemble of substitute models, which are assumed to be similar to the

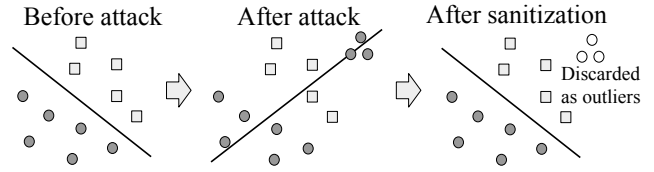


Fig. 14: Data Sanitization [41] identifies and discards data poisoning for better model accuracy.

victim’s model. Any attack that negatively affects the substitute models will presumably attack the victim’s model as well. Given a set of clean data points of different classes, the poisoning algorithm adjusts the clean points to “move closer” to the target within the feature space and form a convex polytope that surrounds it to maximize the chances of the target to be misclassified.

How do we defend against such data poisoning using data sanitization? Figure 14 shows a simple setting where a classifier’s behavior changes after introducing poisoning (top right data points). If the data sanitization can discard the poisonings as outliers, then the model’s accuracy can be restored. Data sanitization techniques [41, 64, 97] have been proposed throughout the years, and a recent study [75] evaluates various defenses by developing attacks and seeing if the defenses work are still effective. For example, a simple defense is to remove points that are far away from class centroids using the L2 distance, which should work for the example in Figure 14. Unfortunately, the conclusion is that no technique can adequately defend against carefully designed attacks. We suspect that data poisoning and sanitization techniques will continue to evolve and compete with each other.

In summary, we covered key approaches in data validation, data cleaning, and data sanitization. Data validation can be performed using visualizations and schema information. Data cleaning has been heavily studied where recent techniques are more tailored to improving model accuracy. Data sanitization has the different flavor of defending against poisoning attacks.

4 Robust Model Training

Even after collecting the right data and cleaning it, data quality may still be an issue during model training. It is widely agreed that real-world datasets are dirty and erroneous despite the data cleaning process. These flaws in datasets can be categorized depending on whether data values are noisy or missing and depending on whether these flaws exist in data features (attributes) or labels.

	Noisy	Missing
Features	Adversarial Learning (Section 4.1)	Data Imputation (Section 4.2)
Labels	Robust Learning (Section 4.3)	Semi-Supervised Learning (Section 4.4)

4.1 Noisy Features

Noisy features are often introduced by adversarial attacks. Among several types of attacks, we focus on the *poisoning attack*, which is known as contamination of the training data, to be aligned with the main theme of this survey. During the training phase of a machine learning model, an adversary tries to poison the training data by injecting maliciously designed data to deceive the training procedure.

Either features or labels or both can be the target of the poisoning attack. The poisoning attack can be done in three ways depending on the capability of adversaries. First, an adversary can randomly perturb the labels, i.e., by assigning other incorrect labels, picked from a random distribution, to a subset of training data. The robust training methods for this type of attacks will be discussed in Section 4.3. Second, an adversary is more powerful and can corrupt the features of the examples possibly determined by analyzing the training algorithm [25]. Third, an adversary may insert adversarial examples into the training data. For more details, the reader can refer to [31].

Various defense strategies have been actively studied for robust training on adversarial examples (e.g., noisy features). Most of the current strategies are not adaptive to all types of attacks, but are effective to only a specific type. We summarize a few well-known, representative strategies in this section.

Most notably, in *adversarial training*, the robustness of a model can be improved using a modified objective function [53],

$$\tilde{J}(\theta, \mathbf{x}, y) = \alpha J(\theta, \mathbf{x}, y) + (1 - \alpha)J(\theta, x + \epsilon \text{sign}(\nabla_x J(\theta, \mathbf{x}, y)), y),$$

which is based on the fast gradient sign method. Here, J is the original loss function, and ∇ represents the gradient operator. By the modified objective function, the model is forced to predict the same class for legitimate and perturbed examples in the same direction.

Knowledge distillation has been shown to be effective for adversarial attacks [94]. *Defensive distillation* is almost the same as typical knowledge distillation, except that the same network architecture is used for both the original network and the distilled network. Specifically, instead of hard labels, where only the true label

has the probability 1 in a probability vector, *soft targets*, which are generated by the original network as the prediction result, are used for training the distilled network. The benefit of using soft targets comes from the additional knowledge found in probability vectors compared to hard class labels [94].

Feature squeezing [137] reduces the degree of freedom available to an adversary by squeezing out unnecessary input features. If the original and squeezed inputs result in substantially different outputs by a model, the corresponding input is determined to be adversarial. A popular squeezing technique for images is reducing the color depth on a pixel level.

Another idea is to detect adversarial examples using separate classification networks [90]. A sub-network, called an *adversary detection network* or simply a detector, is trained to produce an output that indicates the probability of the input being adversarial. For this purpose, a classification network is trained using only non-adversarial examples, and adversarial examples are generated for each example in the training set. Then, the detector is trained using both the original dataset and the corresponding adversarial dataset. MagNet [89] falls into this strategy, and it also contains a reformer that corrects an adversarial example to a similar legitimate example using autoencoders.

4.2 Missing Features

Since missing data can reduce the statistical power and produce biased estimates, data imputation has been an active research topic in statistics and machine learning. Missing features can occur in any types of data, but, in this paper, we focus on *multivariate time-series* data because missing values are unavoidable because of its high input rate and sensor malfunction.

Missing values in multivariate time-series data are ubiquitous in many practical applications ranging from healthcare, geoscience, astronomy, to biology and others. They often inevitably carry missing observations due to various reasons, such as medical events, saving costs, anomalies, inconvenience, and so on. These missing values are usually informative where the missing value and patterns provide rich information about target labels in supervised learning tasks.

We first describe informative missingness. Figure 15 shows MIMIC-III critical care dataset [33]. Starting from the bottom, there are the missing rate of each variable, the correlation between missing rate of each variable and mortality, and the correlation between missing rate of each variable and each ICD-9 diagnosis category. Here, we observe that the values of missing rates are

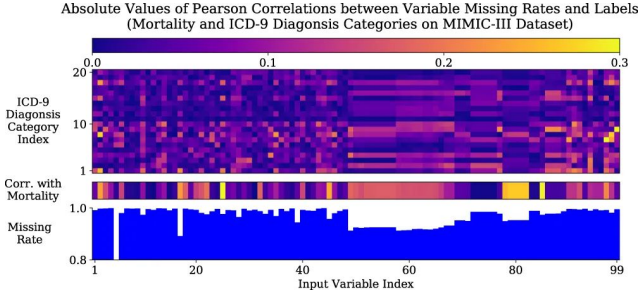


Fig. 15: Informative missingness in the MIMIC-III dataset [33].

correlated with labels, where the values with low missing rates are highly correlated with the labels. In other words, the missing rate of variables of each patient is useful, and this information is more useful for the variables that are more often observed in the dataset.

For existing approaches, a simple solution is to omit the missing data and to perform analysis only on the observed data, but it does not provide good performance when the missing rate is high and the samples are inadequate. Another solution is to fill in the missing values with substituted values, which is known as *data imputation*. However, these methods do not capture variable correlations and may not capture complex patterns to perform imputation. Combining the imputation methods with prediction models often results in a two-step process where imputation and prediction models are separated; the missing patterns are not effectively explored in the prediction model, thus leading to suboptimal analysis results.

GRU-D [33] is a deep learning model based on the gated recurrent unit (GRU) to effectively exploit two representations of informative missingness patterns—masking and time interval. Masking informs the model of which inputs are observed or missing, while time interval encapsulates the input observation patterns. GRU-D captures the observations and their dependencies by applying masking and time interval, which are implemented using a decay term, to the inputs and network states of the GRU, and jointly train all model components through back-propagation. GRU-D not only captures the long-term temporal dependencies of time-series observations, but also utilizes the missing patterns to improve the prediction results.

We elaborate on the two components of GRU-D: *masking* and *time interval*. The value of the missing variable tends to be close to some default value if its last observation happens a long time ago, because the influence of an input variable will fade away over time if the variable has been missing for a while. For example, one medical feature in electronic health records (EHRs) is

only significant in a certain temporal context. As lots of missing patterns are informative and potentially useful in prediction tasks, but unknown and possibly complex, the goal is to learn decay rates from the training data rather than fixing them a priori. The GRU-D model incorporates two different trainable decay mechanisms. For a missing variable, an input decay γ_x is added to decay it over time toward the empirical mean, instead of using the last observation as is. A hidden state decay γ_h in GRU-D has an effect of decaying the extracted features (GRU hidden states) rather than raw input variables directly.

In summary, handling missing values in multivariate time-series data using recurrent neural networks is challenging. GRU-D captures the informative missingness by incorporating masking and time interval directly inside the GRU architecture. GRU-D is shown to significantly outperform non-deep learning methods on synthetic and real-world healthcare datasets.

4.3 Noisy Labels

Because data labeling is done manually in many cases, incorrect or missing labels are, in fact, very common; the proportion of incorrect labels is reported to be 8–38% in several real-world datasets [123]. For example, if you search for Cheetah images, many of them are actually other animals like Jaguars instead of Cheetahs. The reason for these incorrect results includes incorrect annotations. The wrong annotations are caused by human errors. In our example, it may be difficult to distinguish patterns of Cheetahs and Jaguars. So wrong annotations can be caused by such human errors. Another type of error is software error, e.g., during object recognition. In the case where there are many images to annotate, one may use automatic object recognition software. However, the object recognition itself may have errors. Thus, many deep learning techniques have been developed to consider the existence of label noises and errors, which are more critical in deep learning than in conventional machine learning as a deep neural network (DNN) completely memorizes such noises and errors because of its high expressive power.

We explain what kinds of problems occur with noisy labels. In standard supervised learning, training data consists of example and label pairs $\{(x_i, y_i)\}_{i=1}^N$. In a practical setting, however, the label y_i is actually \tilde{y}_i , which means it can be incorrect. If one trains powerful models like VGG-19 on noisy data, the model may simply memorize the noise as well and perform worse on clean data. The goal of the noisy label problem is to train the network as if there are no noisy labels.

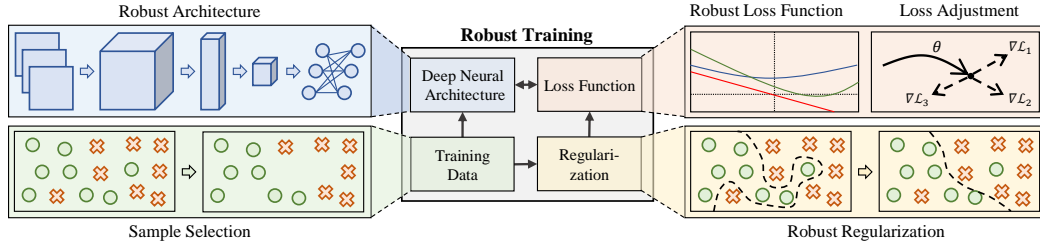


Fig. 16: Robust training categorization [123].

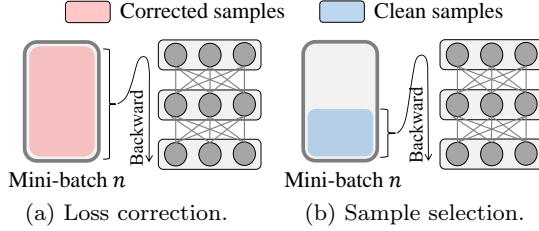


Fig. 17: Two directions of robust training covered in this survey.

Figure 16 from a recent survey [123] shows the categorization of robust training techniques. There are largely four components in the training procedure: deep neural architecture, loss function, input training data, and regularization. For each component, there are relevant robust training techniques. For deep neural architectures, robust architectures have been developed. For training data, various sample selection techniques have been proposed. For the loss function, robust loss functions and loss adjustment techniques have been proposed. More specifically, loss adjustment can be further divided into loss correction [96], loss reweighting [32], and label refurbishment [106]. For regularization, robust regularization techniques have been proposed.

In this survey, we focus on the most representative techniques: sample selection and loss correction techniques as illustrated in Figure 17. *Loss correction* is to correct the loss of *all* samples before a backward step. The representative techniques include Bootstrap [106], F-correction [96], and ActiveBias [32]. *Sample selection* is to select *expectedly clean* samples to update the network. The representative techniques include Decouple [86], MentorNet [68], and Coteaching [57].

We first introduce ActiveBias [32], which is a loss correction technique. ActiveBias performs a forward step on a given mini-batch and computes the sample importance for each sample. There are many statistics for the importance, e.g., variance of predictions. ActiveBias then corrects the loss by multiplying the normalized importance. If a label is noisy, then its importance μ_i decreases. The corrected loss is used to update the network. We then explain sampling selection through

its representative technique Coteaching [57], where the noise rate τ is assumed to be given. It then performs a forward step on a given mini-batch and selects the $(100 - \tau)\%$ low-loss samples as clean samples. The network is updated using the loss of the clean samples.

Although these two methods have improved the robustness to noisy labels, there are limitations of the two methods. Loss correction suffers from accumulated noise due to the large number of false corrections. Since all the examples are used for the training step, false corrections can accumulate for heavily noisy data. On the other hand, sample selection uses only clean samples having low losses and easy to classify. Hence, we may end up ignoring many useful, but hard samples that are classified as unclear.

SELFIE [123] was proposed to overcome the above limitations by using a hybrid of loss correction and sample selection (see Figure 18). SELFIE introduces *refurbishable* samples where labels can be corrected with high precision. The key issue of SELFIE is constructing the refurbishable and clean samples. For the clean samples, SELFIE adopts the small-loss trick [57] and uses the $(100 - \tau)\%$ low-loss samples in the mini-batch. The refurbishable samples are ones that have consistent label predictions. Each label is replaced with the most frequently predicted label during the training step. For example, if an image is predicted mostly as a dog and only sometimes a cat, then the label predictions are considered consistent, and such a cat label is considered noisy and corrected to a dog. Finally, the loss is calculated for the refurbishable samples with correct labels and the clean samples; the samples that are neither refurbishable nor clean are discarded. The advantage of SELFIE is that it minimizes false corrections during the model training by selectively correcting refurbishable samples. As a result, the correction error of refurbishable samples is low. Also as the training progresses, the number of refurbishable samples also increases, so most training samples are exploited in the end.

Prestopping is another technique [123] for avoiding overfitting to noisy labels by early stopping the training of a deep neural network before the noisy labels are severely memorized. The algorithm resumes

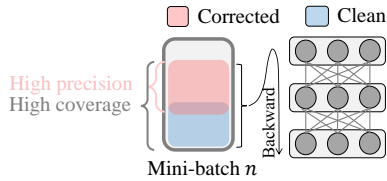


Fig. 18: SELFIE is a hybrid of loss correction and sample selection [123].

training the early-stopped network using a maximal safe set, which maintains a collection of almost certainly true-labeled samples. MORPH [124] further improves Prestopping through a novel concept of *self-transitional learning*, which automatically switches its learning phase at the transition point. The optimal transition point is determined without any supervision such as a true noise rate and a clean validation set, which are usually hard to acquire in real-world scenarios. MORPH rather estimates the noise rate by fitting the loss distribution to a one-dimensional and two-component Gaussian mixture model (GMM).

DivideMix [83] is a recent technique that trains two networks simultaneously. At each epoch, a network models its per-sample loss distribution with a GMM to divide the dataset into a labeled set (mostly clean) and an unlabeled set (mostly noisy), which is then used as training data for the other network (i.e., co-divide). At each mini-batch, a network performs semi-supervised training using an improved MixMatch [24] method, which we cover in the next section. When training on the CIFAR-10 dataset with 40% asymmetric noise, standard training with cross-entropy loss causes the model to overfit and produce over-confident predictions, making the loss difficult to be modeled by the GMM. Also, adding a confidence penalty during the warm up leads to more evenly-distributed loss. Finally, training with DivideMix can effectively reduce the loss for clean samples while keeping the loss larger for most noisy samples.

4.4 Missing Labels

We cover the issue of missing labels where training labels may not exist for either some or all examples. There are largely semi-supervised and unsupervised approaches. In semi-supervised approaches, clean labeled data exists together with unlabeled (or incorrectly labeled) data. The goal is to exploit unlabeled data to improve accuracy as much as possible. Here, the loss is defined as the supervised loss for labeled data plus the unsupervised loss for unlabeled data. The representative techniques include unsupervised loss (e.g., con-

sistency loss) like Mean-Teacher [127] and augmentation techniques like MixMatch [24]. For unsupervised approaches, the representative techniques include self-supervised learning and generative models, and we will cover a self-supervised learning technique called JigsawNet [92].

We now explain Mean-Teacher [127]. Here, the teacher model is the average of consecutive student models. In a training batch, both the student and teacher models evaluate the input with noise in their computations. The softmax output of the student model is compared with the one-hot label using a classification cost. Additionally, the output is compared with the teacher output using the consistency loss. After the weights of the student models are updated via gradient descent, the teacher model weights are updated as an exponential moving average of the student model weights. While both models are usable for prediction, the teacher model is more likely to be correct at the end of training. A training step with *unlabeled* examples is done without the classification cost.

Another semi-supervised technique called MixMatch [24] is very popular. To exploit an unlabeled dataset, MixMatch performs label guessing where stochastic data augmentation is applied to an unlabeled image K times, and each augmented image is fed through the classifier. The average of these K predictions is sharpened by adjusting the distribution’s temperatures. The MixMatch algorithm then mixes both labeled examples and unlabeled example with label guesses. More specifically, when mixing two images, the images are overlaid, and the labels are averaged, following the MixUp augmentation [145].

Now, let us cover unsupervised techniques. Since there are no labels, we need to develop new tasks exploiting labels that can be obtained from the data for free. JigsawNet [92] is one of such techniques. If an image is divided into smaller regions without labels, we can randomize the regions and solve the jigsaw puzzle where we know the correct order and positions, as illustrated in Figure 19. JigsawNet trains a context-free network (CFN) called a siamese-ennead convolutional network to solve the jigsaw task. The trained network can be transferred or fine-tuned to solve the given task using a small amount of labeled data.

To summarize the noisy or missing labels research in Sections 4.3 and 4.4, these problems incur poor generalization on test data. Existing work for noisy labels suffers from either (i) accumulated noise or (ii) partial exploration of training data. Hybrid (e.g., SELFIE) and semi-supervised techniques (e.g., DivideMix) can achieve very high accuracy with noisy training data. Semi-supervised (e.g., MixMatch) and

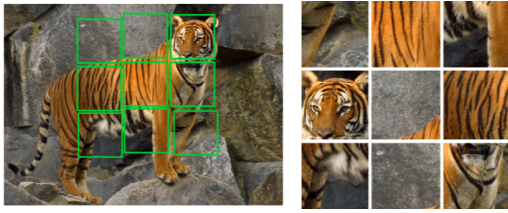


Fig. 19: Example of the jigsaw puzzle task for a given unlabeled image [92].

self-supervised (e.g., JigsawNet) techniques are actively developed to exploit abundant unlabeled data.

5 Fair Model Training

While the previous sections focus on improving model accuracy, we now focus on the issue of fairness where biased data may cause a model to be discriminating and thus unfair. In particular, we discuss how to measure fairness and how to mitigate unfairness. In addition, we discuss a recent trend of how fair and robust techniques are converging. This section extends recent tutorials [134, 82] by the authors.

5.1 Motivation and Issues

Model fairness is becoming a critical issue where bias in the data may result in discriminatory behavior by the model. A famous example is the COMPAS tool by Northpointe, which predicts a defendant’s risk of committing another crime. According to an analysis by ProPublica [15], black defendants are far more likely to be judged as high risk compared to white defendants, which turns out to be inaccurate in practice. This investigation fueled the new research area of algorithmic fairness. There could be multiple reasons why COMPAS discriminates. The training data could be biased where there is more data for certain demographics. Or there can be external factors where the surrounding environment may have caused more crime than race itself. Even the fairness measure can be in question where it does not accurately reflect reality. In general, analyzing fairness can be an extremely complicated issue that involves factors outside the data. A full discussion on fairness and ethics can be found in the recent fair ML book [18], and here we only focus on fairness issues with technical solutions.

5.2 Fairness Measures

Fairness cannot be described by one notion, and there are tens of possible definitions summarized in various

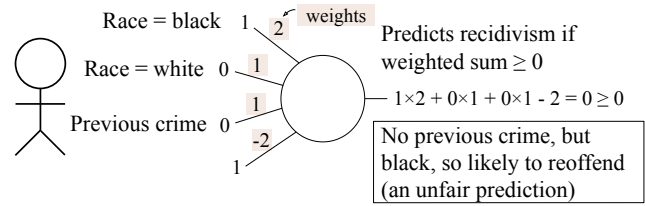


Fig. 20: A perceptron receiving input features and performing a weighted sum. Even a model as simple as a perceptron may show unfairness.

surveys [18, 131, 38, 87] used for predicting crime, hiring, giving loans, and more. We illustrate representative measures using a running example and then categorize them according to reference [18] as shown in Table 1. We use the following notations: Y denotes the label of a sample, \hat{Y} the prediction of a model, and Z is a sensitive attribute like race or gender.

Example 1 We illustrate fairness using the simplest-possible model: a perceptron, which is the most basic unit in a neural network. Suppose the perceptron receives input features as shown in Figure 20. Let us say that we have three features: “Race = black” has a value of one if the person is black or zero otherwise. “Race = white” has one if it is a white person or zero otherwise. “Previous crime” is one if the person has a previous crime and zero otherwise. The last feature is a constant to make the prediction threshold equal to zero. Given an example, we take the weighted sum by multiplying the feature values with the weights [2, 1, 1, -2] and, if the sum is at least zero, the model predicts recidivism and otherwise not. For example, if a white person committed a previous crime, the weighted sum is $0 \times 2 + 1 \times 1 + 1 \times 1 - 1 \times 2 = 0$, which is larger or equal to the threshold zero. The interpretation is that the person previously committed a crime, so is likely to re-offend. A black person who committed a previous crime gets the same prediction. However, for people who did not commit a previous crime, the model starts to discriminate where only a black person is predicted to still re-offend as shown in Figure 20. The prediction is obviously unfair and is shown for illustration purposes to show how even a model as simple as a perceptron can be discriminating.

For our running example, let us assume there are four people: (a) one white person who committed a crime before and committed a crime again, (b) one white person who never committed a crime, (c) one black person who committed a crime before and committed a crime again, and (d) one black person who never committed a crime. In this example, the perceptron correctly predicts for a, b, and c but wrongly predicts for d.

Demographic parity [47,46] requires that sensitive groups must have the same positive prediction rates. The formulation is as follows: $P(\hat{Y} = 1|Z = 0) = P(\hat{Y} = 1|Z = 1)$ where $Z = 0$ (1) means white (black) person or any two sensitive group. $\hat{Y} = 1$ means that the prediction of the model is positive, e.g., predicting recidivism. Demographic parity says that the probability of a prediction being one must be the same for these two groups.

Example 2 In our running example, $P(\hat{Y} = 1|Z = 0) = 0.5$ while $P(\hat{Y} = 1|Z = 1) = 1$, which shows unfairness.

Equalized odds [58,140,22] is defined as $P(\hat{Y} = 1|Z = 0, Y = A) = P(\hat{Y} = 1|Z = 1, Y = A)$, $A \in \{0, 1\}$. That is, we would like to guarantee demographic parity when the label Y is zero (in our example, the person did not commit crime again) and when Y is one (the person committed crime) separately. In other words, equalized odds says that the accuracy conditioned on the true label must be the same for the groups.

Example 3 In our running example, $P(\hat{Y} = 1|Z = 0, Y = 1) = P(\hat{Y} = 1|Z = 1, Y = 1) = 1$, but $P(\hat{Y} = 1|Z = 0, Y = 0) = 0 \neq P(\hat{Y} = 1|Z = 1, Y = 0) = 1$, so there is some unfairness.

Predictive parity [37,43,22] is defined as $P(Y = 1|Z = 0, \hat{Y} = 1) = P(Y = 1|Z = 1, \hat{Y} = 1)$. That is, given that the predictions are positive, we would like the actual likelihood of the label being positive to also be the same. Note that this measure can be extended to other label classes (e.g., $Y = 0, \hat{Y} = 0$).

Example 4 In our running example, $P(Y = 1|Z = 0, \hat{Y} = 1) = 1 \neq P(Y = 1|Z = 1, \hat{Y} = 1) = 0.5$, which shows unfairness.

Interestingly, many statistical fairness measures are equivalent to or variants of the following fairness criteria [18]: independence: $\hat{Y} \perp Z$, separation: $\hat{Y} \perp Z|Y$, and sufficiency: $Y \perp Z|\hat{Y}$. Note that demographic parity is equivalent to independence, equalized odds is equivalent to separation, and predictive parity is equivalent to sufficiency. An impossibility result says that no two fairness criteria can be fully satisfied together (see proofs in [18]).

There are remaining fairness criteria beyond the three above, and we cover the two popular ones: individual fairness and causality fairness.

Individual fairness [46] only uses the classifier for its definition and is defined as $D(f(x), f(x')) \leq d(x, x')$ where d is a distance function among examples, and D is a distance function between outcome distributions. Intuitively, the predictions for similar people must be

similar as well. For example, if two people are similar to each other, then their recidivism rates must be similar as well. Choosing proper distance functions is one of the key challenges in individual fairness.

Causality fairness [73,78,146,91,71] assumes a causal model, which is a diagram of relationships between attributes. In particular, an edge from attribute A to attribute B means that A 's value affects B 's value.

Example 5 In our running example, the race attribute apparently affects crime. However, it could also be the case that race affects the zip code of the person's address (i.e., zip code is known to be correlated to race), which provides an environment for committing more or less crime. A causal graph could thus have three nodes *race*, *zip code*, and *crime* with edges from *race* to *zip code*, *race* to *crime*, and *zip code* to *crime*. One can perform a counterfactual analysis to see if zip code indeed affects crime rates by comparing similar people that live or do not live in that zip code.

5.3 Unfairness Mitigation

Although there are many ways to measure fairness, one would ultimately like to perform *unfairness mitigation* [18,20]. Data bias can be addressed either before, during, or after model training. Each of these approaches are referred to as *pre-processing*, *in-processing*, and *post-processing* approaches respectively. Pre-processing approaches can be viewed as data cleaning, but with a focus on improving fairness. For each approach, we cover representative techniques.

5.3.1 Pre-processing Mitigation

The goal is to fix the unfairness before model training by removing data bias. The advantage is that we may be able to solve the root cause of unfairness within the data. A disadvantage is that it may be tricky to ensure that the model fairness actually improves when we only operate on the data. A naïve approach that does not work is to remove sensitive attributes (also referred to as unawareness) because they are usually correlated with other attributes. For example, removing sensitive attributes like race, income, and gender does not ensure fairness because their values can be inferred using correlated attributes like zip code, credit score, and browsing history, respectively. We cover three natural approaches for pre-processing – repairing data, generating data, and acquiring data – as shown in Figure 21.

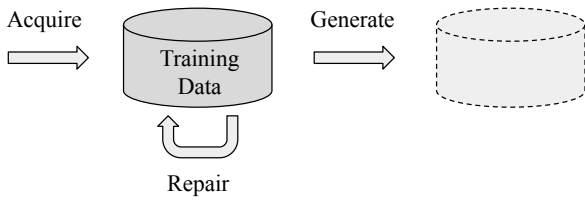


Fig. 21: Pre-processing: preparing unbiased data. There are three approaches: fair data repair, fair data generation, and fair data acquisition.

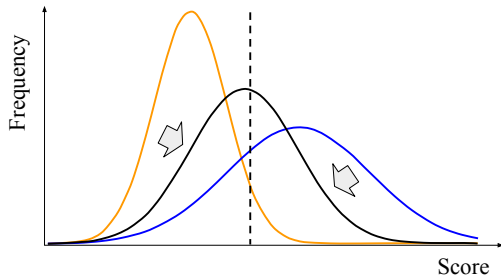


Fig. 22: Repairing distribution by averaging scores of the same percentiles [47].

Fair data repairing We first cover a method [47] that guarantees demographic parity while preserving important statistics like the ranking of data. As a simple example (Figure 22), let us say that we have test scores and distributions for different genders where the data follow normal distributions, but the women’s distribution has a higher mean and larger variance. Now let us say that we want to train a model that uses the test score to make a prediction. If we keep the men and women distributions as they are, then a model using a single threshold is going to be unfair for male versus female and violate demographic parity. Hence, the idea is to combine the two distributions by averaging the scores of the same percentile without losing the ranking information. This method can be extended to more than two sensitive groups.

A more recent system called Cappucin [113] repairs data such that a new causality-based fairness called interventional fairness is satisfied. The key insight is that satisfying interventional fairness can be reduced to satisfying multivalued functional dependencies (MVDs). The authors then propose minimal repair methods for MVDs by reducing the problem to MaxSAT or matrix factorization problems.

Fair Data Generation If there is not enough data to satisfy fairness, an alternative is to generate new data using the available data. A recent method [36] is to generate unbiased data using weak supervision. The input is biased data and an unbiased data that is smaller than the biased data that we have some control on. The idea

is to train a generative model on the bias data except that we are adjusting the example weights such that it is as if the generative model is being trained on unbiased data. Then the generative model generates new data that is unbiased. An example weight reflects how likely the example is part of the biased or unbiased data and can be computed by training a separate classifier for distinguishing the biased data from the unbiased data. The generative model that uses the example weights for training is guaranteed to produce unbiased synthetic data. In addition, GANs [135] have also been used for data generation where a generator competes with two discriminators: one for telling apart real and fake data and another for predicting the sensitive attribute.

Fair Data Acquisition As data is increasingly available, acquiring data from external data sources is also becoming a viable option [35, 16]. A recent approach called Slice Tuner [126] selectively acquires examples with the purpose of maximizing both accuracy and fairness of the trained model (Figure 23). Slice Tuner assumes a set of non-overlapping data slices (e.g., regions), and the fairness measure is equal error rates [131] where the model’s accuracies on different slices must be similar. The key idea is to maintain learning curves of slices, which can be used to predict accuracies on those slices given more data. Slice Tuner then solves a convex optimization problem to determine the amount of data to acquire per slice. Two challenges are that learning curves may be unreliable and that acquiring data for one slice may influence the model’s accuracy on another slice. Slice Tuner solves these problems by iteratively updating the learning curves using a proxy for estimating influence. As a result, a model can obtain better accuracy and fairness compared to various baselines given a fixed budget for data acquisition. Another system called Deepdiver [16] performs data acquisition such that all possible slices contain sufficient amounts of data. Here the slices may overlap with each other, and the objective is to guarantee minimum coverage instead of improving model accuracy or fairness.

5.3.2 In-processing Mitigation

We now cover representative in-processing techniques for unfairness mitigation where the model training is fixed. The advantage is that one can directly optimize accuracy and fairness. On the other hand, the downside is that the model training itself may have to change significantly, which may not be feasible in many applications. There are largely three in-processing approaches as shown in Figure 24. The first is to directly modify the objective function of the model training by adding

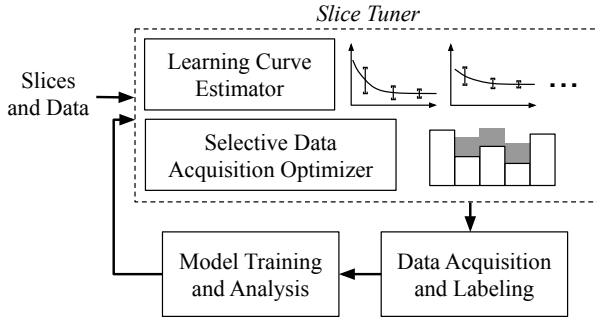


Fig. 23: Slice Tuner [126] is a selective acquisition framework for accurate and fair models where it iteratively estimates learning curves to determine how much data to acquire per data slice.

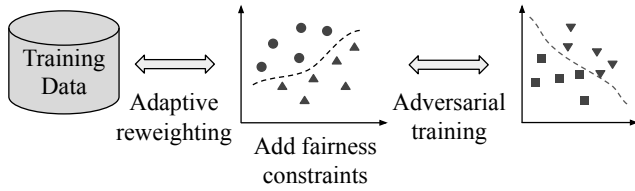


Fig. 24: In-processing: training on biased data. There are three approaches: adding fairness constraints, adversarial training with fair discriminators, and adaptive sample reweighting for fairness.

fairness constraints. The second is to make the model compete with a fairness discriminator via adversarial training. The third is adaptive sample reweighting techniques that re-weight input samples for fairness.

Add Fairness Constraints Directly modifying the model training objective function can be effective in also optimizing for fairness. Zafar et al. propose fairness constraints techniques [141] to use in the objective function of model training to satisfy demographic parity. The focus is on convex margin classifiers like SVMs. However, as the demographic parity constraint is not convex, it cannot be directly added to the objective function. Instead, the idea is to use a proxy that approximates demographic parity and is convex. For the proxy, the authors use the covariance between the sensitive attribute and the signed distance to the decision boundary. Figure 25 provides an intuition why covariance is a good proxy. A limitation of using fairness constraints is that it does not readily generalize to deep neural networks that are not convex. Other optimization techniques [12, 70] for maximizing fairness and accuracy have been proposed as well.

Adversarial Training If one does not want to modify the loss function in the model, another approach is

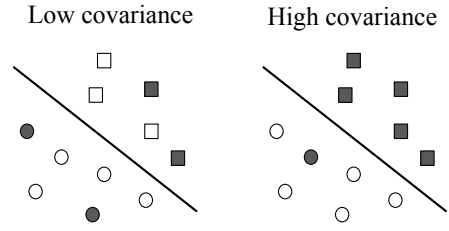


Fig. 25: Suppose we are classifying circles versus squares where the color of the shapes indicate the sensitive group. The left image has little correlation between the sensitive group and being on one side of the decision boundary, which means the covariance is low. The right image, on the other hand, shows a high covariance where just by looking at the sensitive group, one can figure out whether the data point is on which side of the decision boundary. Hence, by minimizing the covariance, one can also make the sensitive attribute more independent of the model predictions and thus satisfy demographic parity as well.

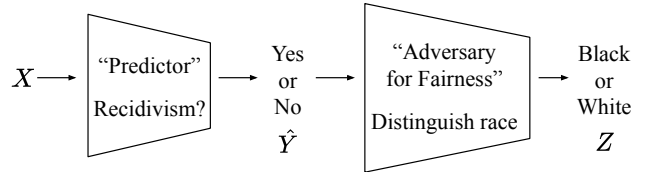


Fig. 26: Adversarial de-biasing [142] competes a classifier with a fair discriminator.

to perform adversarial training with another model for fairness. Adversarial de-biasing [142] is a representative work in this direction. Here the idea is to do adversarial training between a binary classifier and an adversary that tries to infer the sensitive attribute value (Figure 26). For example, the classifier may predict recidivism while the adversary infers the gender of the person based on the classifier predictions. Suppose the fairness measure is demographic parity. A key theoretical result is that, if the adversary optimally predicts the sensitive attribute, but the classifier completely fools the adversary, it means that the model prediction is independent of the sensitive attribute. In our example, recidivism will have nothing to do with the gender. One downside of adversarial debiasing or adversarial training in general is that stability is sometimes an issue where the model training may not easily converge to a single solution.

Adversarial training can also be used to attain both fair and robust training. FR-Train [110] uses a mutual information-based approach to train a model that is both fair and robust. The classifier’s fairness-accuracy tradeoff is harmed when the data is poisoned. FR-train

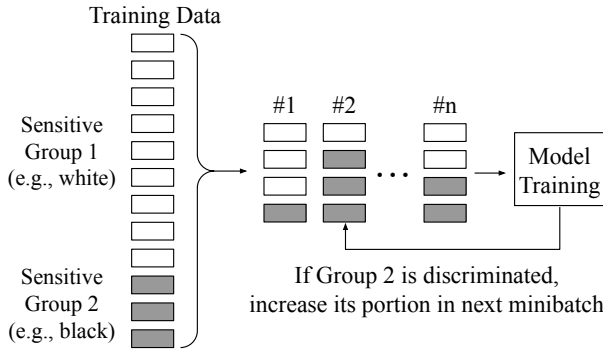


Fig. 27: FairBatch [111] is a batch selection framework for model fairness where sensitive group ratios are adjusted based on intermediate model fairness.

avoids this problem by competing a classifier with two discriminators for fairness and robustness. The robustness discriminator uses a clean validation set that can be constructed using crowdsourcing techniques. We will continue discussing this work in Section 5.4.

Adaptive Reweighting A major downside of the previous methods is that the model training needs to be replaced or modified significantly. A more convenient approach is to only re-weight the samples in order to obtain similar fairness results. FairBatch [111] is a batch selection technique with the purpose of improving fairness. During batch selection, it is common to select a random sample from the training set. Instead, the idea of FairBatch is to adjust the sensitive group ratios within each batch of examples being used for training as illustrated in Figure 27. For example, suppose the training set is biased where a certain sensitive group has very few examples. If an intermediate model shows poor fairness, then the next batch of examples will contain more examples of that sensitive group. How exactly the sensitive group ratio should be adjusted is the technical challenge. OmniFair [143] is a declarative system for supporting group fairness for any model by reweighting samples. While the goals are similar to FairBatch, the specific optimization techniques differ where Omnifair uses a Lagrangian multiplier to translate a constraint optimization problem into an unconstrained optimization problem and leverages a monotonicity property. Other techniques include an adaptive sample reweighting approach that corrects label bias [67] and an adaptive boosting technique for maximizing fairness [66].

5.3.3 Post-processing Mitigation

The final approach for unfairness mitigation is to fix model predictions for fairness (Figure 28). This approach is the only option if the data and model cannot

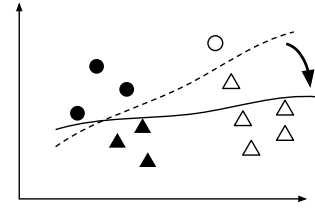


Fig. 28: Post-processing: debiasing a trained model by adjusting the model predictions for fairness.

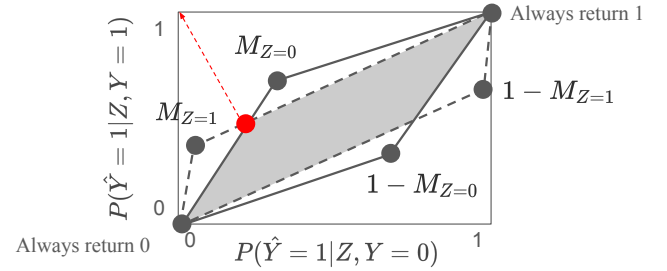


Fig. 29: Post-processing unfairness mitigation [58] involves combining models using randomization to attain the desired fairness.

be modified. However, post-processing usually results in a tradeoff of worse accuracy.

We introduce a representative work [58] that combines models to adjust fairness. The method we explain here assumes equalized odds for binary classifiers, although other settings are supported in the paper as well. We assume a model M for each Z value and then construct the following models: a trivial model that only returns 0, another trivial model that only returns 1, and the “inverted” model $1 - M$, which returns the opposite prediction of M . The idea is to combine these models using randomization such that the fairness criteria is satisfied. Figure 29 illustrates how this combination can be done for $Z = 0$ and $Z = 1$. In each case, we can generate a model with the desired positive prediction rate as long as it is inside the parallelogram. If we generate a model in the intersection of the two parallelograms, we can find a model where $P(\hat{Y} = 1 | Z = 0, Y = A) = P(\hat{Y} = 1 | Z = 1, Y = A)$, $A \in \{0, 1\}$, which is exactly the definition of equalized odds. Among the possible combined models, we then choose the one with the lowest expected loss (i.e., closest to the top left as highlighted in the figure). Other post-processing approaches leverage unlabeled data [39] and calibration [99].

5.4 Convergence with Robustness Techniques

Most recently, we are witnessing a convergence of fairness and robustness techniques. This direction is in-

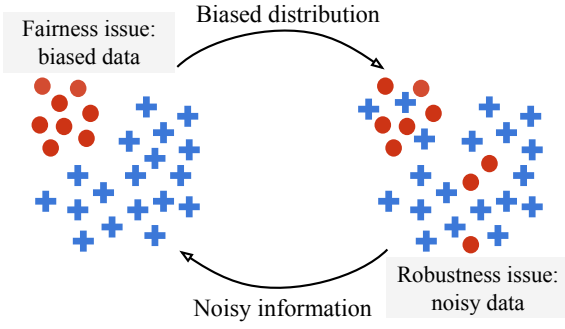


Fig. 30: Fairness and robustness issues may negatively affect each other. Noisy or missing group information may result in inaccurate results after unfairness mitigation. A biased distribution in the data may result in disproportionate accuracies after robust training.

evitable because both techniques address flaws in the data, but one does not subsume the other. Fair training assumes that the data is clean and only focuses on removing its bias. However, the sensitive attribute itself can be noisy or even missing. On the other hand, robust training primarily focuses on improving the overall accuracy, but may result in disproportionate performances between different sensitive groups. Figure 30 illustrates these dynamics. There are three directions for the convergence: making fairness approaches more robust (fairness-oriented), making robust approaches fairer (robust-oriented), and equal mergers of fair and robust training. We summarize the recent research for each of the three approaches.

Fairness-oriented Approaches The first direction of convergence is to make fair training more robust. This research currently has two directions: when the sensitive group information is noisy or entirely missing.

The first scenario may occur if some users may want to hide or mistakenly omit their group memberships. An analysis of fair training results on noisy sensitive group information [133] shows that the true fairness violation on a clean sensitive group can be bounded by a distance between this group and its noisy version. In addition, noise-tolerant fair training techniques [80] have been proposed where the idea is to change the unfairness tolerance to estimate the fairness of the true data distribution.

The second scenario is when the sensitive attribute is fully missing. Here the data collection sometimes does not gather the group information due to various reasons like legal restrictions. Distributionally Robust Optimization (DRO) [121] has been used to improve the model performance for minority sensitive groups without using the group information [59]. The idea is to approximately minimize the worst-case (latent) group

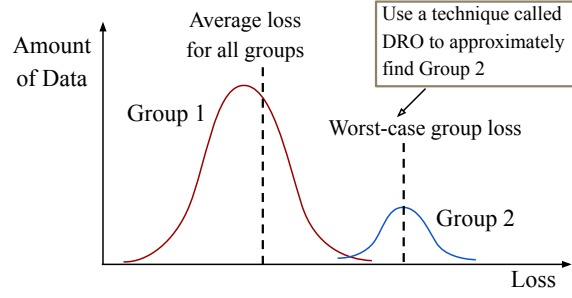


Fig. 31: A DRO-based fair algorithm [59] improves model fairness without using the sensitive group information by identifying the worst-performing samples.

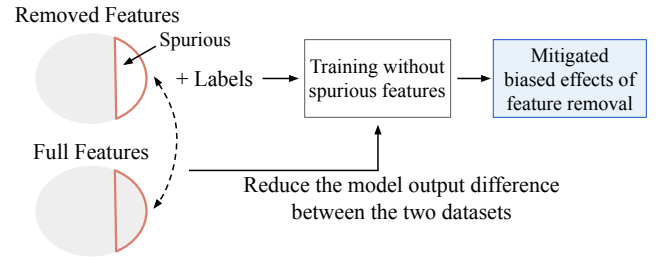


Fig. 32: A self-training technique [72] can mitigate the biased effects of spurious feature removal by also using full-featured data.

loss by identifying the worst-performing samples (Figure 31) and giving them more weight. Adversarially-reweighted learning for fairness [79] makes the assumption that unobserved sensitive attributes are correlated with the features and labels, and performs adversarial training between a classifier versus an adversary that finds less accurate clustered regions and gives more weights on those regions.

Robustness-oriented Approaches Robust training is designed to improve the overall accuracy of a model, but may discriminate groups where some have much worse accuracy than others. There are three directions of research: finding anomalies in the data, training without spurious features, and improving robustness via adversarial training. Fair anomaly detection [144] has been proposed to prevent anomaly detection from discriminating specific groups. The idea is to compete a classifier that finds abnormal data and a discriminator that predicts the sensitive group from the classifier’s prediction. After training, the classifier’s output becomes independent of the sensitive group. Fair training without spurious features [72] addresses the problem of preventing feature removal from being discriminating. A self-training technique is proposed to mitigate accuracy degradation and biased effects (Figure 32). Finally, fair adversarial training [136] prevents adversarial training

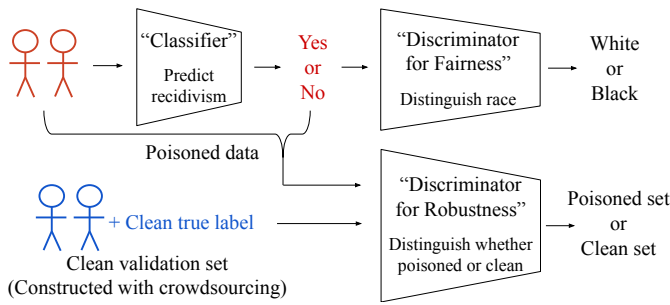


Fig. 33: FR-Train [110] is a mutual information-based approach for achieving both fairness and robustness, which competes one classifier with two discriminators.

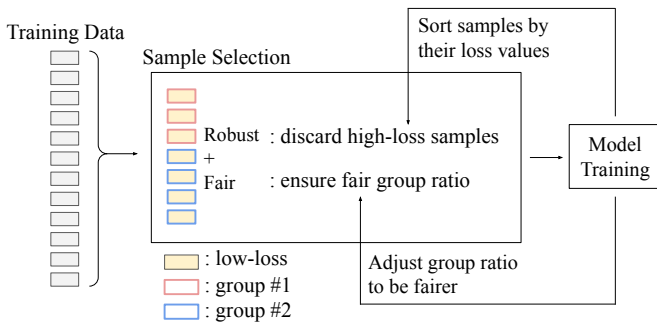


Fig. 34: Adaptive sample selection [112] can be another solution for improving fairness and robustness. The key idea is to utilize only clean and fair samples in training.

from discriminating groups by adding constraints for equalizing accuracy and robustness.

Equal Mergers Robust and fair training can be combined in equal terms as well. One direction is to make the model training fair and robust at the same time. FR-Train [110] is a mutual information-based framework that competes a classifier, discriminator for fairness, and discriminator for robustness to make the classifier fair and robust (Figure 33). A recent sample selection framework [112] adaptively selects training samples for fair and robust model training (Figure 34). This framework does not require modifying the model or leveraging additional clean data. A fairness-aware ERM framework [132] has been proposed based on the observation that group-dependent label noises may reduce both model accuracy and fairness. The solution is to use surrogate loss where the label distribution is corrected based on the noise rates of groups. The surrogate loss better reflects the true loss and thus mitigates the negative effects of group-dependent label noises. Another direction of robust and fair training is to take a role of an adversary and generate attacks that not only reduce accuracy, but also harm fairness. Fairness-targeted poisoning attacks [122] proposes a gradient-based attack

method that finds the optimal attack locations that reduce the fairness the most.

To summarize fair model training, we covered fairness measures, unfairness mitigation techniques, and convergence with robustness techniques. The mitigation can be done before, during, or after model training. Pre-processing is useful when training data can be modified. In-processing is useful when the training algorithm can be modified. Post-processing can be used when we cannot modify the data and model training. The convergence with robustness techniques can be categorized into fair-robust techniques, robust-fair techniques, and equal mergers.

6 Conclusion

As machine learning becomes the new software, collecting data and improving its quality will only become more critical for deep learning. We covered four major topics (data collection, data cleaning and validation, robust model training, and fair model training), which have been studied by different communities, but need to be used together. We believe our survey provides timely guidelines for practitioners in the era of data-centric AI.

References

1. Amazon mechanical turk. <https://www.mturk.com/>. Accessed Jan 15th, 2021.
2. Amazon sagemaker ground truth. <https://aws.amazon.com/sagemaker/groundtruth/>. Accessed Jan 15th, 2021.
3. Data age 2025. <https://www.seagate.com/our-story/data-age-2025/>.
4. Data-centric ai competition. <https://https-deeplearning-ai.github.io/data-centric-comp>.
5. Facets – visualization for ml datasets. <https://pair-code.github.io/facets/>. Accessed Jan 15th, 2021.
6. Gcp ai platform data labeling service. <https://cloud.google.com/ai-platform/data-labeling/docs>. Accessed Jan 15th, 2021.
7. Principles for ai ethics. <https://research.samsung.com/artificial-intelligence>. Accessed Jan 15th, 2021.
8. Responsible ai practices. <https://ai.google/responsibilities/responsible-ai-practices>. Accessed Jan 15th, 2021.
9. Responsible ai principles from microsoft. <https://www.microsoft.com/en-us/ai/responsible-ai>. Accessed Jan 15th, 2021.
10. Software 2.0. <https://medium.com/@karpathy/software-2-0-a64152b37c35>.
11. Trusting ai. <https://www.research.ibm.com/artificial-intelligence/trusted-ai/>. Accessed Jan 15th, 2021.
12. Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna M. Wallach. A reductions approach to fair classification. In *ICML*, pages 60–69, 2018.

13. Pulkit Agrawal, Rajat Arya, Aanchal Bindal, Sandeep Bhatia, Anupriya Gagneja, Joseph Godlewski, Yucheng Low, Timothy Muss, Mudit Manu Paliwal, Sethu Raman, Vishrut Shah, Bochoa Shen, Laura Sugden, Kaiyu Zhao, and Ming-Chuan Wu. Data platform for machine learning. In *SIGMOD*, pages 1803–1816, 2019.
14. Saleema Amershi, Andrew Begel, Christian Bird, Robert DeLine, Harald C. Gall, Ece Kamar, Nachiappan Nagappan, Besmira Nushi, and Thomas Zimmermann. Software engineering for machine learning: a case study. In *ICSE*, pages 291–300, 2019.
15. J. Angwin, J. Larson, S. Mattu, and L. Kirchner. Machine bias: There’s software used across the country to predict future criminals. And its biased against blacks., 2016.
16. Abolfazl Asudeh, Zhongjun Jin, and H. V. Jagadish. Assessing and remedying coverage for a given dataset. In *ICDE*, pages 554–565, 2019.
17. Stephen H. Bach, Daniel Rodriguez, Yintao Liu, Chong Luo, Haidong Shao, Cassandra Xia, Souvik Sen, Alexander Ratner, Braden Hancock, Houman Alborzi, Rahul Kuchhal, Christopher Ré, and Rob Malkin. Snorkel drybell: A case study in deploying weak supervision at industrial scale. In *SIGMOD*, pages 362–375, 2019.
18. Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning*. fairmlbook.org, 2019. <http://www.fairmlbook.org>.
19. Denis Baylor, Eric Breck, Heng-Tze Cheng, Noah Fiedel, Chuan Yu Foo, Zakaria Haque, Salem Haykal, Mustafa Ispir, Vihan Jain, Levent Koc, Chiu Yuen Koo, Lukasz Lew, Clemens Mewald, Akshay Naresh Modi, Neoklis Polyzotis, Sukriti Ramesh, Sudip Roy, Steven Euijong Whang, Martin Wicke, Jarek Wilkiewicz, Xin Zhang, and Martin Zinkevich. TFX: A tensorflow-based production-scale machine learning platform. In *KDD*, pages 1387–1395, 2017.
20. Rachel K. E. Bellamy, Kuntal Dey, Michael Hind, et al. AI fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias. *IBM Journal of Research and Development*, 2019.
21. Yoav Benjamini and Daniel Yekutieli. The control of the false discovery rate in multiple testing under dependency. *The Annals of Statistics*, 29(4):1165 – 1188, 2001.
22. Richard Berk, Hoda Heidari, Shahin Jabbari, Michael Kearns, and Aaron Roth. Fairness in criminal justice risk assessments: The state of the art, 2017.
23. David Berthelot, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Kihyuk Sohn, Han Zhang, and Colin Raffel. Remixmatch: Semi-supervised learning with distribution matching and augmentation anchoring. In *ICLR*, 2020.
24. David Berthelot, Nicholas Carlini, Ian J. Goodfellow, Nicolas Papernot, Avital Oliver, and Colin Raffel. Mixmatch: A holistic approach to semi-supervised learning. In *NeurIPS*, pages 5050–5060, 2019.
25. Battista Biggio, Igino Corona, Davide Maiorca, Blaine Nelson, Nedim Srdic, Pavel Laskov, Giorgio Giacinto, and Fabio Roli. Evasion attacks against machine learning at test time. In *ECML PKDD*, pages 387–402. Springer, 2013.
26. Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In *COLT*, pages 92–100, New York, NY, USA, 1998. ACM.
27. Matthias Boehm, Iulian Antonov, Sebastian Baunsgaard, Mark Dokter, Robert Ginhör, Kevin Innerebner, Florian Klezin, Stefanie N. Lindstaedt, Arnab Phani, Benjamin Rath, Berthold Reinwald, Shafaq Siddiqui, and Sebastian Benjamin Wrede. Systemds: A declarative machine learning system for the end-to-end data science lifecycle. In *CIDR*, 2020.
28. Eric Breck, Martin Zinkevich, Neoklis Polyzotis, Steven Whang, and Sudip Roy. Data validation for machine learning. In *MLSys*, 2019.
29. Dan Brickley, Matthew Burgess, and Natasha F. Noy. Google dataset search: Building a search engine for datasets in an open web ecosystem. In *WWW*, pages 1365–1375, 2019.
30. Michael J. Cafarella, Alon Y. Halevy, Hongrae Lee, Jayant Madhavan, Cong Yu, Daisy Zhe Wang, and Eugene Wu. Ten years of webtables. *PVLDB*, 11(12):2140–2149, 2018.
31. Anirban Chakraborty, Manaar Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopadhyay. Adversarial attacks and defences: A survey. *CoRR*, abs/1810.00069, 2018.
32. Haw-Shiuan Chang, Erik G. Learned-Miller, and Andrew McCallum. Active bias: Training more accurate neural networks by emphasizing high variance samples. In *NeurIPS*, pages 1002–1012, 2017.
33. Zhengping Che, Sanjay Purushotham, Kyunghyun Cho, David Sontag, and Yan Liu. Recurrent neural networks for multivariate time series with missing values. *Nature Scientific Reports*, 8(1):6085, 2018.
34. Andrew Chen, Andy Chow, Aaron Davidson, Arjun DCunha, Ali Ghodsi, Sue Ann Hong, Andy Konwinski, Clemens Mewald, Siddharth Murching, Tomas Nykodym, Paul Ogilvie, Mani Parkhe, Avesh Singh, Fen Xie, Matei Zaharia, Richard Zang, Juntao Zheng, and Corey Zumar. Developments in mlflow: A system to accelerate the machine learning lifecycle. In *DEEM@SIGMOD*, pages 5:1–5:4, 2020.
35. Irene Y. Chen, Fredrik D. Johansson, and David A. Sontag. Why is my classifier discriminatory? In *NeurIPS*, pages 3543–3554, 2018.
36. Kristy Choi, Aditya Grover, Trisha Singh, Rui Shu, and Stefano Ermon. Fair generative modeling via weak supervision. In *ICML*, pages 1887–1898, 2020.
37. Alexandra Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big Data*, 5(2):153–163, 2017.
38. Alexandra Chouldechova and Aaron Roth. A snapshot of the frontiers of fairness in machine learning. *Commun. ACM*, 63(5):82–89, 2020.
39. Evgenii Chzhenn, Christophe Denis, Mohamed Hebiri, Luca Oneto, and Massimiliano Pontil. Leveraging labeled and unlabeled data for consistent fair binary classification. In *NeurIPS*, pages 12739–12750, 2019.
40. Andrew Cotter, Heinrich Jiang, and Karthik Sridharan. Two-player games for efficient non-convex constrained optimization. In *ALT*, pages 300–332, 2019.
41. Gabriela F. Cretu, Angelos Stavrou, Michael E. Locasto, Salvatore J. Stolfo, and Angelos D. Keromytis. Casting out demons: Sanitizing training data for anomaly sensors. In *IEEE S&P*, pages 81–95, 2008.
42. Ekin D. Cubuk, Barret Zoph, Dandelion Mané, Vijay Vasudevan, and Quoc V. Le. Autoaugment: Learning augmentation strategies from data. In *CVPR*, pages 113–123, 2019.
43. William Dieterich, Christina Mendoza, and Tim Brennan. Compas risk scales: Demonstrating accuracy equity and predictive parity. Technical report, Northpoint Inc, 2016.

44. Mohamad Dolatshah, Mathew Teoh, Jiannan Wang, and Jian Pei. Cleaning crowdsourced labels using oracles for statistical classification. *PVLDB*, 12(4):376–389, 2018.
45. Xin Luna Dong and Theodoros Rekatsinas. Data integration and machine learning: A natural synergy. In *KDD*, pages 3193–3194, 2019.
46. Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard S. Zemel. Fairness through awareness. In *ITCS*, pages 214–226, 2012.
47. Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and removing disparate impact. In *KDD*, pages 259–268, 2015.
48. Raul Castro Fernandez, Ziawasch Abedjan, Famién Koko, Gina Yuan, Samuel Madden, and Michael Stonebraker. Aurum: A data discovery system. In *ICDE*, pages 1001–1012, 2018.
49. Dean P. Foster and Robert A. Stine. Alpha-investing: a procedure for sequential control of expected false discoveries. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(2):429–444, 2008.
50. Karan Goel, Albert Gu, Yixuan Li, and Christopher Ré. Model patching: Closing the subgroup performance gap with data augmentation. In *ICLR*, 2021.
51. Ian J. Goodfellow. NIPS 2016 tutorial: Generative adversarial networks. *CoRR*, abs/1701.00160, 2017.
52. Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014.
53. Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015.
54. J. Gordon. Introducing tensorflow hub: A library for reusable machine learning modules in tensorflow., 2018.
55. Stefan Grafberger, Julia Stoyanovich, and Sebastian Schelter. Lightweight inspection of data preprocessing in native machine learning pipelines. In *CIDR*, 2021.
56. Alon Y. Halevy, Flip Korn, Natalya Fridman Noy, Christopher Olston, Neoklis Polyzotis, Sudip Roy, and Steven Euijong Whang. Goods: Organizing google’s datasets. In *SIGMOD*, pages 795–806, 2016.
57. Bo Han, Quanming Yao, Xingrui Yu, Gang Niu, Miao Xu, Weihua Hu, Ivor W. Tsang, and Masashi Sugiyama. Co-teaching: Robust training of deep neural networks with extremely noisy labels. In *NeurIPS*, pages 8536–8546, 2018.
58. Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In *NIPS*, pages 3315–3323, 2016.
59. Tatsunori B. Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. Fairness without demographics in repeated loss minimization. In Jennifer G. Dy and Andreas Krause, editors, *ICML*, volume 80, pages 1934–1943. PMLR, 2018.
60. Kim M. Hazelwood, Sarah Bird, David M. Brooks, Soumith Chintala, Utku Diril, Dmytro Dzhulgakov, Mohamed Fawzy, Bill Jia, Yangqing Jia, Aditya Kalro, James Law, Kevin Lee, Jason Lu, Pieter Noordhuis, Misha Smelyanskiy, Liang Xiong, and Xiaodong Wang. Applied machine learning at facebook: A datacenter infrastructure perspective. In *HPCA*, pages 620–629, 2018.
61. Dan Hendrycks, Norman Mu, Ekin Dogus Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. In *ICLR*, 2020.
62. Geon Heo, Yuji Roh, Seonghyeon Hwang, Dayun Lee, and Steven Euijong Whang. Inspector gadget: A data programming-based labeling system for industrial images. *PVLDB*, 2021.
63. M. Hermann, J. and Del Baso. Meet michelangelo: Uber’s machine learning platform, 2017.
64. Victoria J. Hodge and Jim Austin. A survey of outlier detection methodologies. *Artif. Intell. Rev.*, 22(2):85–126, 2004.
65. Ihab F. Ilyas and Xu Chu. *Data Cleaning*. ACM, 2019.
66. Vasileios Iosifidis and Eirini Ntoutsi. Adafair: Cumulative fairness adaptive boosting. In *CIKM*, pages 781–790, 2019.
67. Heinrich Jiang and Ofir Nachum. Identifying and correcting label bias in machine learning. In *AISTATS*, pages 702–712, 2020.
68. Lu Jiang, Zhengyuan Zhou, Thomas Leung, Li-Jia Li, and Li Fei-Fei. Mentornet: Learning data-driven curriculum for very deep neural networks on corrupted labels. In *ICML*, pages 2309–2318, 2018.
69. Faisal Kamiran and Toon Calders. Data preprocessing techniques for classification without discrimination. *Knowl. Inf. Syst.*, 33(1):1–33, 2011.
70. Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. Fairness-aware classifier with prejudice remover regularizer. In *ECML PKDD*, pages 35–50, 2012.
71. Aria Khademi, Sanghack Lee, David Foley, and Vasant Honavar. Fairness in algorithmic decision making: An excursion through the lens of causality. In *WWW*, pages 2907–2914, 2019.
72. Fereshte Khani and Percy Liang. Removing spurious features can hurt accuracy and affect groups disproportionately. In *FAccT*, pages 196–205. ACM, 2021.
73. Niki Kilbertus, Mateo Rojas-Carulla, Giambattista Parascandolo, Moritz Hardt, Dominik Janzing, and Bernhard Schölkopf. Avoiding discrimination through causal reasoning. In *NeurIPS*, pages 656–666, 2017.
74. Hoon Kim, Kangwook Lee, Gyeongjo Hwang, and Changho Suh. Crash to not crash: Learn to identify dangerous vehicles using a simulator. In *AAAI*, pages 978–985, 2019.
75. Pang Wei Koh, Jacob Steinhardt, and Percy Liang. Stronger data poisoning attacks break data sanitization defenses. *CoRR*, abs/1811.00741, 2018.
76. Sanjay Krishnan, Jiannan Wang, Eugene Wu, Michael J. Franklin, and Ken Goldberg. Activeclean: Interactive data cleaning for statistical modeling. *PVLDB*, 9(12):948–959, 2016.
77. Karol Kurach, Mario Lucic, Xiaohua Zhai, Marcin Michalski, and Sylvain Gelly. The GAN landscape: Losses, architectures, regularization, and normalization. *CoRR*, abs/1807.04720, 2018.
78. Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In *NeurIPS*, pages 4066–4076. 2017.
79. Preethi Lahoti, Alex Beutel, Jilin Chen, Kang Lee, Flavien Prost, Nithum Thain, Xuezhi Wang, and Ed Chi. Fairness without demographics through adversarially reweighted learning. In *NeurIPS*, 2020.
80. Alexandre Louis Lamy and Ziyuan Zhong. Noise-tolerant fair classification. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett, editors, *NeurIPS*, pages 294–305, 2019.
81. Doris Jung Lin Lee and Aditya G. Parameswaran. The case for a visual discovery assistant: A holistic solution for accelerating visual data exploration. *IEEE Data Eng. Bull.*, 41(3):3–14, 2018.

82. Jae-Gil Lee, Yuji Roh, Hwanjun Song, and Steven Euijong Whang. Machine learning robustness, fairness, and their convergence. In *KDD*, pages 4046–4047, 2021.
83. Junnan Li, Richard Socher, and Steven C. H. Hoi. Dividemix: Learning with noisy labels as semi-supervised learning. In *ICLR*, 2020.
84. Peng Li, Xi Rao, Jennifer Blase, Yue Zhang, Xu Chu, and Ce Zhang. CleanML: A benchmark for joint data cleaning and machine learning. *CoRR*, abs/1904.09483, 2019.
85. Peng Li, Xi Rao, Jennifer Blase, Yue Zhang, Xu Chu, and Ce Zhang. Cleanml: A benchmark for joint data cleaning and machine learning [experiments and analysis]. In *ICDE*, 2021.
86. Eran Malach and Shai Shalev-Shwartz. Decoupling “when to update” from “how to update”. In *NIPS*, pages 960–970, 2017.
87. Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *CoRR*, abs/1908.09635, 2019.
88. Leonel Aguilar Melgar, David Dao, Shaoduo Gan, Nezihe Merve Gürel, Nora Hollenstein, Jiawei Jiang, Bojan Karlas, Thomas Lemmin, Tian Li, Yang Li, Xi Rao, Johannes Rausch, Cédric Renggli, Luka Rimanic, Maurice Weber, Shuai Zhang, Zhikuan Zhao, Kevin Schawinski, Wentao Wu, and Ce Zhang. Ease.ml: A lifecycle management system for machine learning. In *CIDR*, 2021.
89. Dongyu Meng and Hao Chen. Magnet: A two-pronged defense against adversarial examples. In Bhavani M. Thuraisingham, David Evans, Tal Malkin, and Dongyan Xu, editors, *ACM SIGSAC*, pages 135–147, 2017.
90. Jan Hendrik Metzen, Tim Genewein, Volker Fischer, and Bastian Bischoff. On detecting adversarial perturbations. In *ICLR*, 2017.
91. Razieh Nabi and Ilya Shpitser. Fair inference on outcomes. In *AAAI*, pages 1931–1940, 2018.
92. Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *ECCV*, pages 69–84, 2016.
93. Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE TKDE*, 22(10):1345–1359, 2010.
94. Nicolas Papernot, Patrick D. McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In *IEEE SP*, pages 582–597, 2016.
95. Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *NeurIPS*, pages 8024–8035. Curran Associates, Inc., 2019.
96. Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making deep neural networks robust to label noise: A loss correction approach. In *CVPR*, pages 2233–2241, 2017.
97. Andrea Paudice, Luis Muñoz-González, András György, and Emil C. Lupu. Detection of adversarial training examples in poisoning attacks through anomaly detection. *CoRR*, abs/1802.03041, 2018.
98. Nikos Pelekis, Christos Ntrigkogiass, Panagiotis Tampakis, Stylianos Sideridis, and Yannis Theodoridis. Her-moupolis: A trajectory generator for simulating generalized mobility patterns. In *ECML PKDD*, pages 659–662, 2013.
99. Geoff Pleiss, Manish Raghavan, Felix Wu, Jon M. Kleinberg, and Kilian Q. Weinberger. On fairness and calibration. In *NIPS*, pages 5680–5689, 2017.
100. Neoklis Polyzotis, Sudip Roy, Steven Euijong Whang, and Martin Zinkevich. Data management challenges in production machine learning. In *SIGMOD*, pages 1723–1726, 2017.
101. Neoklis Polyzotis, Sudip Roy, Steven Euijong Whang, and Martin Zinkevich. Data lifecycle challenges in production machine learning: A survey. *SIGMOD Rec.*, 47(2):17–28, 2018.
102. Alexander Ratner, Stephen H. Bach, Henry Ehrenberg, Jason Fries, Sen Wu, and Christopher Ré. Snorkel: Rapid training data creation with weak supervision. *PVLDB*, 11(3):269–282, November 2017.
103. Alexander Ratner, Stephen H. Bach, Henry R. Ehrenberg, Jason A. Fries, Sen Wu, and Christopher Ré. Snorkel: rapid training data creation with weak supervision. *VLDB J.*, 29(2-3):709–730, 2020.
104. Alexander J. Ratner, Henry R. Ehrenberg, Zeshan Hussain, Jared Dunnmon, and Christopher Ré. Learning to compose domain-specific transformations for data augmentation. In *NIPS*, pages 3239–3249, 2017.
105. Sergey Redyuk, Zoi Kaoudi, Volker Markl, and Sebastian Schelter. Automating data quality validation for dynamic data ingestion. In *EDBT*, pages 61–72, 2021.
106. Scott E. Reed, Honglak Lee, Dragomir Anguelov, Christian Szegedy, Dumitru Erhan, and Andrew Rabinovich. Training deep neural networks on noisy labels with bootstrapping. In *ICLR*, 2015.
107. Theodoros Rekatsinas, Xu Chu, Ihab F. Ilyas, and Christopher Ré. Holoclean: Holistic data repairs with probabilistic inference. *PVLDB*, 10(11):1190–1201, 2017.
108. Francesco Ricci, Lior Rokach, and Bracha Shapira, editors. *Recommender Systems Handbook*. Springer, 2015.
109. Yuji Roh, Geon Heo, and Steven Euijong Whang. A survey on data collection for machine learning: a big data - AI integration perspective. *IEEE TKDE*, 2019.
110. Yuji Roh, Kangwook Lee, Steven Euijong Whang, and Changho Suh. FR-Train: A mutual information-based approach to fair and robust training. In *ICML*, 2020.
111. Yuji Roh, Kangwook Lee, Steven Euijong Whang, and Changho Suh. Fairbatch: Batch selection for model fairness. In *ICLR*. OpenReview.net, 2021.
112. Yuji Roh, Kangwook Lee, Steven Euijong Whang, and Changho Suh. Sample selection for fair and robust training. In *NeurIPS*, 2021.
113. Babak Salimi, Luke Rodriguez, Bill Howe, and Dan Suciu. Interventional fairness: Causal database repair for algorithmic fairness. In *SIGMOD*, pages 793–810, 2019.
114. S. Schelter, Joos-Hendrik Böse, Johannes Kirschnick, T. Klein, and Stephan Seufert. Automatically tracking metadata and provenance of machine learning experiments. In *Workshop on ML Systems at NIPS*, 2017.
115. Sebastian Schelter, Dustin Lange, Philipp Schmidt, Meltem Celikel, Felix Bießmann, and Andreas Grafberger. Automating large-scale data quality verification. *Proc. VLDB Endow.*, 11(12):1781–1794, 2018.
116. Sebastian Schelter, Tammo Rukat, and Felix Biessmann. JENGA - A framework to study the impact of data errors on the predictions of machine learning models. In *EDBT*, pages 529–534, 2021.

117. D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips, Dietmar Ebner, Vinay Chaudhary, Michael Young, Jean-François Crespo, and Dan Dennison. Hidden technical debt in machine learning systems. In *NIPS*, pages 2503–2511, 2015.
118. Burr Settles. *Active Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2012.
119. Ali Shafahi, W. Ronny Huang, Mahyar Najibi, Octavian Suci, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. In *NeurIPS*, pages 6106–6116, 2018.
120. Victor S. Sheng, Foster J. Provost, and Panagiotis G. Ipeirotis. Get another label? improving data quality and data mining using multiple, noisy labelers. In *KDD*, pages 614–622, 2008.
121. Aman Sinha, Hongseok Namkoong, and John C. Duchi. Certifying some distributional robustness with principled adversarial training. In *ICLR*, 2018.
122. David Solans, Battista Biggio, and Carlos Castillo. Poisoning attacks on algorithmic fairness. In Frank Hutter, Kristian Kersting, Jefrey Lijffijt, and Isabel Valera, editors, *ECML PKDD*, volume 12457, pages 162–177. Springer, 2020.
123. Hwanjun Song, Minseok Kim, and Jae-Gil Lee. SELFIE: refurbishing unclean samples for robust deep learning. In *ICML*, pages 5907–5915, 2019.
124. Hwanjun Song, Minseok Kim, Dongmin Park, Yooju Shin, and Jae-Gil Lee. Robust learning by self-transition for handling noisy labels. In *KDD*, pages 1490–1500, 2021.
125. Michael Stonebraker and El Kindi Rezig. Machine learning and big data: What is important? *IEEE Data Eng. Bull.*, 2019.
126. Ki Hyun Tae and Steven Euijong Whang. Slice tuner: A selective data acquisition framework for accurate and fair machine learning models. In *SIGMOD*, pages 1771–1783. ACM, 2021.
127. Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In *NIPS*, pages 1195–1204, 2017.
128. Ignacio G. Terrizzano, Peter M. Schwarz, Mary Roth, and John E. Colino. Data wrangling: The challenging journey from the wild to the lake. In *CIDR*, 2015.
129. Isaac Triguero, Salvador García, and Francisco Herrera. Self-labeled techniques for semi-supervised learning: taxonomy, software and empirical study. *Knowl. Inf. Syst.*, 42(2):245–284, 2015.
130. Manasi Vartak, Sajjadur Rahman, Samuel Madden, Aditya G. Parameswaran, and Neoklis Polyzotis. SEEDB: Efficient data-driven visualization recommendations to support visual analytics. *PVLDB*, 8(13):2182–2193, 2015.
131. Suresh Venkatasubramanian. Algorithmic fairness: Measures, methods and representations. In *PODS*, page 481, 2019.
132. Jialu Wang, Yang Liu, and Caleb Levy. Fair classification with group-dependent label noise. In Madeleine Clare Elish, William Isaac, and Richard S. Zemel, editors, *FAccT*, pages 526–536. ACM, 2021.
133. Serena Wang, Wenshuo Guo, Harikrishna Narasimhan, Andrew Cotter, Maya R. Gupta, and Michael I. Jordan. Robust optimization for fairness with noisy protected groups. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *NeurIPS*, 2020.
134. Steven Euijong Whang and Jae-Gil Lee. Data collection and quality challenges for deep learning. *Proc. VLDB Endow.*, 13(12):3429–3432, 2020.
135. Depeng Xu, Shuhan Yuan, Lu Zhang, and Xintao Wu. Fairgan: Fairness-aware generative adversarial networks. In *IEEE Big Data*, pages 570–575, 2018.
136. Han Xu, Xiaorui Liu, Yaxin Li, Anil K. Jain, and Jiliang Tang. To be robust or to be fair: Towards fairness in adversarial training. In Marina Meila and Tong Zhang, editors, *ICML*, volume 139, pages 11492–11501. PMLR, 2021.
137. Weilin Xu, David Evans, and Yanjun Qi. Feature squeezing: Detecting adversarial examples in deep neural networks. In *NDSS*, 2018.
138. David Yarowsky. Unsupervised word sense disambiguation rivaling supervised methods. In *ACL*, pages 189–196, Stroudsburg, PA, USA, 1995. Association for Computational Linguistics.
139. Sangdo Yun, Dongyoon Han, Sanghyuk Chun, Seong Joon Oh, Youngjoon Yoo, and Junsuk Choe. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *ICCV*, pages 6022–6031, 2019.
140. Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez-Rodriguez, and Krishna P. Gummadi. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *WWW*, pages 1171–1180. ACM, 2017.
141. Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez-Rodriguez, and Krishna P. Gummadi. Fairness constraints: Mechanisms for fair classification. In *AISTATS*, pages 962–970, 2017.
142. Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial learning. In *AIES*, pages 335–340, 2018.
143. Hantian Zhang, Xu Chu, Abolfazl Asudeh, and Shamkant B. Navathe. Omnifair: A declarative system for model-agnostic group fairness in machine learning. In *SIGMOD*, pages 2076–2088, 2021.
144. Hongjing Zhang and Ian Davidson. Facct. pages 138–148. ACM, 2021.
145. Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *ICLR*, 2018.
146. Junzhe Zhang and Elias Bareinboim. Fairness in decision-making - the causal explanation formula. In *AAAI*, 2018.
147. Yi Zhang and Zachary G. Ives. Finding related tables in data lakes for interactive data science. In *SIGMOD*, pages 1951–1966, 2020.
148. Zheguang Zhao, Lorenzo De Stefani, Emanuel Zraggen, Carsten Binnig, Eli Upfal, and Tim Kraska. Controlling false discoveries during interactive data exploration. In *SIGMOD*, pages 527–540, 2017.
149. Yan Zhou and Sally A. Goldman. Democratic co-learning. In *IEEE ICTAI*, pages 594–602, 2004.
150. Zhi-Hua Zhou and Ming Li. Tri-training: Exploiting unlabeled data using three classifiers. *IEEE TKDE*, 17(11):1529–1541, November 2005.
151. Chen Zhu, W. Ronny Huang, Hengduo Li, Gavin Taylor, Christoph Studer, and Tom Goldstein. Transferable clean-label poisoning attacks on deep neural nets. In *ICML*, pages 7614–7623, 2019.
152. Xiaojin Zhu. Semi-Supervised Learning Literature Survey. Technical report, Computer Sciences, University of Wisconsin-Madison, 2005.