
How Data-Related AI Research can Support Technical Solutions for Regulatory Compliance

Danilo Brajovic^{1,2} David A. Kreplin¹ Marco F. Huber^{1,2}

¹Fraunhofer IPA ²Stuttgart University

danilo.brajovic@ipa.fraunhofer.de

david.kreplin@hs-heilbronn.de marco.huber@ieee.org

Abstract

Ensuring high-quality, representative, and secure datasets is crucial for compliance with emerging regulatory frameworks, such as the European AI Act. In this paper, we survey five data-centric challenges: inherent data quality, application-specific data suitability, data sufficiency, dataset variability, and data security, and link each to both established and emerging methods and research. We then propose a workflow that aligns best practices from machine learning (ML) research with regulatory requirements, showing how each step can be operationalized to meet the “relevant, representative, error-free” criteria. Our analysis highlights key opportunities for regulators to refine their mandates and for ML researchers to conduct follow-up research.

1 Introduction

High-risk AI systems depend critically on training data that is “relevant, representative, error-free, and complete” (European AI Act Art. 10 (3)). Yet, practitioners often satisfy these mandates with coarse governance checklists and simple metrics (e.g., label-error rate, missing-value fraction), leaving a gap between regulatory intent and technical practice. Concurrently, AI research has produced methods such as data attribution that can directly address legal requirements but remain underused in compliance workflows. In this paper, we bridge that gap by:

1. Proposing a regulation-aligned taxonomy of five ML data challenges (inherent data quality; application-specific data suitability; sufficiency; variability; and security),
2. Mapping established and emerging methods to this taxonomy, and
3. Sketching an exemplary workflow for “regulatable ML.” Our analysis highlights underutilized areas, especially data attribution and valuation, that can address all five challenges and provide strong motivation for follow-up research.

2 Related Work

Data quality in ML encompasses regulatory requirements, documentation, quantitative metrics, and other methodologies. We review major contributions in each area and highlight gaps.

Regulatory Frameworks The European AI Act [18] mandates that training, validation, and test sets be “relevant, representative, error-free, and complete” (Art. 10 (3)), but leaves the technical implementation open. These criteria are further detailed in standards, such as those from the ISO/IEC JTC 1/SC 42 committee for AI. Most notably ISO/IEC 5259 [35] specifies 24 data-quality metrics

for analytics and ML. ISO/IEC 5259 organizes these metrics along two dimensions (inherent and context-dependent data quality) and operationalizes the notions of relevance, representativeness, accuracy, and completeness, albeit at a relatively high level of abstraction. Definitions of each metric are provided in Appendix C.

Data Documentation and Transparency Standardized documentation frameworks, such as Data Cards [65] and Datasheets for Datasets [26], record a dataset’s origin, composition, labeling process, and limitations. While these efforts support regulatory compliance by improving transparency, they do not engage directly with legal mandates or prescribe quantitative metrics for quality.

Data Quality Metrics Surveys on data quality metrics [95, 11, 55, 29, 64] collate measures, often aligned with ISO/IEC 5259, that quantify dimensions such as completeness, consistency, and timeliness. Although taxonomies vary, a common dichotomy separates inherent and context-dependent dimensions. Few works explicitly tie metric selection to regulatory requirements, leaving a gap between theoretical measures and compliance.

Data-Centric AI Summaries and Recommendations Jakubik et al. [36] distinguish between dataset extension (collection) and refinement (quality improvement), illustrating commercial tooling for each. Zha et al. [90] organize the data lifecycle into training-data development, inference-data development, and maintenance, providing resources and tools. Hammoudeh and Lowd [34] survey methods for training-data influence with a technical focus, and Yu et al. [89] review dataset distillation techniques, highlighting applications in continual and federated learning, privacy, and robustness. Practitioner-oriented recommendations by Lones [48], Orr and Crawford [61], and Zhao et al. [94] advocate for dataset diversification, rigorous quality checks, and transparent documentation. However, all these works do not consider regulatory criteria explicitly.

2.1 Generative AI and its Impact

Generative AI (GenAI) can both improve existing data practices for ML and be influenced by the data itself, potentially exacerbating existing challenges. GenAI can enhance data practices by augmenting scarce datasets (e.g., rare disease cases), improving fairness (e.g., by generating diverse synthetic profiles), and simulating variability (e.g., edge cases for autonomous systems) [50, 5]. However, this can also introduce artifacts, amplify biases, or degrade performance by erasing real-world distribution nuances, leading to errors and outputs that undermine safety in applications like finance or healthcare [73, 8]. Additionally, GenAI’s reliance on vast, often opaque pretraining corpora raises ethical and legal concerns [88, 49].

Due to the immense impact of GenAI, regulators have felt the necessity to separately regulate it, resulting in specific rules for General Purpose AI, which includes GenAI. For such systems, the AI Act [18] additionally requires documented training sources (Art. 53(d)) and copyright compliance policies (Art. 53(c)).

2.2 Interim Conclusion

While prior work addresses legal mandates, documentation standards, metrics, and methodologies, these threads remain largely siloed. Moreover, the rapid emergence of GenAI both amplifies these challenges, by introducing novel data-generation and usage patterns, while it depends critically on data quality itself. In the next section, we introduce key terminology to bridge this gap.

3 Taxonomy of ML Data Challenges

Both ML research and regulatory texts use varied, and often opaque, terminology to describe data challenges, even when discussing the same issues. For example, the EU AI Act speaks of *completeness*, which can refer both to domain coverage and to the completeness of feature sets. Similarly, terms like *appropriateness* and *representativeness* appear interchangeably across studies to describe a dataset’s fitness for its intended use case.

To eliminate this ambiguity, we propose a unified taxonomy of five categories: INHERENT DATA QUALITY, APPLICATION-SPECIFIC DATA SUITABILITY, DATA SUFFICIENCY, DATA VARIABILITY,

CHALLENGE	EXAMPLES
INHERENT DQ	Annotation errors; duplicate records; missing values; ambiguous labels
SUITABILITY	Mismatched distributions; class imbalance; inappropriate train/test splits
SUFFICIENCY	Insufficient data; excessive data leading to high cost or privacy concerns
VARIABILITY	Temporal drift; evolving user behavior or data-collection methods
SECURITY	Data breaches; unauthorized access or modifications; data poisoning

Table 1: Taxonomy of data challenges in ML, with examples.

and DATA SECURITY. Table 1 summarizes each category with examples. Where existing terms are precise and unambiguous, we adopt them; otherwise, we introduce new labels to ensure clarity.

Inherent Data Quality Inherent data quality encompasses all metrics that can be computed with respect to the available data without any further knowledge. Key challenges include ensuring label and annotation accuracy, identifying and correcting noise and errors such as incorrect values, inconsistencies, and duplicate data [59, 10]. These issues can degrade model accuracy, skew the training process, and inflate dataset size.

Relationship to other sources: ISO/IEC 5259 [35] defines inherent data quality similarly. GDPR Article 5(1)(d) requires that personal data be accurate, and the AI Act (Art. 10) mandates error-free datasets.

Application-Specific Data Suitability Application-specific data suitability refers to the degree to which a dataset aligns with the specific conditions of its intended deployment environment. Challenges include when training data diverges from operational reality, for example, through biased sampling that creates class imbalances [12, 54] or through data-splits that fail to mirror the true input distribution. Such mismatches rest on assumptions about real-world conditions.

Relationship to other sources: The AI Act (Art. 10) emphasizes that datasets must be *representative*. Other works use terms like “appropriateness” to capture similar concerns.

Data Sufficiency Data sufficiency addresses whether a dataset contains an adequate, but not excessive, amount of information. Too little data can hamper model generalization, especially for rare events or classes; too much data can raise costs, latency, and privacy risks. Furthermore, acquiring and labeling data often entails significant expense. Techniques like transfer learning [96], and active learning [74] maximize data utility.

Relationship to other sources: GDPR Article 5(1)(c) mandates the principle of “data minimization,” requiring that personal data be “adequate, relevant and limited to what is necessary” for the purposes for which they are processed. By analogy, the EU AI Act mandates that training, validation, and testing datasets be “relevant” to the AI system’s intended use.

Data Variability Data variability describes changes in data over time, such as evolving user behaviors, shifting environments, or updated collection methods, causing drift between training and real-world data. Static models may fail in dynamic settings, reducing reliability. Drift detection [6] and adaptive methods address this.

Relationship to other sources: AI Act Article 15 requires monitoring for high-risk systems to address drift [18]. GDPR Article 5(1)(d) mandates that personal data must be kept up-to-date to maintain its relevance and integrity [17].

Data Security AI systems are prone to both classical security aspects related to data access and novel attack schemes such as data poisoning [80] or membership inference attacks [75]. Classical privacy concerns involve safeguarding sensitive information from unauthorized access and ensuring compliance with regulations such as GDPR [17]. Security, on the other hand, focuses on protecting data from breaches, malicious attacks, and unauthorized modifications. Encryption, secure data storage, and access controls are measures to maintain data integrity and security. From the ML side, techniques like data anonymization [56], differential privacy [4] and measures against data poisoning [80] can help protect individual privacy while maintaining data utility.

	INHERENT	SUITABILITY	SUFFICIENCY	VARIABILITY	SECURITY
Data Validation	✓				
Drift Detection				✓	
Feature Selection	✓	✓			✓
Active Learning		✓	✓		
Core-Sets	✓		✓		
Data Augmentation		✓	✓		✓
Synthetic Data	✓	✓	✓	✓	✓
Data Difficulty	✓		✓		
Data Distillation			✓		✓
Memorization	✓		✓		✓
Explainable AI	✓	✓	✓		✓
Data Attribution	✓	✓	✓	✓	✓
Sample Size	✓	✓	✓		

Table 2: Data Challenges and corresponding methods and research areas. The top methods are more established whereas the bottom ones require more research.

Relationship to other sources: Regulatory frameworks (e.g., AI Act Art. 15; GDPR Art. 32) prescribe technical and organizational measures to safeguard data integrity and confidentiality. In ML literature, this area is often termed privacy-preserving ML.

4 Methodology

The aim of this paper is to map emerging and established research efforts onto our five data challenges. To this end, we conducted a structured literature review. First, we queried Google Scholar, arXiv, and the proceedings of major AI/ML venues (e.g., NeurIPS, AAAI, DMLR) for studies on data-centric practices and regulatory requirements. Second, we extracted methods catalogued in recent surveys and overviews [48, 90, 36] and designated these as *established*. Finally, we identified methods from our search that were not covered in those summaries, labeled them as *emerging*, and mapped every method to our taxonomy of data challenges.

5 Mapping Established Methods to Regulatory Needs

We now show how common techniques address regulatory requirements. Table 2 (top) summarizes the mapping. For each method, we illustrate *opportunities* to address the identified data quality issues. We also discuss potential *challenges* associated with implementing these methods in practice.

Data Validation Techniques and Metrics Data validation aims to ensure that training datasets are accurate, consistent, and free of errors. Statistical techniques can detect and correct label mistakes [59, 60], while tools like TensorFlow Data Validation (TFDV) flag anomalies and generate descriptive statistics [10]. Beyond error detection, several works introduce quantitative metrics for data quality. For instance, Mitchell et al. [55] review measures such as Euclidean distance and Kullback–Leibler divergence, and Zhao et al. [93] analyze over 100 machine learning datasets to propose a framework for evaluating reliability, validity, and diversity. Other studies examine dimensions such as accuracy, completeness, consistency, and timeliness. We provide a comprehensive comparison of these metrics in Appendix C [31, 95, 11].

Opportunities: Data validation techniques play a crucial role in detecting irregularities within datasets, which can significantly impact model performance [59]. By identifying and correcting these errors, data validation enhances dataset accuracy and reliability, ultimately leading to improved model performance and increased trustworthiness. This addresses INHERENT data quality.

Challenges: One critical challenge in data validation is distinguishing between errors and edge cases, particularly when automatic corrections are applied. Automatically fixing errors can inadvertently alter data instances that were originally correct or overlook complex errors that require human judgment.

Drift Detection Drift Detection refers to the recognition of shifts between the data a model was trained on and real-world data. Identifying such drifts can be achieved by tracking performance changes over time or, from a data perspective, by comparing the distribution of the training data to that of the real-world data [28].

Opportunities: Drift detection directly addresses data VARIABILITY by identifying shifts in live data. This information can be utilized to take various actions, such as stopping the model, triggering a fail-safe mode, or requiring human intervention to reassess or retrain the model.

Challenges: Although drift detection is crucial, implementing it effectively can be challenging. Identifying the right metrics for detecting drift and establishing thresholds for action can be complex.

Feature Selection Feature Selection has been well addressed in the literature [78]. Jakubik et al. [36] summarize it as one of the two major components of data-centric AI. They distinguish between methods aimed at improving feature quality, which include removing irrelevant features, and methods for creating or acquiring new relevant features.

Opportunities: Feature selection can reduce the size of the dataset while also enhancing its representativity. It addresses both INHERENT data quality and SUITABILITY; removing unnecessary features increases the inherent value of the data, while certain features may only be essential for specific tasks, thereby improving the contextual relevance of the dataset. Furthermore, the choice of features may present SECURITY issues, as they might reveal sensitive information.

Challenges: Although feature selection improves efficiency and relevance, it can also introduce or amplify biases.

Active Learning Active Learning aims to optimize model performance under a fixed annotation budget by iteratively selecting the most informative examples for labeling. Given a small labeled set D and a large unlabeled pool U , the algorithm may query an oracle (e.g., a human annotator) for up to b labels. Active Learning strategies often focus on those unlabeled instances with the highest predictive uncertainty [74].

Opportunities: Active learning efficiently utilizes labeling resources by focusing on the most informative data instances, thereby addressing the challenge of SUFFICIENCY. Furthermore, active learning supports the selection of the most relevant and diverse data instances for a specific task, enhancing the representativity of the dataset for that application. This, in turn, addresses SUITABILITY, ensuring that the model is trained on data that accurately reflects the target domain.

Challenges: Implementing active learning can be computationally intensive, as it requires iterative model training and evaluation. Additionally, the effectiveness of active learning heavily depends on the choice of the query strategy, which may not be universally optimal for all types of data and tasks.

Core-Sets Core-Sets are small subset $S \subseteq D$ of a full dataset D such that the learning algorithm achieves similar performance when trained on this subset as it would on the entire dataset [20]. A popular method for finding core-sets is the use of clustering techniques, where representative samples are selected based on their distances to cluster centroids.

Opportunities: Core-sets significantly reduce the volume of data required for training without compromising model performance. This is particularly beneficial for large datasets, as it minimizes storage and computational needs while retaining the essential characteristics of the data. By ensuring that the dataset remains manageable and efficient, core-sets address the challenge of data SUFFICIENCY. Conversely, if a smaller dataset can achieve similar performance to a larger one, the smaller dataset should be favored. The data size can be considered an INHERENT property.

Challenges: Core-sets might not capture all the nuances of the original data, particularly in scenarios where rare events or minority classes are important. This limitation can lead to reduced performance in applications where such events are critical.

Data Augmentation Data Augmentation methods enhance the quality and diversity of training data by artificially increasing the size of datasets. Techniques like Synthetic Minority Over-sampling Technique (SMOTE) [12] can be effective in addressing data representativity issues, particularly class imbalances.

Opportunities: Data augmentation techniques contribute to creating a more diverse and representative dataset. By artificially increasing the data size, these methods effectively address class imbalances and enhance the model’s generalizability, leading to improved performance and fairness. While data augmentation primarily addresses data SUFFICIENCY, it can also improve SUITABILITY for specific tasks and may provide partial defenses against overfitting and model vulnerability to SECURITY breaches by diversifying the training data.

Challenges: A recent article titled *The Good, the Bad, and the Ugly Sides of Data Augmentation* summarizes the challenges well [46]. Artificially modified or balanced data may not represent real-world scenarios, potentially introducing noise or artifacts that could negatively affect model performance [16].

Synthetic Data Synthetic Data refers to the creation of artificial data for ML. This approach is particularly advantageous in scenarios where data collection is difficult, expensive, or where privacy concerns are significant [5, 50]. Typically, techniques such as Generative Adversarial Networks (GANs) or stable diffusion are employed to generate this data.

Opportunities: Synthetic data enables the generation of large, diverse datasets without collecting sensitive real-world samples, thereby reducing cost and time while safeguarding privacy and ensuring regulatory compliance. By filling gaps, e.g., underrepresented classes or rare scenarios, and simulating varied conditions, it addresses INHERENT data quality, SUITABILITY and VARIABILITY, and mitigates SECURITY risks tied to real data.

Challenges: Generating synthetic data can be challenging, especially for images, because it often requires training large models like GANs [30] or diffusion models [68]. Additionally, synthetic data can introduce another layer of bias and a content gap [85]. Obtaining performance comparable to or superior to real data remains a significant hurdle in many applications [5].

6 Mapping Emerging Methods to Regulatory Needs

Next, we map emerging directions onto our taxonomy. We classify a method as “emerging” if it was identified in our review but is not yet covered by existing surveys. Table 2 (bottom) provides an overview of these methods and the primary data challenges they address.

6.1 Data Difficulty

Meding et al. [53] describe dichotomous data difficulty and show that many datasets, such as ImageNet, suffer from imbalanced data difficulty. There are many data instances in the test set that are never classified correctly (called *impossible*) and many that are always classified correctly (called *trivial*). They show that models can be better compared on the remaining instances. A similar conclusion can be drawn for label errors. Northcutt et al. [59] show that larger models tend to be favored on datasets with label errors, while smaller models might actually outperform them when evaluated on a dataset without errors.

Opportunities: Understanding dichotomous data difficulty enhances the precision of model performance evaluation. Performance across these distinct samples could be assessed specifically, potentially through extending labels with metadata that indicates the difficulty level of individual samples. This approach could also help identify areas of the dataset that are underrepresented. If certain concepts that tend to be more difficult can be identified, it may suggest the need for gathering additional training data. Thus, this understanding can serve as both an INHERENT data quality metric and a means to address data SUFFICIENCY.

Challenges: The practicability across diverse domains requires additional empirical study.

6.2 Data Distillation

Data Distillation is a relatively new field in ML, first introduced by Wang et al. [86]. Conceptually related to core-sets, the goal of data distillation is to condense a large dataset into a smaller one that maintains similar performance. The key distinction is that distilled data consists of synthetic images, which are often not recognizable by humans. A popular method for creating these synthetic images is gradient matching [92].

Opportunities: Similar to core-sets and synthetic data, data distillation can reduce the required dataset size, thereby addressing the data SUFFICIENCY issue. Furthermore, as Yu et al. [89] point out, it can also contribute to resolving privacy, security, and robustness challenges, thus addressing SECURITY challenges.

Challenges: Although data distillation holds significant potential and has been evaluated on real-world examples, such as medical data [45], its practical applications, particularly in meeting regulatory demands, remain underexplored.

6.3 Memorization

Memorization occurs when a model heavily relies on unique training instances to make predictions. For example, in a facial recognition dataset, a single image of a person with a distinctive tattoo may be memorized, improving accuracy for that individual but risking privacy through membership inference attacks [87]. Feldman [21] defines a training instance as unique if its removal reduces the model’s ability to classify it correctly. Estimating memorization involves training models with and without the instance, though efficient methods exist [21, 38]. Jiang et al. [38] further use memorization as a measure to categorize the structure of a dataset and show that mislabeled instances are harder to memorize.

Opportunities: Memorization scores identify underrepresented regions and label errors, and, additionally, have proven useful in data pruning [77], thereby addressing data SUFFICIENCY and INHERENT data quality. They also flag privacy risks, enhancing SECURITY.

Challenges: A primary issue is the difficulty of translating memorization scores into actionable data collection strategies. Memorization scores identify unique training samples critical for generalization but do not inherently specify what additional data to collect.

6.4 Explainable AI

Several works explore how explainable artificial intelligence (XAI) can be used as a tool for improving ML models. While early works have demonstrated that XAI techniques, such as saliency maps, provide insight into the inner functionality of ML models [67, 51], several recent works address the question of how useful this is for debugging ML models [1, 2, 3, 13, 19, 47, 71, 57, 23, 40, 69], the certification of AI systems [25], and for building trust among users [62, 81, 82, 91, 52, 41, 72]. Although these works primarily focus on model debugging, practical recommendations often address data issues. A well-known example is provided by Ribeiro et al. [67], where the AI model recognizes wolves based on the presence of snow in the background rather than the animal’s features. Consequently, the training data has to be extended with more images of wolves without snow. This highlights how current research on XAI can be leveraged to address data quality challenges in ML.

Opportunities: XAI can serve as a tool for identifying biases and inconsistencies within datasets [70]. By understanding the model’s decision-making process, data scientists can pinpoint problematic areas in the dataset and take corrective actions, such as collecting more representative data or rebalancing the dataset. This approach addresses both INHERENT data quality and SUITABILITY and serves as an indicator of data SUFFICIENCY. Furthermore, it could reveal SECURITY issues by indicating whether the model has learned to replicate specific data instances and even reduce liability risks [24].

Challenges: Interpreting the explanations provided by XAI tools can be difficult, requiring expertise to ensure that the correct actions are taken based on the insights provided, often, the obtained insights are opaque [25, 24]. There is also the risk of automation bias, which may lead to overlooking other important aspects of model performance and data quality. Finally, although the application of XAI for data seems straightforward, it is underrepresented in practice [25]. This underrepresentation is further underscored by its absence in existing literature reviews, indicating that its implementation is more complex than it might appear.

6.5 Data Attribution

Data Attribution can be interpreted as data-centric XAI method. It comprises two closely related tasks: instance-level attribution and data valuation. In instance-level attribution, each test instance j is linked back to every training instance i by computing an influence score $M_{i,j}$. A high positive

score indicates that removing sample i would likely cause instance j to be misclassified, whereas a high negative score means that sample i itself contributes to j 's misclassification. The full set of scores forms an attribution matrix $M \in \mathbb{R}^{n \times m}$; since this matrix can grow very large, practitioners typically display only the top five positive and top five negative influencers. Efficient algorithms such as TRAK [63] achieve this without retraining the model for every training–test pair.

Data valuation assigns each training example i a single scalar value v_i that reflects its overall impact on model performance. The value of instance i can be interpreted as the average attribution score $v_i = \sum_{k=0}^m M_{i,k}$. Although data values can be computed from an attribution matrix, e.g., created by TRAK, it is currently more common to rely on sampling-based methods. Ghorbani and Zou [27] apply the Shapley value from cooperative game theory to data valuation. In recent years, many methods have been introduced to speed up the computation of data values, most of which aim to find better and faster approximations for the Shapley value [84, 43, 44, 37, 76].

Opportunities: In theory, data attribution and valuation can address the full spectrum of data-quality challenges. By discarding low-value examples, they enable the detection of label errors and the removal of noisy data, thereby improving INHERENT data quality. Selecting high-value samples tailored to specific domains or environmental conditions further supports domain transfer and enhances SUITABILITY and VARIABILITY. Efficient training-set reduction by eliminating low-value data also contributes to DATA SUFFICIENCY by ensuring that only the most informative samples are retained. More recently, data valuation has been adopted in data marketplaces, where prospective buyers estimate and acquire valuable datasets without direct access to the raw content [79], thus further addressing DATA SUFFICIENCY. Finally, high-value instances tend to be unique or rare, which can expose models to privacy vulnerabilities and underscore critical SECURITY considerations.

Challenges: Despite the rapid growth of XAI research, data attribution remains underexplored and is not yet widely adopted by practitioners [58]. Although data valuation is often cited as a tool for dataset understanding [90], its concrete benefits, limitations, and implications for representativeness are rarely analyzed. The complex dependencies between training and test instances make raw attribution values hard to interpret, deriving actionable security recommendations from a ranked list of influencers is nontrivial (see Fig. 1 of the Appendix).

6.6 Sample-size Estimation

Determining appropriate sample sizes for both test and training sets is crucial for reliable model evaluation and efficient data collection.

Several theoretical approaches have been proposed for test-set sizing. Guyon [33] suggest that the optimal fraction r reserved for validation (or test) should scale inversely with the square root of the number of model parameters $|\theta|$: $r = \frac{1}{\sqrt{|\theta|}}$.

An alternative formulation asks what absolute number of test samples n is needed to estimate error rates with statistical significance. Guyon et al. [32] gives the rule of thumb $n \approx \frac{100}{p}$, where p is the expected error rate of the best recognizer (e.g., a human). Additional formulations are discussed in Appendix D.

When test sets are small, high variance in performance estimates can undermine confidence. Bouthillier et al. [9] recommend running multiple evaluation trials to stabilize metrics in this case. Although the standard deviation is high with small test sets, repeated runs can compensate this when gathering more data is infeasible (see Figure 2 of the Appendix for an example).

For training-set sizing, simple heuristics (e.g. the “one-in-ten” rule of ten examples per parameter) lack solid justification. A more informative approach is to plot model performance against the fraction of training data, which typically exhibits logarithmic gains [22, 83]. Early additions yield large improvements, while returns diminish beyond a certain point (e.g., performance surpasses 90% at 50% of data, with less than a 5% gain from doubling in Appendix B).

Opportunities: Determining the appropriate test-set size ensures that performance estimates are both statistically sound and reflective of real-world conditions, thereby addressing SUITABILITY. Estimating the amount of training data helps avoid unnecessary collection and annotation costs, addressing data SUFFICIENCY. Finally, the total dataset size constitutes an INHERENT quality attribute of the data.

Challenges: In domains with scarce or costly data, gathering an ideal number of samples may simply be infeasible. Even when test or training sets are sufficiently large, models can still pick up spurious correlations. Moreover, despite foundational work on sample-size estimation predating the 2000s, systematic and widely adopted estimation methods seem underutilized in practice today.

7 An Exemplary Workflow for Regulated ML

So far, we have positioned both established and emerging data-centric methods in our taxonomy and highlighted promising research directions. We now propose an exemplary, five-step workflow as a provisional blueprint for developing a fully “regulatable” ML process:

1. Define Objectives & Scope

- (a) Enumerate use-case requirements and edge cases.
- (b) Document data obligations (e.g., privacy, fairness, robustness).
- (c) Record acceptance criteria.

2. Plan Acquisition & Splits

- (a) Derive the required test-set size (e.g. $n \approx 100/p$ for error rate p or $r = 1/\sqrt{|\theta|}$ for model complexity θ).
- (b) Collect roughly $5n$ samples and stratify into train ($3n$), test (n), and holdout (n) subsets.

3. Validate Inherent Quality & Document

- (a) Run automated checks (e.g. CleanLab for label noise).
- (b) Produce a Data Card recording provenance, quality metrics, and any corrections applied.

4. Advanced Assessment & Remediation

- (a) On the holdout set, perform drift detection, data-difficulty analysis, and attribution/valuation studies (cf. Sec. 6).
- (b) Identify coverage gaps and remediate via targeted relabeling, augmentation, synthetic data generation, or dataset distillation.

5. Final Audit & Continuous Monitoring

- (a) Merge the holdout data back into training, finalize all documentation, and archive an audit trail.
- (b) In production, deploy drift-and-bias monitors and schedule periodic revalidation to ensure ongoing compliance.

8 Summary, Limitations & Future Work

In this paper, we introduced a regulation-aligned taxonomy of five ML data challenges, systematically mapped both established and emerging data-centric methods to these challenges, and sketched a five-step workflow for “regulatable ML”. While these methods offer powerful techniques to enhance the quality, safety, and regulatory compliance of AI systems, they require further empirical validation and integration into practical toolchains and workflows.

Looking ahead, researchers can refine and validate these methods in real-world settings, closing the gap between theoretical promise and operational impact. Regulators, where technically feasible, can explore incorporating elements of this framework into emerging AI oversight regimes. Notably, data attribution and valuation, despite their unique ability to tackle all five identified challenges, remain underutilized in both industry and policy. This offers potential for follow-up work, motivating new tools, benchmarks, and standards.

Acknowledgements

This paper is funded in parts by the German Federal Ministry of Research, Technology and Space under the project “KI-Fogger”.

With funding from the:



References

- [1] Reduan Achibat, Maximilian Dreyer, Ilona Eisenbraun, Sebastian Bosse, Thomas Wiegand, Wojciech Samek, and Sebastian Lapuschkin. From attribution maps to human-understandable explanations through concept relevance propagation, 2023. URL <http://arxiv.org/pdf/2206.03208>.
- [2] Julius Adebayo, Michael Muelly, Ilaria Lliccardi, and Been Kim. Debugging tests for model explanations. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS ’20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- [3] Julius Adebayo, Michael Muelly, Hal Abelson, and Been Kim. Post hoc explanations may be ineffective for detecting unknown spurious correlation, 2022. URL <https://arxiv.org/abs/2212.04629>.
- [4] Samah Baraheem and Zhongmei Yao. A survey on differential privacy with machine learning and future outlook, 2022. URL <https://arxiv.org/abs/2211.10708>.
- [5] André Bauer, Simon Trapp, Michael Stenger, Robert Leppich, Samuel Kounev, Mark Leznik, Kyle Chard, and Ian Foster. Comprehensive exploration of synthetic data generation: A survey, 2024. URL <https://arxiv.org/abs/2401.02524>.
- [6] Firas Bayram, Bestoun S. Ahmed, and Andreas Kasser. From concept drift to model degradation: An overview on performance-aware drift detectors. *Knowledge-Based Systems*, 245: 108632, 2022. ISSN 0950-7051. doi: <https://doi.org/10.1016/j.knosys.2022.108632>. URL <https://www.sciencedirect.com/science/article/pii/S0950705122002854>.
- [7] Claudia Beleites, Ute Neugebauer, Thomas Bocklitz, Christoph Krafft, and Jürgen Popp. Sample size planning for classification models, 1 2013. ISSN 00032670.
- [8] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, S. Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen A. Creel, Jared Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren E. Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas F. Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, O. Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suir P. Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avani Narayan, Deepak Narayanan, Benjamin Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, J. F. Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher

- Potts, Aditi Raghunathan, Robert Reich, Hongyu Ren, Frieda Rong, Yusuf H. Roohani, Camilo Ruiz, Jack Ryan, Christopher R’e, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishna Parasuram Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei A. Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. On the opportunities and risks of foundation models. *ArXiv*, 2021. URL <https://crfm.stanford.edu/assets/report.pdf>.
- [9] Xavier Bouthillier, Pierre Delaunay, Mirko Bronzi, Assya Trofimov, Brennan Nichyporuk, Justin Szeto, Nazanin Mohammadi Sepahvand, Edward Raff, Kanika Madan, Vikram Voleti, et al. Accounting for variance in machine learning benchmarks, 2021.
 - [10] Eric Breck, Neoklis Polyzotis, Sudip Roy, Steven Euijong Whang, and Martin Zinkevich. Data validation for machine learning, 2019. URL <https://research.google/pubs/pub47967/>.
 - [11] Lukas Budach, Moritz Feuerpfeil, Nina Ihde, Andrea Nathansen, Nele Noack, Hendrik Patzlaff, Felix Naumann, and Hazar Harmouch. The effects of data quality on machine learning performance. 7 2022. URL <http://arxiv.org/abs/2207.14529>.
 - [12] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. Smote: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16:321–357, June 2002. ISSN 1076-9757. doi: 10.1613/jair.953. URL <http://dx.doi.org/10.1613/jair.953>.
 - [13] Julien Colin, Thomas FEL, Remi Cadene, and Thomas Serre. What i cannot predict, i do not understand: A human-centered evaluation framework for explainability methods. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, *Advances in Neural Information Processing Systems*, volume 35, pages 2832–2845. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/13113e938f2957891c0c5e8df811dd01-Paper-Conference.pdf.
 - [14] Thomas G. Dietterich. Approximate statistical tests for comparing supervised classification learning algorithms, oct 1998. ISSN 0899-7667. URL <https://doi.org/10.1162/089976698300017197>.
 - [15] Alexander Dubbs. Test set sizing via random matrix theory, 2024. URL <https://arxiv.org/abs/2112.05977>.
 - [16] Yotam Elor and Hadar Averbuch-Elor. To smote, or not to smote?, 2022. URL <https://arxiv.org/abs/2201.08528>.
 - [17] European Parliament and Council. Regulation (EU) 2016/679 of the European Parliament and of the Council, 2016. URL <https://data.europa.eu/eli/reg/2016/679/oj>.
 - [18] European Parliament and Council. Regulation (eu) 2024/1689 of the european parliament and of the council of 13 june 2024 laying down harmonised rules on artificial intelligence and amending regulations (ec) no 300/2008, (eu) no 167/2013, (eu) no 168/2013, (eu) 2018/858, (eu) 2018/1139 and (eu) 2019/2144 and directives 2014/90/eu, (eu) 2016/797 and (eu) 2020/1828 (artificial intelligence act) (text with eea relevance), 2024. ISSN 52021PC0206. URL <http://data.europa.eu/eli/reg/2024/1689/oj>.
 - [19] Thomas Fel, Agustin Picard, Louis Béthune, Thibaut Boissin, David Vigouroux, Julien Colin, Rémi Cadène, and Thomas Serre. Craft: Concept recursive activation factorization for explainability. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2711–2721, June 2023.
 - [20] Dan Feldman. Core-sets: An updated survey, 2020. ISSN 19424795. - Name core set dates back to 2005.
 - [21] Vitaly Feldman. Does learning require memorization? a short tale about a long tail, 2020. ISSN 07378017. URL <https://arxiv.org/abs/1906.05271>. - Long Tail examples
Mixture distribution.

- [22] Rosa L. Figueroa, Qing Zeng-Treitler, Sasikiran Kandula, and Long H. Ngo. Predicting sample size required for classification performance, 2012. ISSN 14726947.
- [23] Raymond Fok and Daniel S. Weld. In search of verifiability: Explanations rarely enable complementary performance in ai-advised decision making. *CoRR*, abs/2305.07722, 2023. doi: 10.48550/ARXIV.2305.07722. URL <https://doi.org/10.48550/arXiv.2305.07722>.
- [24] Benjamin Fresz, Elena Dubovitskaya, Danilo Brajovic, Marco Huber, and Christian Horz. How should ai decisions be explained? requirements for explanations from the perspective of european law, 2024.
- [25] Benjamin Fresz, Vincent Philipp Göbels, Safa Omri, Danilo Brajovic, Andreas Aichele, Janika Kutz, Jens Neuhüttler, and Marco F. Huber. The contribution of xai for the safe development and certification of ai: An expert-based analysis, 2024. URL <https://arxiv.org/abs/2408.02379>.
- [26] Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna M. Wallach, Hal Daumé III, and Kate Crawford. Datasheets for datasets, 2018. URL <http://arxiv.org/abs/1803.09010>.
- [27] Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning, 2019.
- [28] Igor Goldenberg and Geoffrey I. Webb. Survey of distance measures for quantifying concept drift and shift in numeric data. *Knowledge and Information Systems*, 60:591–615, 8 2019. ISSN 02193116. doi: 10.1007/s10115-018-1257-z.
- [29] Youdi Gong, Guangzhen Liu, Yunzhi Xue, Rui Li, and Lingzhong Meng. A survey on dataset quality in machine learning. *Information and Software Technology*, 162, 10 2023. ISSN 09505849. doi: 10.1016/j.infsof.2023.107268.
- [30] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014. URL <https://arxiv.org/abs/1406.2661>.
- [31] Nitin Gupta, Hima Patel, Shazia Afzal, Naveen Panwar, Ruhi Sharma Mittal, Shanmukha Guttula, Abhinav Jain, Lokesh Nagalapatti, Sameep Mehta, Sandeep Hans, Pranay Lohia, Aniya Aggarwal, and Diptikalyan Saha. Data quality toolkit: Automatic assessment of data quality and remediation for machine learning datasets. 8 2021. URL <http://arxiv.org/abs/2108.05935>.
- [32] I. Guyon, J. Makhoul, R. Schwartz, and V. Vapnik. What size test set gives good error rate estimates? *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(1):52–64, 1998. doi: 10.1109/34.655649.
- [33] Isabelle M Guyon. A scaling law for the validation-set training-set size ratio. 1997. URL <https://api.semanticscholar.org/CorpusID:16194090>.
- [34] Zayd Hammoudeh and Daniel Lowd. Training data influence analysis and estimation: a survey. *Machine Learning*, 113(5):2351–2403, March 2024. ISSN 1573-0565. doi: 10.1007/s10994-023-06495-7. URL <http://dx.doi.org/10.1007/s10994-023-06495-7>.
- [35] ISO/IEC 5259. Artificial intelligence — data quality for analytics and machine learning (ml). Standard, International Organization for Standardization, Geneva, CH, 2023. URL <https://www.iso.org/standard/81088.html>.
- [36] Johannes Jakubik, Michael Vössing, Niklas Kühl, Jannis Walk, and Gerhard Satzger. Data-centric artificial intelligence, 2024. URL <https://arxiv.org/abs/2212.11854>.
- [37] Kevin Fu Jiang, Weixin Liang, James Zou, and Yongchan Kwon. Opendataval: a unified benchmark for data valuation, 2023. URL <https://arxiv.org/abs/2306.10577>.
- [38] Ziheng Jiang, Chiyuan Zhang, Kunal Talwar, and Michael C. Mozer. Characterizing structural regularities of labeled data in overparameterized models, 2020. ISSN 2331-8422. URL <http://arxiv.org/abs/2002.03206>.

- [39] V. Roshan Joseph. Optimal ratio for data splitting, 8 2022. ISSN 19321872.
- [40] Serhiy Kandul, Vincent Micheli, Juliane Beck, Markus Kneer, Thomas Burri, Francois Fleuret, and Markus Christen. Explainable ai: A review of the empirical literature. *SSRN Electronic Journal*, 01 2023. doi: 10.2139/ssrn.4325219.
- [41] Sunnie S. Y. Kim, Nicole Meister, Vikram V. Ramaswamy, Ruth Fong, and Olga Russakovsky. HIVE: Evaluating the human interpretability of visual explanations. In *European Conference on Computer Vision (ECCV)*, 2022.
- [42] Frank Konietzschke, Karima Schwab, and Markus Pauly. Small sample sizes: A big data problem in high-dimensional data analysis, 3 2021. ISSN 14770334.
- [43] Yongchan Kwon and James Zou. Beta shapley: a unified and noise-reduced data valuation framework for machine learning, 2022. URL <https://arxiv.org/abs/2110.14049>.
- [44] Yongchan Kwon and James Zou. Data-oob: Out-of-bag estimate as a simple and efficient data value, 2023. URL <https://arxiv.org/abs/2304.07718>.
- [45] Guang Li, Ren Togo, Takahiro Ogawa, and Miki Haseyama. Soft-label anonymous gastric x-ray image distillation. In *2020 IEEE International Conference on Image Processing (ICIP)*, page 305–309. IEEE, October 2020. doi: 10.1109/icip40778.2020.9191357. URL <http://dx.doi.org/10.1109/ICIP40778.2020.9191357>.
- [46] Chi-Heng Lin, Chiraag Kaushik, Eva L. Dyer, and Vidya Muthukumar. The good, the bad and the ugly sides of data augmentation: An implicit spectral regularization perspective, 2024. URL <https://arxiv.org/abs/2210.05021>.
- [47] Yi-Shan Lin, Wen-Chuan Lee, and Z. Berkay Celik. What do you see? evaluation of explainable artificial intelligence (xai) interpretability through neural backdoors. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, KDD ’21*, page 1027–1035, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450383325. doi: 10.1145/3447548.3467213. URL <https://doi.org/10.1145/3447548.3467213>.
- [48] Michael A. Lones. How to avoid machine learning pitfalls: a guide for academic researchers, 2021. URL <https://arxiv.org/abs/2108.02497>.
- [49] Shayne Longpre, Robert Mahari, Naana Obeng-Marnu, William Brannon, Tobin South, Katy Gero, Sandy Pentland, and Jad Kabbara. Data authenticity, consent, and provenance for ai are all broken: what will it take to fix them?, 2024. URL <https://arxiv.org/abs/2404.12691>.
- [50] Yingzhou Lu, Minjie Shen, Huazheng Wang, Xiao Wang, Capucine van Rechem, Tianfan Fu, and Wenqi Wei. Machine learning for synthetic data generation: A review, 2024. URL <https://arxiv.org/abs/2302.04062>.
- [51] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions, 2017.
- [52] Shuai Ma, Ying Lei, Xinru Wang, Chengbo Zheng, Chuhan Shi, Ming Yin, and Xiaojuan Ma. Who should i trust: Ai or myself? leveraging human and ai correctness likelihood to promote appropriate trust in ai-assisted decision-making. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems, CHI ’23*, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9781450394215. doi: 10.1145/3544548.3581058. URL <https://doi.org/10.1145/3544548.3581058>.
- [53] Kristof Meding, Luca M. Schulze Buschoff, Robert Geirhos, and Felix A. Wichmann. Trivial or impossible – dichotomous data difficulty masks model differences (on imagenet and beyond), 2021. URL <http://arxiv.org/abs/2110.05922>.
- [54] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning, 2022. URL <https://arxiv.org/abs/1908.09635>.
- [55] Margaret Mitchell, Alexandra Sasha Luccioni, Nathan Lambert, Marissa Gerchick, Angelina McMillan-Major, Ezinwanne Ozoani, Nazneen Rajani, Tristan Thrush, Yacine Jernite, and Douwe Kiela. Measuring data, 12 2022. URL <http://arxiv.org/abs/2212.05129>.

- [56] Stéphane Monteiro, Diogo Oliveira, João António, Filipe Sá, Cristina Wanzeller, Pedro Martins, and Maryam Abbasi. Data anonymization: Techniques and models. In José Luís Reis, Marisa Del Rio Araujo, Luís Paulo Reis, and José Paulo Marques dos Santos, editors, *Marketing and Smart Technologies*, pages 73–84, Singapore, 2024. Springer Nature Singapore. ISBN 978-981-99-0333-7.
- [57] Romy Müller. How explainable ai affects human performance: A systematic review of the behavioural consequences of saliency maps. *International Journal of Human–Computer Interaction*, 0(0):1–32, 2024. doi: 10.1080/10447318.2024.2381929. URL <https://doi.org/10.1080/10447318.2024.2381929>.
- [58] Elisa Nguyen, Evgenii Kortukov, Jean Y. Song, and Seong Joon Oh. Exploring practitioner perspectives on training data attribution explanations, 2023. URL <https://arxiv.org/abs/2310.20477>.
- [59] Curtis G. Northcutt, Anish Athalye, and Jonas Mueller. Pervasive label errors in test sets destabilize machine learning benchmarks, 3 2021. URL <http://arxiv.org/abs/2103.14749>.
- [60] Curtis G. Northcutt, Lu Jiang, and Isaac L. Chuang. Confident learning: Estimating uncertainty in dataset labels, 2022. URL <https://arxiv.org/abs/1911.00068>.
- [61] Will Orr and Kate Crawford. Building better datasets: Seven recommendations for responsible design from dataset creators. *Journal of Data-centric Machine Learning Research*, 2024. URL <https://openreview.net/forum?id=6bd8BrRKTW>.
- [62] Andrea Papenmeier, Gwenn Englebienne, and Christin Seifert. How model accuracy and explanation fidelity influence user trust, 2019. URL <https://arxiv.org/abs/1907.12652>.
- [63] Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and Aleksander Madry. Trak: Attributing model behavior at scale, 2023. URL <https://arxiv.org/abs/2303.14186>.
- [64] Maria Priestley, Fionntán O'Donnell, and Elena Simperl. A survey of data quality requirements that matter in ml development pipelines. *Journal of Data and Information Quality*, 15, 6 2023. ISSN 19361963. doi: 10.1145/3592616.
- [65] Mahima Pushkarna, Andrew Zaldivar, and Oddur Kjartansson. Data cards: Purposeful and transparent dataset documentation for responsible ai, 4 2022. URL <http://arxiv.org/abs/2204.01075>.
- [66] S. J. Raudys and A. K. Jain. Small sample size effects in statistical pattern recognition: Recommendations for practitioners and open problems. volume 1, pages 417–423. Publ by IEEE, 1990. ISBN 0818620625. doi: 10.1109/icpr.1990.118138.
- [67] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. " why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144, 2016.
- [68] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022. URL <https://arxiv.org/abs/2112.10752>.
- [69] Yao Rong, Tobias Leemann, Thai-Trang Nguyen, Lisa Fiedler, Peizhu Qian, Vaibhav Unhelkar, Tina Seidel, Gjergji Kasneci, and Enkelejda Kasneci. Towards human-centered explainable ai: A survey of user studies for model explanations. *IEEE Trans. Pattern Anal. Mach. Intell.*, 46(4):2104–2122, nov 2023. ISSN 0162-8828. doi: 10.1109/TPAMI.2023.3331846. URL <https://doi.org/10.1109/TPAMI.2023.3331846>.
- [70] Nina Schaaf, Omar de Mitri, Hang Beom Kim, Alexander Windberger, and Marco F. Huber. Towards measuring bias in image classification, 2021. URL <https://arxiv.org/abs/2107.00360>.

- [71] Max Schemmer, Patrick Hemmer, Maximilian Nitsche, Niklas Kühl, and Michael Vössing. A meta-analysis of the utility of explainable artificial intelligence in human-ai decision-making. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '22, page 617–626, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450392471. doi: 10.1145/3514094.3534128. URL <https://doi.org/10.1145/3514094.3534128>.
- [72] Max Schemmer, Niklas Kuehl, Carina Benz, Andrea Bartos, and Gerhard Satzger. Appropriate reliance on ai advice: Conceptualization and the effect of explanations. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*, IUI '23, page 410–422, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701061. doi: 10.1145/3581641.3584066. URL <https://doi.org/10.1145/3581641.3584066>.
- [73] Mohamed El Amine Seddik, Suei-Wen Chen, Soufiane Hayou, Pierre Youssef, and Merouane Debbah. How bad is training on synthetic data? a statistical analysis of language model collapse, 2024. URL <https://arxiv.org/abs/2404.05090>.
- [74] Burr Settles. Active learning literature survey, 2010. ISSN 00483931.
- [75] Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models, 2017. URL <https://arxiv.org/abs/1610.05820>.
- [76] Rachael Hwee Ling Sim, Xinyi Xu, and Kian Hsiang Low. Data valuation in machine learning: "ingredients", strategies, and open challenges. In *International Joint Conference on Artificial Intelligence*, 2022. URL <https://api.semanticscholar.org/CorpusID:249319573>.
- [77] Ben Sorscher, Robert Geirhos, Shashank Shekhar, Surya Ganguli, and Ari S. Morcos. Beyond neural scaling laws: beating power law scaling via data pruning, 6 2022. URL <http://arxiv.org/abs/2206.14486>.
- [78] Dipti Theng and Kishor K. Bhoyar. Feature selection techniques for machine learning: a survey of more than two decades of research, 3 2024. ISSN 02193116.
- [79] Zhihua Tian, Jian Liu, Jingyu Li, Xinle Cao, Ruoxi Jia, Jun Kong, Mengdi Liu, and Kui Ren. Private data valuation and fair payment in data marketplaces, 2023. URL <https://arxiv.org/abs/2210.08723>.
- [80] Zhiyi Tian, Lei Cui, Jie Liang, and Shui Yu. A comprehensive survey on poisoning attacks and countermeasures in machine learning. *ACM Comput. Surv.*, 55(8), dec 2022. ISSN 0360-0300. doi: 10.1145/3551636.
- [81] Amy Turner, Meenakshi Kaushik, Mu-Ti Huang, and Srikar Varanasi. Calibrating trust in ai-assisted decision making, 2020. URL <https://www.ischool.berkeley.edu/projects/2020/calibrating-trust-ai-assisted-decision-making>.
- [82] Mor Vered, Tali Livni, Piers Douglas Lionel Howe, Tim Miller, and Liz Sonenberg. The effects of explanations on automation bias. *Artificial Intelligence*, 322:103952, 2023. ISSN 0004-3702. doi: <https://doi.org/10.1016/j.artint.2023.103952>. URL <https://www.sciencedirect.com/science/article/pii/S000437022300098X>.
- [83] Tom Viering and Marco Loog. The shape of learning curves: a review, 3 2021. URL <http://arxiv.org/abs/2103.10948>.
- [84] Jiachen T. Wang and Ruoxi Jia. Data banzhaf: A robust data valuation framework for machine learning, 2023. URL <https://arxiv.org/abs/2205.15466>.
- [85] Ruyu Wang, Sabrina Schmedding, and Marco F. Huber. Improving the effectiveness of deep generative data, 2023. URL <https://arxiv.org/abs/2311.03959>.
- [86] Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A. Efros. Dataset distillation, 2020. URL <https://arxiv.org/abs/1811.10959>.
- [87] Jiaheng Wei, Yanjun Zhang, Leo Yu Zhang, Ming Ding, Chao Chen, Kok-Leong Ong, Jun Zhang, and Yang Xiang. Memorization in deep learning: A survey, 2024. URL <https://arxiv.org/abs/2406.03880>.

- [88] Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zac Kenton, Sasha Brown, Will Hawkins, Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks, William Isaac, Sean Legassick, Geoffrey Irving, and Iason Gabriel. Ethical and social risks of harm from language models, 2021. URL <https://arxiv.org/abs/2112.04359>.
- [89] Ruonan Yu, Songhua Liu, and Xinchao Wang. Dataset distillation: A comprehensive review, 2023. URL <https://arxiv.org/abs/2301.07014>.
- [90] Daochen Zha, Zaid Pervaiz Bhat, Kwei-Herng Lai, Fan Yang, Zhimeng Jiang, Shaochen Zhong, and Xia Hu. Data-centric artificial intelligence: A survey, 2023. URL <https://arxiv.org/abs/2303.10158>.
- [91] Yunfeng Zhang, Q. Vera Liao, and Rachel K. E. Bellamy. Effect of confidence and explanation on accuracy and trust calibration in ai-assisted decision making. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, FAT* '20*, page 295–305, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450369367. doi: 10.1145/3351095.3372852. URL <https://doi.org/10.1145/3351095.3372852>.
- [92] Bo Zhao, Konda Reddy Mopuri, and Hakan Bilen. Dataset condensation with gradient matching, 2021. URL <https://arxiv.org/abs/2006.05929>.
- [93] Dora Zhao, Jerone T A Andrews, Orestis Papakyriakopoulos, and Alice Xiang. Position: Measure dataset diversity, don’t just claim it, 2024.
- [94] Dora Zhao, Jerone T. A. Andrews, Orestis Papakyriakopoulos, and Alice Xiang. Position: Measure dataset diversity, don’t just claim it, 2024. URL <https://arxiv.org/abs/2407.08188>.
- [95] Yuhao Zhou, Fengjiao Tu, Kewei Sha, Junhua Ding, and Haihua Chen. A survey on data quality dimensions and tools for machine learning. 6 2024. URL <http://arxiv.org/abs/2406.19614>.
- [96] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning, 2020. URL <https://arxiv.org/abs/1911.02685>.

Appendix

A Data Attribution Example



Figure 1: Example of data attributions computed with TRAK. The leftmost image is the test sample (a cat). Above it are the five training images with the highest positive influence: adding them to the training set increases the model’s confidence in correctly classifying the test image. Below are the five images with the strongest negative influence: their presence tends to reduce classification confidence. Some attributions can be counterintuitive (e.g. a dog image showing positive influence on the cat sample), highlighting challenges in interpreting influence scores.

B Removing Train and Test Data

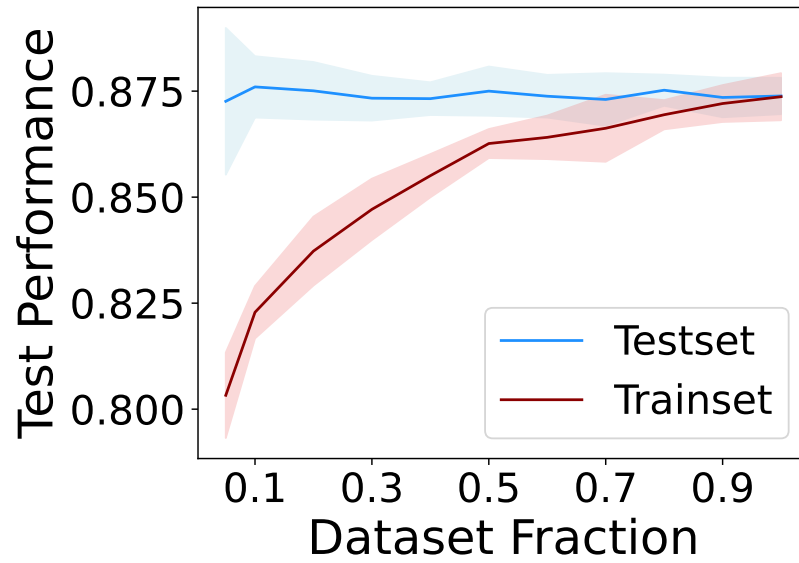


Figure 2: **Blue:** Test accuracy and deviation (y -axis) when training 10 models on full training data with increasing test size (x -axis). **Red:** Test performance and deviation (y -axis) with increasing size of training data (x -axis) on full test set.

C Data Quality Metrics from Research

Table 3: Examples of different data measurements from [55].

		DISTANCE		DENSITY	DIVERSITY		TENDENCY		ASSOCIATION
Physical ences	Sci-	Length		Mass-per- volume	Biodiversity		Mean, Mode	Median,	Correlation
General Data Measures		Euclidean Dis- tance		Data Density	Gini Diversity		Burstiness		
		Cosine Similar- ity		KNN Density	Vendi Score				
		Earth Mover’s Distance							
		Kullback- Leibler Diver- gence							
Modality- Specific Data Measures		Word Mover’s Distance (lan- guage)		Information Density (lan- guage)	Text Diversity (language)		Perplexity (language)		Pointwise Mutual Information
		Levenshtein Distance (lan- guage)		Idea Density (language)	Lexical Diver- sity (language)		Fit to Zipf’s Law (language)		
		Inception Dis- tance (vision)		The Incep- tion score (vision)	Image Diversity (vision)				
					Subset Diver- sity (vision)				

Table 4: Data Quality dimensions, metrics, descriptions, and examples from [95].

DIMENSION	METRICS	DESCRIPTION	EXAMPLES
Intrinsic	Correctness	A record in a dataset is free of errors.	Before starting a mailing campaign, the correctness of the attributes "postal code" shall be evaluated, and even small deviations shall be penalized because a deviation of only 1% (the postal codes 80000 and 79200) hinders the delivery of a mailing.
		Data is correctly labeled if it is a labeled record.	In the medical domain, an informal phrase of "lack of feeling" should be labeled as "numbness".
	Duplication	Measures if the same instances repeat in the dataset, especially in both the training and test datasets.	If a record in a medical concept training dataset is "Hunger – don't want to eat", and there is exactly the same record in a test dataset, then the record is considered as an overlapped record in the two datasets.
	Trustworthiness	Defines how factual the source that provides the information is. It can be subjectively evaluated, such as indicating the level on a scale, or the data can go through fact-check algorithms.	For a medical concept dataset, it should be obtained directly from the hospital's system, which undergoes regular data quality checks and is maintained according to industry standards.
Contextual	Class imbalance	Evaluates if the distribution of examples across the known classes is biased or skewed.	Most of the contemporary works on class imbalance fall into the imbalance ratios ranging from 1:4 up to 1:100. The imbalance ratio may range from 1:1000 up to 1:5000 for extreme class imbalance problems.
	Completeness	A complete dataset should include as few missing values as possible.	A medical insurance dataset must include a customer's birthdate, otherwise the medical consumption forecast model performance will be hindered.
	Comprehensiveness	A dataset contains all representative samples from the population.	In a medical text classification task, the training dataset should contain sufficient labeled medical texts covering all the conditions, symptoms, and treatments.
	Unbiasedness	Refers to whether the data used for ML training has a distribution bias or historical bias.	Photo recognition software does not recognize the facial expressions of ethnic minorities, or electronic soap dispensers that do not respond to darker skin tones because the training image datasets have an insufficient representation of some geographic regions.
	Variety	Requires each validation dataset and the test dataset to contain a significant amount of new data compared to the corresponding training dataset.	The percentage of the overlapped data between a test/validation dataset and its corresponding training dataset should be as low as possible, such as less than 10%.
Representational	Conformity	Measures how much the data conforms to the conventions for capturing information in a certain manner, including machine-readable data structures and formats for capturing specific attributes.	In a text classification task, a dataset of textual documents is labeled with sentiment (positive, negative, neutral). The labels should be encoded following a standardized set of categories, and all data processing, such as removing punctuation, converting text to lowercase, and tokenizing sentences, should be made to the whole dataset.
	Consistency	Requires data to be presented in the same format and to be compatible with previous data.	In an image classification task, if one dataset uses pixel values in the range [0, 255], while another dataset scales pixel values to the range [0, 1], this will cause inconsistency in model training and predictions.
Accessibility	Sufficiency	High data sufficiency ensures that data is readily accessible with defined user permissions for access and modifications.	In a healthcare ML application for diagnosing diseases from medical images, user entitlements are managed through strict access controls, allowing only authorized medical professionals and data scientists to access the images.

Table 5: Mapping of EU AI Act (Art. 10) dataset quality requirements to ISO/IEC 5259 [35] definitions.

METRIC	DESCRIPTION	CATEGORY
Accuracy	The degree to which data correctly describes the “real world” object or event.	Inherent
Completeness	The extent to which data has no missing values.	Inherent
Relevance	The extent to which data is applicable and useful for the intended purpose.	Inherent
Representativeness	The degree to which data accurately reflects the broader population or phenomenon.	Inherent

Table 6: Recommendations for Reliability and Validity from Zhao et al. [93].

TOPIC	RECOMMENDATION	DESCRIPTION
Reliability	Inter-annotator agreement	An established method for assessing reliability, particularly in crowdsourcing, is through inter-annotator agreement. This method often entails multiple annotators labeling an instance, with the final label determined by a majority vote. Another method to gauge inter-annotator reliability is by employing statistical measures of agreement. We find that some text datasets provide quantitative metrics to quantify inter-annotator agreement, such as Fleiss’s κ or Cohen’s κ . While consensus methods are employed in both text and image datasets, quantitative metrics for inter-annotator agreement are reported exclusively in text datasets. We recommend that image dataset curators also incorporate these statistical measures when evaluating crowdsourced labels.
	Test-retest reliability	Another approach that dataset collectors can adopt is the test-retest method. In education, this method involves administering the same test twice over a period, with consistent results indicating reliability. This principle is particularly relevant when assessing the reliability of collection methods like web scraping. For instance, curators can reapply the same methodology to recollect instances, validating whether the recollected dataset maintains the same diversity properties. Nonetheless, a lack of reliability from these tests does not necessarily imply that the collection methodology inadequately captures diversity. Changes in the underlying data distribution over time can influence the results. For example, when evaluating linguistic diversity using data scraped from Reddit, major societal events, such as elections, can unexpectedly alter the distribution. Even in such cases, measuring test-retest reliability remains valuable for gaining insights into potential shifts in data distributions.
Validity	Convergent validity: Cross-dataset generalization	Commonly employed to evaluate "dataset bias," cross-dataset generalization enables researchers to compare datasets. By utilizing existing datasets with similar structures (e.g., label taxonomy, modality) and constructs of diversity, collectors can train on their dataset and test on existing datasets or vice versa, comparing relevant metrics such as accuracy. Model performance can also be assessed against standard train-test splits from the same dataset. If the models perform similarly in both cross-dataset and same-dataset scenarios, it suggests that the datasets have similar distributions for the target variable, indicating correlated constructs of diversity. Model performance can also be assessed against standard train-test splits from the same dataset. If the models perform similarly in both cross-dataset and same-dataset scenarios, it suggests that the datasets have similar distributions for the target variable, indicating correlated constructs of diversity. However, a constraint of employing cross-dataset generalization is the necessity for congruent taxonomies (for the target variable) and comparable distributions across datasets.
	Convergent validity: Comparing existing diversity metrics	Dataset collectors can leverage established metrics for measuring data diversity. For instance, the Vendi Score, drawing inspiration from ecology and quantum statistical mechanics, has been introduced as a measure of diversity within image and text dataset categories. Curators can demonstrate how their collection process aligns with such recognized diversity metrics. Given that diversity metrics depend on the embedding space employed, datasets should be benchmarked across a multiplicity of spaces optimized for the definition of diversity selected by the dataset curators.
	Discriminant validity	Discriminant validity assesses whether measurements for theoretically unrelated constructs yield unrelated results. Consider the initial Visual Question Answer (VQA) dataset, which aimed to collect diverse and interesting questions and answers, encompassing question types such as "What is ...", "How many ...", and "Do you see a ...". If diversity is defined by the types of questions asked, it should have no relation to other factors, such as gender distribution. Prior works identified language biases in how questions and answers are formulated in the VQA dataset. For instance, based on the dataset construction, a model predicting "Yes" whenever the question begins with "Do you see a ..." can achieve high accuracy without considering the image in question. This suggests potential low discriminant validity for the given measure, highlighting the importance of applying discriminant validity to mitigate construction biases during dataset creation.

D Suggested Split Ratios for Train-Test Splits

Regarding split ratios, several works attempt to find theoretical justifications for optimal splits. Recently, Dubbs [15] derived ideal splits for linear regression models and independent Gaussian distributions. Guyon [33] suggest that *the fraction of patterns reserved for the validation set should be inversely proportional to the square root of the number of free adjustable parameters*. Let $|\theta|$ be the number of adjustable parameters in the model, then the optimal test ratio r is given by:

$$r = \frac{1}{\sqrt{|\theta|}} .$$

A similar formula is suggested by Joseph [39]. Instead of the free parameters of the model, they suggest using the number of parameters in a linear regression model that explains the data well. Let $|\beta|$ be the number of parameters in such a linear regression model, then the optimal ratio r is given by:

$$r = \frac{1}{\sqrt{|\beta| + 1}} .$$

Another approach is to reverse the question and ask *what size test set gives good error rate estimates* [32]. Guyon et al. [32] propose various formulas for this. The simplest form suggests that the number of test samples n should be approximately:

$$n \approx \frac{100}{p}$$

where p is the *expected error rate of the best recognizer*, i.e., a human. This approach is related to sample size estimation based on statistical significance, commonly used in medical and psychological research. For example, to evaluate the effect of medication against a certain disease and the target is an error margin of ϵ with a confidence level of $\sigma = 0.99$, the necessary sample size n is given by:

$$n \approx \frac{z_{0.99}^2 \times p(1 - p)}{\epsilon^2} .$$

Here, p is the a-priori known occurrence rate of the disease, and $z_{0.99}$ is the value for the selected confidence interval derived from the normal distribution. This method can also be applied to ML and has been discussed in several works from the medical domain [7, 14, 66, 42]. In the context of ML, the interpretation could be as follows: if a model is trained to predict the disease with a randomly drawn test set of size n , there is a 99% confidence that the real-world model performance is within 1% of the error rate on the test set (for $\epsilon = 0.01$ and $\sigma = 0.99$).