

Title:**Towards Improved Datasets in Machine Learning: Hypotheses on Pollution, Poison, and the Role of Misspellings****Jubo Zhang**

digbigAI, Independent Research Lab, Newark, CA

jerryzhang668@gmail.com

Abstract:

The quality and relevance of training data are critical determinants of the performance of machine learning models. This paper proposes three hypotheses concerning the composition of datasets: (1) **Pollution**: The introduction of heterogeneous data sources—such as multiple languages or mixed-domain content—can impair model performance; (2) **Poison**: The presence of spurious correlations, false factors, and low-quality data within datasets may lead to degraded performance or erroneous outputs; and (3) **Misspelling Inclusion**: Intentional incorporation of misspelled inputs can improve a model’s robustness to real-world noisy data. We further propose the integration of automated tools and specialized AI modules to detect, manage, and remediate these issues. Our discussion synthesizes existing literature with novel hypotheses, highlighting strategies for ensuring robust model training and deployment.

1. Introduction

As machine learning (ML) systems grow in scale and reach, the integrity of the data used to train them becomes a critical concern. Larger datasets are more likely to include conflicting, mislabeled, or noisy examples [1]. In this paper, we introduce three key hypotheses regarding dataset composition:

- **Pollution**: mixing multilingual or cross-domain data sources reduces performance on focused tasks [2].
- **Poisoning**: unintentional inclusion of spurious correlations or low-quality labels can undermine model learning [3,4].
- **Misspelling**: deliberately including common user-generated errors may improve model robustness [5,6].

We propose concrete methods for identifying and mitigating harmful data while leveraging noise that may enhance generalization.

2. Hypotheses and Framework

2.1 Hypothesis One: Dataset Pollution

Hypothesis: Combining unrelated domains or languages in a single dataset pollutes training, leading to performance degradation in specific target tasks.

Discussion:

Multilingual or mixed-domain datasets can introduce ambiguous patterns and inconsistent semantics, making it harder for models to build coherent representations. Research in cross-lingual learning and domain-specific modeling shows that such pollution often harms single-task performance [2,7].

2.2 Hypothesis Two: Dataset Poisoning

Hypothesis: Even small amounts of mislabeling, spurious correlations, or anomalous entries in a dataset can poison training outcomes.

Discussion:

Data poisoning is often unintentional. Common causes include annotation errors, corrupted files, or selection bias. Studies on robust learning show that noisy labels and incorrect signals can disproportionately degrade accuracy, especially in deep models [3,4,8].

2.3 Hypothesis Three: Controlled Misspelling Inclusion

Hypothesis: Including systematic misspellings in the training data improves a model's ability to handle noisy inputs at inference time.

Discussion:

Unlike other forms of noise, misspellings are natural and frequent in real-world input. Several studies show that when models are exposed to this noise type during training, their tolerance to unseen typographical errors increases [5,6]. However, this inclusion must be controlled — random noise may still cause harm unless carefully modeled [9].

3. Strategies for Detection and Remediation

3.1 Pollution Detection Tools

- **Language and Domain Classifiers:** Use pretrained NLP tools to separate inputs by language and topical domain [10].
 - **Fine-Tuning on Clean Subsets:** Adapt general-purpose models to cleaner, task-specific data through domain-specific fine-tuning.
 - **Stream Monitoring:** Track shifts in language and domain content during dataset construction to avoid gradual drift.
-

3.2 Poison Remediation Tools

- **Outlier Detection:** Use clustering and anomaly detection to flag potentially harmful examples.
 - **Cross-Source Validation:** Use overlapping datasets to catch contradictions and inconsistencies [3].
 - **Human-in-the-Loop Validation:** Leverage active learning to prioritize review of uncertain or novel data [11].
-

3.3 Misspelling-Driven Robustness Techniques

- **Synthetic Misspelling Augmentation:** Automatically generate typos based on real-world patterns and introduce them during training [5,9].
 - **Dual-Set Training:** Train jointly on clean and noisy data to improve generalization without harming accuracy.
 - **Pre-Inference Correction:** Use pretrained spelling correction models to normalize user input during deployment [6].
-

4. Experimental Proposals

To test our hypotheses, we recommend four experiment types:

- **Pollution Comparison:** Evaluate model performance when trained on clean vs. polluted (multilingual or cross-domain) data [2].
 - **Poison Injection Studies:** Intentionally inject errors to measure model resilience and the effectiveness of cleaning tools [3,4].
 - **Misspelling Stress Tests:** Use natural and synthetic noisy inputs to benchmark robustness improvements [5].
 - **Tool Impact Evaluation:** Compare results from models trained with vs. without detection/remediation tools.
-

5. Future Work

Future efforts should focus on:

- Addressing other real-world noise types such as grammar mistakes and code-switching.
 - Developing dataset quality scoring metrics that correlate with downstream performance.
 - Integrating curation tools into end-to-end training pipelines.
 - Evaluating the long-term benefits of these interventions in production systems.
-

6. Conclusion

We propose three hypotheses on the impact of data composition in machine learning. While pollution and poison degrade model reliability, controlled inclusion of realistic misspellings may improve robustness. Together with automated tools and curation strategies, these insights could help improve the resilience of ML systems across noisy and diverse environments.

References

[1] E. Bhardwaj et al. Machine Learning Data Practices through a Data Curation Lens: An Evaluation Framework. *arXiv:2405.02703v1*, 2024.

- [2] A. Conneau et al. Unsupervised Cross-lingual Representation Learning at Scale. In *ACL*, 2020.
- [3] C. Northcutt, A. Athalye, and J. Mueller. Pervasive Label Errors in Test Sets Destabilize Machine Learning Benchmarks. In *NeurIPS*, 2021.
- [4] B. Biggio and F. Roli. Wild Patterns: Ten Years After the Rise of Adversarial Machine Learning. *Pattern Recognition*, 2018.
- [5] Y. Belinkov and Y. Bisk. Synthetic and Natural Noise Both Break Neural Machine Translation. In *ICLR*, 2018.
- [6] D. Pruthi, B. Dhingra, and Z. C. Lipton. Combating Adversarial Misspellings with Robust Word Recognition. In *ACL*, 2019.
- [7] O. Xhelili, Y. Liu, H. Schütze. Breaking the Script Barrier in Multilingual Pre-Trained Language Models with Transliteration-Based Post-Training Alignment. In *EMNLP*, 2024.
- [8] A. Ghosh, H. Kumar, P.S. Sastry. Robust Loss Functions under Label Noise for Deep Neural Networks. In *AAAI*, 2017.
- [9] S. Eger et al. Text Processing Like Humans Do: Visually Attacking and Shielding NLP Systems. In *NAACL*, 2020.
- [10] G. Lample and A. Conneau. Cross-lingual Language Model Pretraining. In *NeurIPS*, 2019.
- [11] B. Settles. Active Learning Literature Survey. *University of Wisconsin–Madison*, 2009.