

---

# The Pitfalls of Memorization: When Memorization Hinders Generalization

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Neural networks often learn simple explanations that fit the majority of the data  
2 while memorizing exceptions that deviate from these explanations. This leads to  
3 poor generalization when the learned explanations are spurious. In this work, we  
4 formalize *the interplay between memorization and generalization*, showing that  
5 spurious correlations, when combined with memorization, can reduce the training  
6 loss to zero, leaving no incentive to learn robust, generalizable patterns. To address  
7 this issue, we introduce *memorization-aware training* (MAT). MAT leverages the  
8 flip side of memorization by using held-out predictions to adjust a model’s logits,  
9 guiding it towards learning robust patterns that remain invariant from training to  
10 test, thereby enhancing generalization under distribution shifts.

## 11 1 Introduction

12 Neural networks can learn simple explanations that work for the majority of their training data  
13 (Geirhos et al., 2020; Shah et al., 2020; Dherin et al., 2022). These models might then treat minority  
14 examples—those that do not conform to the learned explanation—as exceptions (Zhang et al., 2021).  
15 This becomes particularly problematic if the learned explanation is spurious, meaning it does not  
16 hold in general or is not representative of the true data distribution (Idrissi et al., 2022; Sagawa et al.,  
17 2020; Pezeshki et al., 2021; Puli et al., 2023).

18 Empirical Risk Minimization (ERM), the standard learning algorithm for neural networks, can  
19 exacerbate this issue. ERM enables neural networks to quickly capture spurious correlations and,  
20 with sufficient capacity, memorize the remaining examples rather than learning the true patterns that  
21 explain the entire dataset. This could be dangerously misleading, as a model that appears to excel in  
22 most cases may have actually captured a spurious correlation. Combined with memorization of the  
23 remaining minority examples, a neural network can **fully mask its failure** to grasp the true patterns  
24 in the data, giving a false sense of reliability and robustness.

25 Identifying whether a model with nearly perfect accuracy on the training data has learned generalizable  
26 patterns or merely relies on a mix of spurious correlations and memorization is critical. The answer  
27 lies in the model’s performance on held-out data, particularly on minority examples. Metrics such as  
28 held-out average accuracy or more fine-grained group accuracies can help us identify a better model.  
29 A question that arises is: *How can one use held-out performance signals to proactively guide a model  
30 toward learning generalizable patterns?*

31 Traditionally, held-out performance signals are primarily used for hyperparameter tuning and model  
32 selection. However, in this work, we propose a novel approach that leverages these signals more strate-  
33 gically to guide the learning process. But we first need to precisely understand when memorization  
34 can hinder generalization. Towards this goal, our paper makes the following contributions:

- 35 • *Formalizing the interplay between memorization and spurious correlations:* We study how  
36 memorization affects generalization in an interpretable setup, revealing that spurious correlations  
37 *alone* do not cause poor generalization in neural networks. Instead, it is the combination of  
38 spurious correlations with memorization that leads to this problem. Our analysis shows that models  
39 trained with ERM tend to rely on spurious features for the majority of the data while memorizing  
40 exceptions, achieving zero training loss but failing to generalize on minority examples.
- 41 • *Introducing memorization-aware training (MAT):* MAT is a novel learning algorithm that leverages  
42 the flip side of memorization by using held-out predictions to adjust a model’s logits during  
43 training. This adjustment guides the model toward learning invariant features that generalize  
44 better under distribution shifts. Unlike ERM, which relies on the i.i.d. assumption, MAT is built  
45 upon an alternative assumption that takes into account the instability of spurious correlations  
46 across different data distributions.

47 The main body of the paper examines our first contribution: the link between memorization and  
48 generalization, showing how their interaction impacts a model’s ability to learn robust patterns versus  
49 spurious correlations through controlled experiments. For more on this interaction in other tasks, such  
50 as regression, and how memorization-aware training (MAT), our second contribution, can improve  
51 generalization, see the appendix (Sections C and A). For related work, refer to Section D.

## 52 2 The Interplay between Memorization and Spurious Correlations in ERM

53 **Problem Setup and Preliminaries.** We consider a standard supervised learning setup for a  $K$ -class  
54 classification problem. The data consists of input-label pairs  $(\mathbf{x}, y)$ , where  $\mathbf{x}$  is the input vector and  
55  $y \in \{1, \dots, K\}$  is the class label. Let  $p(y | \mathbf{x})$  denote the training data distribution and let  $a$  denote  
56 any attribute or combination of attributes within  $\mathbf{x}$  that may or may not be relevant for predicting the  
57 target  $y$ .

58 The objective is to learn a model  $\hat{p}(y | \mathbf{x})$  that accurately estimates  $p(y | \mathbf{x})$ . Given input  $\mathbf{x}$ , let  
59  $f(\mathbf{x}; \mathbf{w}) \in \mathbb{R}^K$  be the output logits of a model, where each  $f_k(\mathbf{x}; \mathbf{w})$  represents the logit for class  $k$ .  
60 The estimated conditional probability  $\hat{p}(y = k | \mathbf{x}; \mathbf{w})$  is computed using a softmax function with  
61 temperature  $\tau > 0$ :

$$\hat{p}^{\text{tr}}(y = k | \mathbf{x}; \mathbf{w}) = \frac{\exp(f_k(\mathbf{x}; \mathbf{w})/\tau)}{\sum_{j=1}^K \exp(f_j(\mathbf{x}; \mathbf{w})/\tau)}.$$

62 In order to generalize well under the i.i.d. assumption that  $p(y, \mathbf{x})$  is invariant between training  
63 and test sets, empirical risk minimization (ERM) seeks to minimize the following regularized cross-  
64 entropy loss over a training dataset  $D^{\text{tr}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ :

$$\mathcal{L}^{\text{ERM}} = \frac{1}{n} \sum_{i=1}^n l(y_i, \hat{p}^{\text{tr}}(y | \mathbf{x}_i; \mathbf{w})) + \frac{\lambda}{2} \|\mathbf{w}\|^2,$$

65 where  $l(y_i, \hat{p}^{\text{tr}}(y | \mathbf{x}_i; \mathbf{w})) = -\sum_{k=1}^K \mathbb{I}(y_i = k) \log \hat{p}^{\text{tr}}(y = k | \mathbf{x}_i; \mathbf{w})$  is the cross-entropy loss, and  
66  $\frac{\lambda}{2} \|\mathbf{w}\|^2$  is weight-decay regularization.

67 In cases where there is a distribution shift between training and test, in the presence of spurious  
68 correlations, the i.i.d. assumption breaks down, and when combined with memorization, ERM can  
69 result in poor generalization. We study such scenario in the following section.

70 **Memorization Can Exacerbate Spurious Correlations** Adapting the frameworks introduced in  
71 [Sagawa et al. \(2020\)](#) and [Puli et al. \(2023\)](#), we now study the interplay between memorization and  
72 spurious correlations in an interpretable setup.

73 **Setup 2.1** (Spurious correlations and memorization). *Consider a binary classification problem with*  
74 *labels  $y \in \{-1, +1\}$  and an unknown spurious attribute  $a \in \{-1, +1\}$ . Each input  $\mathbf{x} \in \mathbb{R}^{d+2}$*   
75 *is given by  $\mathbf{x} = (x_y, \gamma x_a, \epsilon)$ , where  $x_y \in \mathbb{R}$  is a core feature dependent only on  $y$ ,  $x_a \in \mathbb{R}$  is a*  
76 *spurious feature dependent only on  $a$ , and  $\epsilon \in \mathbb{R}^d$  are noise features uncorrelated with both  $y$  and  $a$ .*  
77 *The scalar  $\gamma \in \mathbb{R}$  modulates the rate at which the model learns to rely on the spurious feature  $x_a$ ,*  
78 *effectively acting as a scaling factor that increases the feature’s learning rate relative to the core*  
79 *feature  $x_y$ . The attribute  $a$  is considered spurious because it is correlated with the labels  $y$  at training*

80 but has no correlation with  $y$  at test time, potentially leading to poor generalization if the model  
 81 relies on  $x_a$ . Specifically, the data generation process is defined as:

$$\begin{array}{c}
 \begin{array}{c}
 \textcircled{y} \quad \textcircled{a} \\
 \downarrow \quad \downarrow \\
 \textcircled{x_y} \quad \textcircled{\gamma x_a} \\
 \downarrow \quad \downarrow \\
 \textcircled{\epsilon}
 \end{array} \\
 \mathbf{x} := \left( \begin{array}{c} x_y \\ \gamma x_a \\ \epsilon \end{array} \right) \in \mathbb{R}^{d+2}
 \end{array}
 \quad
 \begin{array}{l}
 \begin{cases}
 x_y \sim \mathcal{N}(y, \sigma_y^2) \\
 x_a \sim \mathcal{N}(a, \sigma_a^2) \\
 \epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2 \mathbf{I})
 \end{cases}
 \end{array}
 \quad
 \begin{array}{l}
 a = \begin{cases} y & \text{w.p. } \rho, \\ -y & \text{w.p. } 1 - \rho, \end{cases} \\
 \rho = \begin{cases} \rho^{\text{tr}}, & (\text{train}), \\ 0.5, & (\text{test}). \end{cases}
 \end{array}
 \end{array}$$

82 To better understand this setup, one can think of a classification task between cow and camel  
 83 images. In this example,  $\mathbf{x}$  represents the pixel data,  $y \in \{\text{cow, camel}\}$  are the class labels, and  
 84  $a \in \{\text{grass, sand}\}$  are the background labels. Here,  $x_y$  represents the pixels associated with the  
 85 animal itself (either cow or camel),  $x_a$  represents the pixels associated with the background (grass  
 86 or sand), and  $\epsilon$  represents irrelevant pixels that are specific to each individual example. The key  
 87 assumption is that the joint distribution of class labels and attribute labels differs between training  
 88 and test datasets, i.e.,  $p^{\text{tr}}(a, y) \neq p^{\text{te}}(a, y)$ . For example, in the training set, most cows (camels)  
 89 might appear on grass (sand), while in the test set, cows (camels) appear equally on each background.

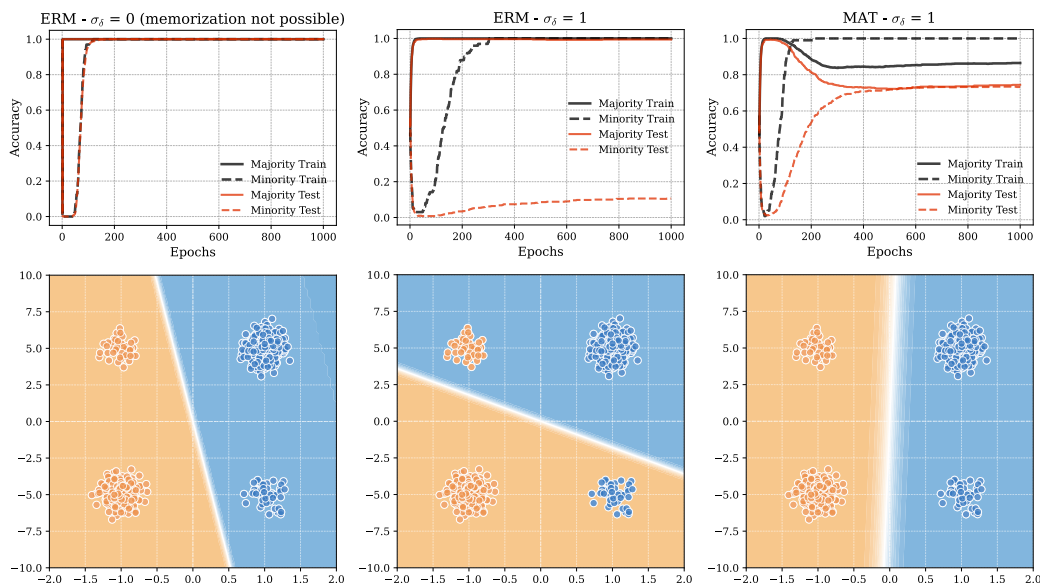


Figure 1: Illustration of two scenarios in the interpretable classification setup involving spurious correlations and memorization. The left panel represents a scenario without input noise ( $\sigma_\epsilon \rightarrow 0$ ), where memorization is not possible. In this case, the model trained with ERM initially learns the spurious feature  $x_a$  serving the majority, but eventually adjusts the decision boundary to the core feature  $x_y$ , resulting in good generalization on minority test examples. The middle and right panels depict a scenario with input noise ( $\sigma_\epsilon \gg 0$ ), where memorization is possible. Here, the model trained with ERM fails to generalize as it memorizes exceptions using the noise features  $\epsilon$  leaving no more incentive for the model to learn the core feature. In contrast, the model trained with MAT (Appendix A) learns the invariant features, and generalizes well even in the presence of noise.

90 **Illustrative Scenarios.** We first empirically study a configuration of the above setup where  $\gamma = 5$   
 91 making the spurious feature easier and faster for the model to learn while being only 90% correlated  
 92 with the class label, i.e.,  $\rho^{\text{tr}} = 0.9$ . In contrast, the core feature  $x_y$  is 100% correlated with  $y$ , but due  
 93 to a smaller norm, it is learned more slowly. Here we consider two cases:

94 1. **Noiseless input  $\Rightarrow$  Spurious Features but No Memorization  $\Rightarrow$  ERM generalizes well.** Figure  
 95 1-(left) presents a case where there are no input noise features ( $\sigma_\epsilon \rightarrow 0$ ). As training progresses,  
 96 the model first learns  $x_a$  due to its larger norm, resulting in perfect accuracy on the majority  
 97 examples. Once the model achieve nearly perfect accuracy on the majority examples, it starts  
 98 to learn the minority examples. At this point, the model must adjust its decision boundary to  
 99 place more emphasis on the core feature  $x_y$ , ultimately achieving perfect generalization on both  
 100 majority and minority examples.

101 2. **Noisy input  $\Rightarrow$  Spurious Features + Memorization  $\Rightarrow$  ERM fails to generalize.**  
 102 Figure 1-(middle) presents a similar setup to the former, but this time with input noise features  
 103 ( $\sigma_\epsilon \gg 0$ ). Again, initially, the model learns the spurious feature  $x_a$ . However, unlike Case 1,  
 104 the noise features  $\epsilon$  provides the model an opportunity to memorize minority examples directly.  
 105 As a result, the model achieves zero training loss by memorizing minority examples using the  
 106 noise dimensions instead of learning to rely on the core feature  $x_y$ . Consequently, the model fails  
 107 to adjust its decision boundary to align with  $x_y$ , and does not generalize on held-out minority  
 108 examples. We argue that most real-world scenarios resemble this case rather than the former case.

109 These results illustrate that the combination of spurious correlations and memorization creates a  
 110 ‘loophole’ for the model. When memorization happens, there is **no more incentives** for the model to  
 111 learn the true, underlying patterns necessary for robust generalization.

112 **Theoretical Analysis.** We now provide a formal analysis to formalize our empirical observations.  
 113 Complete proofs are provided in Appendix E.

114 **Theorem 2.2** (Memorization Exacerbates Spurious Correlations). *Consider a binary classification*  
 115 *problem under the setup described in Setup 2.1, where a linear model  $f(\mathbf{x}; \mathbf{w}) = \mathbf{x}^\top \mathbf{w}$  is trained*  
 116 *using Empirical Risk Minimization (ERM). Let  $\hat{\mathbf{w}}_{ERM} = (\hat{w}_y, \hat{w}_a, \hat{\mathbf{w}}_\epsilon) \in \mathbb{R}^{d+2}$  denote the learned pa-*  
 117 *rameters, where  $\hat{w}_y, \hat{w}_a \in \mathbb{R}$  correspond to the core feature  $x_y$  and spurious feature  $x_a$ , respectively,*  
 118 *and  $\hat{\mathbf{w}}_\epsilon \in \mathbb{R}^d$  corresponds to the noise features  $\epsilon$ .*

119 Assume the following asymptotic conditions hold:

$$\lambda \rightarrow 0^+, n \rightarrow \infty, \lambda\sqrt{n} \rightarrow \infty,$$

120 where  $\lambda > 0$  is the weight decay regularization parameter, and  $n$  is the number of training samples.  
 121 These conditions ensure that ERM converges to the maximum likelihood estimator. For a training  
 122 dataset  $\mathcal{D}^r$  generated under  $\rho^r > 0.5$ , where  $\rho^r$  is the probability that  $a = y$  at training time, the  
 123 following results hold:

124 The ERM-trained classifier  $\hat{y}_{ERM}(\mathbf{x}) = \text{sign}(\mathbf{x}^\top \hat{\mathbf{w}}_{ERM})$  achieves perfect accuracy on all training  
 125 examples:

$$p(\hat{y}_{ERM}(\mathbf{x}) = y \mid \mathbf{x} \in \mathcal{D}^r) \rightarrow 1.$$

126 For held-out (test) examples, denote the classifier as  $\hat{y}_{ERM}^{ho}(\mathbf{x})$ . Then:

127 1. **Noiseless Input Case:** In the noiseless case where the noise variance  $\sigma_\epsilon \rightarrow 0^+$ , the ERM-  
 128 trained classifier converges to a classifier that relies solely on the core feature  $x_y$ . For a  
 129 random test point  $\mathbf{x}$ :

$$p(\hat{y}_{ERM}^{ho}(\mathbf{x}) = y) \rightarrow 1.$$

130 2. **Noisy Input Case:** Suppose  $d \gg \log n$  (where  $d$  is the dimension of the noise features) and  
 131  $\gamma \gg \sigma_\epsilon \sqrt{d/m}$ , where  $m := \rho^r n$  is the number of majority samples in the training set. Then,  
 132 at test time, the ERM-trained classifier  $\hat{y}_{ERM}(\mathbf{x})$  relies pathologically on the spurious feature  
 133  $x_a$ . For a random test point  $\mathbf{x}$ :

$$p(\hat{y}_{ERM}^{ho}(\mathbf{x}) = a) \rightarrow 1.$$

134 The condition  $d \gg \log n$  ensures that noise features from different samples are approximately  
 135 orthogonal, and  $\gamma \gg \sigma_\epsilon \sqrt{d/m}$  guarantees that the spurious feature  $x_a$  is learned faster by gradient  
 136 descent than other features.

### 137 3 Discussion

138 In our first contribution, we showed that spurious features *alone* do not solely cause poor generaliza-  
 139 tion. Instead, memorization features remove the incentive for the model to learn the true underlying  
 140 patterns from minority cases. However, to achieve our main goal of learning generalizable patterns, it  
 141 is crucial to provide the model with feedback on its failures. Held-out performance, which is free  
 142 of memorization, offers a way to achieve this. To address this, we propose MAT (Memorization-  
 143 Aware Training), a method that adjusts model logits during training to encourage the learning of  
 144 generalizable features. Details of MAT are provided in Appendix A. In addition, we introduce and  
 145 analyze three types of memorization—bad, good, and ugly—highlighting their effects and relevance  
 146 in different scenarios of benign and malign overfitting. This discussion is elaborated in Appendix C.

147 **References**

- 148 Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization.  
149 *arXiv preprint arXiv:1907.02893*, 2019.
- 150 Devansh Arpit, Stanisław Jastrzębski, Nicolas Ballas, David Krueger, Emmanuel Bengio, Maxinder S  
151 Kanwal, Tegan Maharaj, Asja Fischer, Aaron Courville, Yoshua Bengio, et al. A closer look at  
152 memorization in deep networks. In *International conference on machine learning*, pp. 233–242.  
153 PMLR, 2017.
- 154 Peter L Bartlett, Philip M Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear  
155 regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020.
- 156 Mikhail Belkin, Daniel Hsu, Siyuan Ma, and Soumik Mandal. Reconciling modern machine-learning  
157 practice and the classical bias–variance trade-off. *Proceedings of the National Academy of Sciences*,  
158 116(32):15849–15854, 2019a.
- 159 Mikhail Belkin, Alexander Rakhlin, and Alexandre B Tsybakov. Does data interpolation contradict  
160 statistical optimality? In *The 22nd International Conference on Artificial Intelligence and Statistics*,  
161 pp. 1611–1619. PMLR, 2019b.
- 162 Simone Bombari and Marco Mondelli. How spurious features are memorized: Precise analysis for  
163 random and ntk features. In *Forty-first International Conference on Machine Learning*, 2024.
- 164 Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced  
165 datasets with label-distribution-aware margin loss. *Advances in neural information processing  
166 systems*, 32, 2019.
- 167 Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and  
168 Chiyuan Zhang. Quantifying memorization across neural language models. *arXiv preprint  
169 arXiv:2202.07646*, 2022.
- 170 Guillem Collell, Drazen Prelec, and Kaustubh Patil. Reviving threshold-moving: a simple plug-in  
171 bagging ensemble for binary and multiclass imbalanced data. *arXiv preprint arXiv:1606.08698*,  
172 2016.
- 173 Elliot Creager, Jörn-Henrik Jacobsen, and Richard Zemel. Environment inference for invariant  
174 learning. In *International Conference on Machine Learning*, pp. 2189–2200. PMLR, 2021.
- 175 Nikolay Dageev, Brett D Roads, Xiaoliang Luo, Daniel N Barry, Kaustubh R Patil, and Bradley C  
176 Love. A too-good-to-be-true prior to reduce shortcut reliance. *Pattern recognition letters*, 166:  
177 164–171, 2023.
- 178 B. Dherin, M. Munn, M. Rosca, and D. Barrett. Why neural networks find simple solutions: the  
179 many regularizers of geometric complexity. *Neural Information Processing Systems*, 2022. doi:  
180 10.48550/arXiv.2209.13083.
- 181 Vitaly Feldman and Chiyuan Zhang. What neural networks memorize and why: Discovering the long  
182 tail via influence estimation. *Advances in Neural Information Processing Systems*, 33:2881–2891,  
183 2020.
- 184 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias  
185 Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine  
186 Intelligence*, 2(11):665–673, 2020.
- 187 Zongbo Han, Zhipeng Liang, Fan Yang, Liu Liu, Lanqing Li, Yatao Bian, Peilin Zhao, Bingzhe Wu,  
188 Changqing Zhang, and Jianhua Yao. Umix: Improving importance weighting for subpopulation  
189 shift via uncertainty-aware mixup. *Advances in Neural Information Processing Systems*, 35:  
190 37704–37718, 2022.
- 191 Reyhane Askari Hemmat, Mohammad Pezeshki, Florian Bordes, Michal Drozdal, and Adriana  
192 Romero-Soriano. Feedback-guided data synthesis for imbalanced classification. *arXiv preprint  
193 arXiv:2310.00158*, 2023.

- 194 Badr Youbi Idrissi, Martin Arjovsky, Mohammad Pezeshki, and David Lopez-Paz. Simple data  
195 balancing achieves competitive worst-group-accuracy. *CLear*, 2022. URL <https://arxiv.org/abs/2110.14503>.  
196
- 197 Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis  
198 Kalantidis. Decoupling representation and classifier for long-tailed recognition. *arXiv preprint*  
199 *arXiv:1910.09217*, 2019.
- 200 Byungju Kim and Junmo Kim. Adjusting decision boundary for class imbalanced learning. *IEEE*  
201 *Access*, 8:81674–81685, 2020.
- 202 Byungju Kim, Hyunwoo Kim, Kyungsu Kim, Sungjin Kim, and Junmo Kim. Learning not to learn:  
203 Training deep neural networks with biased data. In *Proceedings of the IEEE/CVF conference on*  
204 *computer vision and pattern recognition*, pp. 9012–9020, 2019.
- 205 Polina Kirichenko, Pavel Izmailov, and Andrew Gordon Wilson. Last layer re-training is sufficient  
206 for robustness to spurious correlations. *arXiv preprint arXiv:2204.02937*, 2022.
- 207 David Krueger, Nicolas Ballas, Stanislaw Jastrzebski, Devansh Arpit, Maxinder S Kanwal, Tegan  
208 Maharaj, Emmanuel Bengio, Asja Fischer, and Aaron Courville. Deep nets don’t learn via  
209 memorization. 2017.
- 210 Tyler LaBonte, Vidya Muthukumar, and Abhishek Kumar. Towards last-layer retraining for group  
211 robustness with fewer annotations. *Advances in Neural Information Processing Systems*, 36, 2024.
- 212 Yi Li and Nuno Vasconcelos. Repair: Removing representation bias by dataset resampling. In  
213 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 9572–  
214 9581, 2019.
- 215 Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa,  
216 Percy Liang, and Chelsea Finn. Just train twice: Improving group robustness without training  
217 group information. In *International Conference on Machine Learning*, pp. 6781–6792. PMLR,  
218 2021.
- 219 Sheng Liu, Xu Zhang, Nitesh Sekhar, Yue Wu, Prateek Singhal, and Carlos Fernandez-Granda.  
220 Avoiding spurious correlations via logit correction. *arXiv preprint arXiv:2212.01433*, 2022.
- 221 Pratyush Maini, Saurabh Garg, Zachary Lipton, and J Zico Kolter. Characterizing datapoints via  
222 second-split forgetting. *Advances in Neural Information Processing Systems*, 35:30044–30057,  
223 2022.
- 224 Pratyush Maini, Michael C Mozer, Hanie Sedghi, Zachary C Lipton, J Zico Kolter, and Chiyuan  
225 Zhang. Can neural network memorization be localized? *arXiv preprint arXiv:2307.09542*, 2023.
- 226 Neil Mallinar, James Simon, Amirhesam Abedsoltan, Parthe Pandit, Misha Belkin, and Preetum  
227 Nakkiran. Benign, tempered, or catastrophic: Toward a refined taxonomy of overfitting. *Advances*  
228 *in Neural Information Processing Systems*, 35:1182–1195, 2022.
- 229 Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit, and  
230 Sanjiv Kumar. Long-tail learning via logit adjustment. *arXiv preprint arXiv:2007.07314*, 2020.
- 231 Vidya Muthukumar, Kailas Vodrahalli, Vignesh Subramanian, and Anant Sahai. Harmless interpo-  
232 lation of noisy data in regression. *IEEE Journal on Selected Areas in Information Theory*, 1(1):  
233 67–83, 2020.
- 234 Vaishnavh Nagarajan, Anders Andreassen, and Behnam Neyshabur. Understanding the failure modes  
235 of out-of-distribution generalization. *arXiv preprint arXiv:2010.15775*, 2020.
- 236 Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep  
237 double descent: Where bigger models and more data hurt. *Journal of Statistical Mechanics: Theory*  
238 *and Experiment*, 2021(12):124003, 2021.
- 239 Junhyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. Learning from failure:  
240 De-biasing classifier from biased classifier. *Advances in Neural Information Processing Systems*,  
241 33:20673–20684, 2020.

- 242 Sung Min Park, Kristian Georgiev, Andrew Ilyas, Guillaume Leclerc, and A. Madry. Trak: Attributing  
243 model behavior at scale. *International Conference on Machine Learning*, 2023. doi: 10.48550/  
244 arXiv.2303.14186.
- 245 Mohammad Pezeshki, Oumar Kaba, Yoshua Bengio, Aaron C Courville, Doina Precup, and Guil-  
246 laume Lajoie. Gradient starvation: A learning proclivity in neural networks. *Advances in Neural*  
247 *Information Processing Systems*, 34:1256–1272, 2021.
- 248 Mohammad Pezeshki, Diane Bouchacourt, Mark Ibrahim, Nicolas Ballas, Pascal Vincent, and David  
249 Lopez-Paz. Discovering environments with xrm. *arXiv preprint arXiv:2309.16748*, 2023.
- 250 Ahlad Manas Puli, Lily Zhang, Yoav Wald, and Rajesh Ranganath. Don’t blame dataset shift!  
251 shortcut learning due to gradients and cross entropy. *Advances in Neural Information Processing*  
252 *Systems*, 36:71874–71910, 2023.
- 253 Jiawei Ren, Cunjun Yu, Xiao Ma, Haiyu Zhao, Shuai Yi, et al. Balanced meta-softmax for long-tailed  
254 visual recognition. *Advances in neural information processing systems*, 33:4175–4186, 2020.
- 255 Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for  
256 robust deep learning. In *International conference on machine learning*, pp. 4334–4343. PMLR,  
257 2018.
- 258 Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust  
259 neural networks for group shifts: On the importance of regularization for worst-case generalization.  
260 *arXiv preprint arXiv: 1911.08731*, 2019.
- 261 Shiori Sagawa, Aditi Raghunathan, Pang Wei Koh, and Percy Liang. An investigation of why  
262 overparameterization exacerbates spurious correlations. In *International Conference on Machine*  
263 *Learning*, pp. 8346–8356. PMLR, 2020.
- 264 Avi Schwarzschild, Zhili Feng, Pratyush Maini, Zachary C Lipton, and J Zico Kolter. Rethinking  
265 llm memorization through the lens of adversarial compression. *arXiv preprint arXiv:2404.15146*,  
266 2024.
- 267 Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The  
268 pitfalls of simplicity bias in neural networks. *Advances in Neural Information Processing Systems*,  
269 33:9573–9585, 2020.
- 270 Cory Stephenson, Suchismita Padhy, Abhinav Ganesh, Yue Hui, Hanlin Tang, and SueYeon Chung.  
271 On the geometry of generalization and memorization in deep neural networks. *arXiv preprint*  
272 *arXiv:2105.14602*, 2021.
- 273 Christos Tsirigotis, Joao Monteiro, Pau Rodriguez, David Vazquez, and Aaron C Courville. Group  
274 robust classification without any group information. *Advances in Neural Information Processing*  
275 *Systems*, 36, 2024.
- 276 Yoav Wald, Gal Yona, Uri Shalit, and Yair Carmon. Malign overfitting: Interpolation and invariance  
277 are fundamentally at odds. In *NeurIPS 2022 Workshop on Distribution Shifts: Connecting Methods*  
278 *and Applications*, 2022.
- 279 Yao-Yuan Yang, Chi-Ning Chou, and Kamalika Chaudhuri. Understanding rare spurious correlations  
280 in neural networks. *arXiv preprint arXiv:2202.05189*, 2022.
- 281 Zitong Yang, Michal Lukasik, Vaishnavh Nagarajan, Zonglin Li, Ankit Rawat, Manzil Zaheer,  
282 Aditya K Menon, and Sanjiv Kumar. Resmem: Learn what you can and memorize the rest.  
283 *Advances in Neural Information Processing Systems*, 36, 2024.
- 284 Huaxiu Yao, Yu Wang, Sai Li, Linjun Zhang, Weixin Liang, James Zou, and Chelsea Finn. Improving  
285 out-of-distribution robustness via selective augmentation. In *International Conference on Machine*  
286 *Learning*, pp. 25407–25437. PMLR, 2022.
- 287 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep  
288 learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115,  
289 2021.

290 Michael Zhang, Nimit S Sohoni, Hongyang R Zhang, Chelsea Finn, and Christopher Ré. Correct-n-  
291 contrast: A contrastive approach for improving robustness to spurious correlations. *arXiv preprint*  
292 *arXiv:2203.01517*, 2022.



293 **A Memorization-Aware Training (MAT)**

294 As exemplified in Section 2, the i.i.d. assumption underlying ERM is violated in the presence of  
 295 *spurious correlations* between the label  $y$  and certain attributes  $a$ . If a classifier  $\hat{y}(\mathbf{x})$  relies on these  
 296 unstable correlations, it may fail to generalize to test data where  $p^{\text{tr}}(y, a) \neq p^{\text{te}}(y, a)$ . To address this  
 297 distribution shift, we propose an alternative assumption.

298 **A.1 An Alternative to the i.i.d. Assumption**

299 We assume that any predictive path involving  $a$  is unreliable because  $p(y, a)$  changes between training  
 300 and test. To focus on the stable relationship independent of  $a$ , we introduce an *invariant quantity* that  
 301 remains consistent across both distributions. Specifically, we assume that:

$$\frac{p^{\text{tr}}(y | \mathbf{x})}{\sum_a p^{\text{tr}}(y | a)p^{\text{tr}}(a | \mathbf{x})} \propto \frac{p^{\text{te}}(y | \mathbf{x})}{\sum_a p^{\text{te}}(y | a)p^{\text{te}}(a | \mathbf{x})}, \quad (1)$$

302 where the term  $\sum_a p(y | a)p(a | \mathbf{x})$  represents the conditional probability of  $y$  given  $\mathbf{x}$  when passing  
 303 through  $a$  as an intermediate attribute.

304 **Deriving a New Learning Algorithm.** To derive a new learning algorithm based on this assumption,  
 305 we aim to express  $p^{\text{tr}}(y | \mathbf{x})$  in terms of  $p^{\text{te}}(y | \mathbf{x})$ . Starting from the assumption in Equation 1, we  
 306 have:

$$\frac{p^{\text{tr}}(y | \mathbf{x})}{\sum_a p^{\text{tr}}(y | a)p^{\text{tr}}(a | \mathbf{x})} \propto \frac{p^{\text{te}}(y | \mathbf{x})}{\sum_a p^{\text{te}}(y | a)p^{\text{te}}(a | \mathbf{x})}, \quad (\text{by assumption}) \quad (2)$$

$$\Rightarrow \frac{p^{\text{tr}}(y | \mathbf{x})}{\sum_a p^{\text{tr}}(y | a)p^{\text{tr}}(a | \mathbf{x})} \propto \frac{p^{\text{te}}(y | \mathbf{x})}{p^{\text{te}}(y)}, \quad (\text{assuming } y \perp a \text{ in test set}) \quad (3)$$

$$\Rightarrow p^{\text{tr}}(y | \mathbf{x}) \propto p^{\text{te}}(y | \mathbf{x}) \sum_a p^{\text{tr}}(y | a)p^{\text{tr}}(a | \mathbf{x}), \quad (\text{assuming } p^{\text{te}}(y) \sim \mathcal{U}) \quad (4)$$

$$\Rightarrow p^{\text{tr}}(y | \mathbf{x}) \propto p^{\text{te}}(y | \mathbf{x}) p_a^{\text{tr}}(y | \mathbf{x}), \quad (\text{change of variable}) \quad (5)$$

$$\Rightarrow p^{\text{tr}}(y | \mathbf{x}) = \frac{p^{\text{te}}(y | \mathbf{x}) p_a^{\text{tr}}(y | \mathbf{x})}{\sum_{y'} p^{\text{te}}(y' | \mathbf{x}) p_a^{\text{tr}}(y' | \mathbf{x})}, \quad (\text{normalization so it sums to 1}) \quad (6)$$

307 where  $p_a^{\text{tr}}(y | \mathbf{x})$  is the correction term accounting for the prediction of label  $y$  given  $\mathbf{x}$  that goes  
 308 through  $a$ :

$$p_a^{\text{tr}}(y | \mathbf{x}) := \sum_a p^{\text{tr}}(y | a)p^{\text{tr}}(a | \mathbf{x}).$$

309 Instead of directly estimating  $p^{\text{tr}}(y | \mathbf{x})$ , we estimate  $p^{\text{te}}(y | \mathbf{x})$  using a softmax on the logits of a  
 310 model. Thus, the expression for  $\hat{p}^{\text{tr}}(y = k | \mathbf{x})$  is:

$$\hat{p}^{\text{tr}}(y = k | \mathbf{x}) = \frac{\exp(f_k(\mathbf{x}; \mathbf{w}) + \log p_a^{\text{tr}}(y = k | \mathbf{x}))}{\sum_{y'} \exp(f_{y'}(\mathbf{x}; \mathbf{w}) + \log p_a^{\text{tr}}(y' | \mathbf{x}))}, \quad (\text{see Lemma E.1}). \quad (7)$$

311 Equation 7 proposes adjusting the logits of a model to account for the fact that  $p_a(y | \mathbf{x})$  is unreliable  
 312 and varies from training to test. This formulation is related to prior work on logit adjustment (Kang  
 313 et al., 2019; Menon et al., 2020; Ren et al., 2020; Liu et al., 2022; Tsirigotis et al., 2024), but differs  
 314 in how the adjustment is computed.

315 **A.2 Estimating  $p_a(y | \mathbf{x})$  Using Held-Out Predictions**

316 In Section 2, we showed that a model trained with ERM under Setup 2.1 tends to rely heavily on  
 317 spurious attributes when evaluated on held-out data. Specifically, for a given input  $\mathbf{x}$ , the predicted  
 318 label  $\hat{y}_{\text{ERM}}^{\text{ho}}$  aligns almost exclusively with the spurious attribute  $a$ , implying  $p(\hat{y}_{\text{ERM}}^{\text{ho}} = a | \mathbf{x}) \rightarrow 1$ .  
 319 This implies that for a specific  $a = a^*$ ,  $p(a^* | \mathbf{x}) \approx 1$  and  $p(a | \mathbf{x}) \approx 0$  for all other  $a \neq a^*$ .

Table 1: Average/worst accuracies comparing methods for environment discovery. We specify access to annotations in training data ( $e^{\text{tr}}$ ) and validation data ( $e^{\text{va}}$ ). Symbol  $\sim$  denotes inferred group annotations by the method, and symbol  $\dagger$  denotes original numbers.

			Waterbirds		CelebA		MNLI		CivilComms	
$e^{\text{tr}}$	$e^{\text{va}}$		Avg	Worst	Avg	Worst	Avg	Worst	Avg	Worst
$\times$	$\checkmark$	ERM	83.8	66.4	95.5	55.1	81.6	72.0	84.3	74.0
$\checkmark$	$\checkmark$	GroupDRO	90.2	86.5	93.1	88.3	80.6	73.4	84.2	73.8
$\times$	$\checkmark$	LC $\dagger$	-	90.5	-	88.1	-	-	-	70.3
$\times$	$\checkmark$	MAT	89.4	88.2	88.0	85.6	TBD	TBD	TBD	TBD
$\times$	$\times$	ERM	83.6	66.4	95.3	58.6	81.8	69.1	81.5	64.7
$\times$	$\sim$	uLA $\dagger$	91.5	86.1	93.9	86.5	-	-	-	-
$\sim$	$\sim$	XRM+GroupDRO $\dagger$	89.3	88.1	91.4	89.1	75.8	72.1	84.0	72.2
$\times$	$\sim$	MAT	TBD	TBD	TBD	TBD	TBD	TBD	TBD	TBD

320 This observation simplifies the estimation of  $p_a^{\text{tr}}(y | \mathbf{x})$  as:

$$p_a^{\text{tr}}(y | \mathbf{x}) = \sum_a p^{\text{tr}}(y | a) p^{\text{tr}}(a | \mathbf{x}) \approx p^{\text{tr}}(y | a^*),$$

321 where  $a^* = \arg \max_a p(\hat{y}_{\text{ERM}}^{\text{ho}} = a | \mathbf{x})$ .

322 To compute  $p^{\text{tr}}(y | a^*)$ , following (Liu et al., 2022), we use the empirical counts from the training  
323 data:

$$p^{\text{tr}}(y | a^*) = \frac{\text{count}(y, a^*)}{\text{count}(a^*)}.$$

324 Thus, the held-out predictions of the ERM model provide a straightforward way to estimate  $p_a^{\text{tr}}(y | \mathbf{x})$ ,  
325 allowing us to adjust the model logits accordingly for improved generalization.

326 Specifically, MAT employs a shared backbone network with three classification heads:

- 327 • **Heads A and B:** These heads are trained on two random non-overlapping splits of the training  
328 data using ERM. Each head provides held-out predictions for the other head’s split, from  
329 which we estimate  $p^{\text{tr}}(y | a^*)$ .
- 330 • **Head C:** This is the main classifier whose logits are adjusted using Equation 7, based on the  
331 held-out predictions from Heads A and B.

332 During training, all three heads—A, B, and C—are updated simultaneously. Heads A and B optimize  
333 only their head parameters. Head C updates its own parameters as well as those of the shared  
334 backbone. To further reinforce reliance on spurious correlations, we employ *Label Flipping* strategy  
335 (Pezeshki et al., 2023) on Heads A and B. Flipping is done according to held-out probabilities and  
336 hence amplifies the biases of the auxiliary classifiers.

## 337 B Experiments

338 We first, conduct experiments to demonstrate the effectiveness of Memorization-Aware Training  
339 (MAT) in improving generalization under subpopulation shift. We then provide a detailed analysis of  
340 the memorization behaviors of models trained with ERM and MAT.

### 341 B.1 Experiments on Subpopulation Shift

342 We evaluate our approach on four datasets under subpopulation shift. In all experiments, we assume  
343 that spurious correlation or environment annotations are not available during training. We consider  
344 two settings: (1) group annotations are available in the validation set for model selection, and (2) no  
345 annotations are available even in the validation set.

346 For evaluation, we report two key metrics on the test set: (1) average test accuracy and (2) worst-group  
347 test accuracy, the latter being computed using ground-truth annotations.

348 Table 1 compares the performance of MAT with several baseline methods, including ERM, GroupDRO  
 349 (Sagawa et al., 2019), and other environment discovery methods like LC (Liu et al., 2022), uLA  
 350 (Tsirigotis et al., 2024) and XRM+GroupDRO (Pezeshki et al., 2023). These methods vary in their  
 351 assumptions about access to annotations, both in training and validation for model selection. For  
 352 instance, ERM does not assume any training group annotations, while GroupDRO has full access to  
 353 group annotations for both training and validation data.

354 In the Waterbirds dataset, MAT demonstrates strong performance with 88.2% worst-group accuracy  
 355 when the ground truth group annotations of the validation set are used for model selection, improving  
 356 substantially over ERM. Similarly, on the CelebA dataset, MAT achieves competitive results, with a  
 357 worst-group accuracy of 85.6%. These results suggest that MAT’s memorization-aware approach  
 358 effectively mitigates overfitting to spurious correlations, particularly in challenging worst-group  
 359 scenarios.

## 360 B.2 Analysis of Memorization Scores

361 To understand the extent of memorization in models trained with ERM, we analyze the distribution  
 362 of memorization scores across subpopulations. We focus on the Waterbird dataset, which includes  
 363 two main classes—Waterbird and Landbird—each divided into majority and minority subpopulations  
 364 based on their background (e.g., Waterbird on water vs. Waterbird on land). This setup allows us to  
 365 investigate how memorization varies with group size and context.

366 The memorization score is derived from the influence function, which measures the effect of each  
 367 training sample on a model’s prediction. Formally, the influence of a training sample  $i$  on a target  
 368 sample  $j$  under a training algorithm  $\mathcal{A}$  is defined as:

$$\text{infl}(\mathcal{A}, \mathcal{D}, i, j) := \hat{p}_{\mathcal{D}}^{(\mathcal{A})}(y_j | \mathbf{x}_j) - \hat{p}_{\mathcal{D}_{-(\mathbf{x}_i, y_i)}}^{(\mathcal{A})}(y_j | \mathbf{x}_j) \quad (8)$$

369 where  $\mathcal{D}$  is the training dataset,  $\mathcal{D}_{-(\mathbf{x}_i, y_i)}$  denotes the dataset with the sample  $(\mathbf{x}_i, y_i)$  removed. The  
 370 memorization score is a specific case of this function where the target sample  $(\mathbf{x}_j, y_j)$  is the same as  
 371 the training sample. It measures the difference between a model’s performance on a training sample  
 372 when that sample is included in the training set (held-in) versus when it is excluded (held-out).

373 Calculating self-influence scores with a naive leave-one-out approach is computationally expensive,  
 374 but recent methods like TRAK (Park et al., 2023) provide an efficient alternative. TRAK approximates  
 375 the data attribution matrix. The diagonal of this matrix directly gives the self-influence scores.

376 Figure 2 depicts the distribution of self-influence scores across subpopulations in the Waterbird  
 377 dataset. We note that minority subpopulations (e.g., Waterbirds on land) show higher self-influence  
 378 scores compared to their majority counterparts (e.g., Waterbirds on water) for a model trained with  
 379 ERM. A model trained with MAT, however, shows a similar distribution of self-influence for both the  
 380 majority and minority examples.

## 381 C Memorization: The Good, the Bad, and the Ugly

382 We showed that the combination of memorization and spurious correlations, rather than spurious  
 383 correlations alone, could be key reason for poor generalization. Neural networks can exploit spurious  
 384 features and memorize exceptions to achieve zero training loss, thereby avoiding learning more  
 385 generalizable patterns. However, an interesting and somewhat controversial question arises: *Is*  
 386 *memorization always bad?*

387 To explore this, we look into a simple regression task to understand different types of memorization  
 388 and their effects on generalization. We argue that the impact of memorization on generalization can  
 389 vary depending on the nature of the data and the model’s learning dynamics, and we categorize these  
 390 types of memorization into three distinct forms.

391 **Setup C.1** (Regression with Memorization). *Let  $x_y \in \mathbb{R}$  be a scalar feature that determines the*  
 392 *true target,  $y^* = f(x_y)$ . Let  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  be a dataset consisting of input-target pairs  $(\mathbf{x}, y)$ .*  
 393 *Define the input vector as  $\mathbf{x} = \text{concat}(x_y, \epsilon) \in \mathbb{R}^{m+1}$ , where  $\epsilon \sim \mathcal{N}(0, \sigma_\epsilon^2 I) \in \mathbb{R}^m$  represents input*  
 394 *noise concatenated with the true feature  $x_y$ . The target is defined as  $y = y^* + \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, \sigma_\epsilon)$*   
 395 *represents additive target noise.*

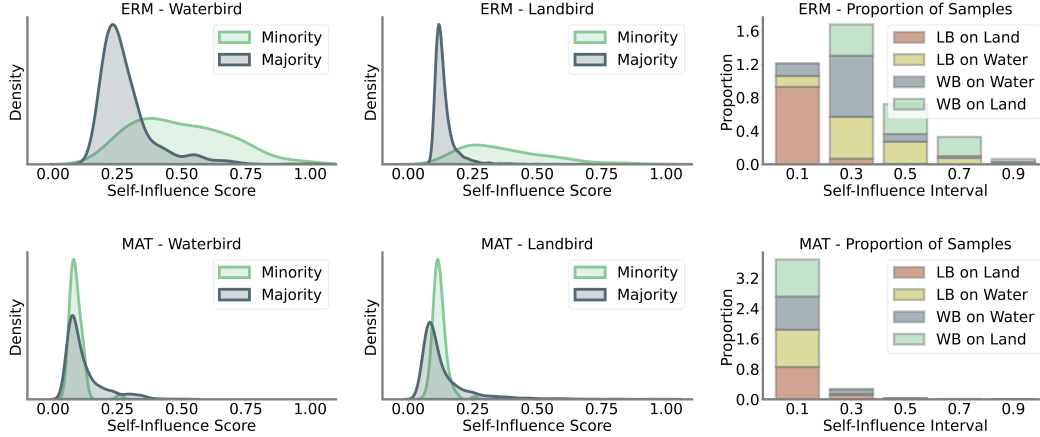


Figure 2: Self-Influence estimation of the Waterbird groups by ERM and MAT. The distribution of self-influence scores is shown for both the majority and minority subpopulations (e.g., Waterbirds on water vs. Waterbirds on land). Models trained with ERM exhibit higher self-influence scores for minority subpopulations, suggesting increased memorization in these groups. In contrast, models trained with MAT show more uniform self-influence distributions across both majority and minority subpopulations. The rightmost plots display the proportion of samples in different self-influence intervals, with MAT producing a more balanced distribution compared to ERM.

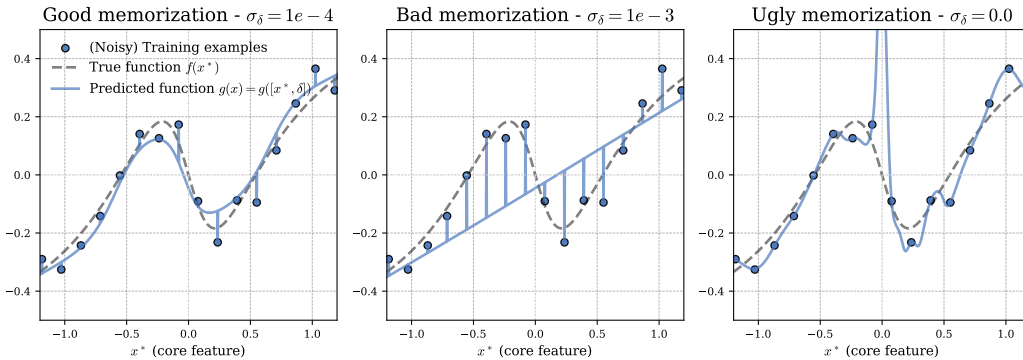


Figure 3: Three types of memorization in regression models trained with different levels of input noise ( $\sigma_\epsilon$ ). The plots show the ERM-trained model  $g(x) = g(x_y, \epsilon)$  (solid blue line) versus the true underlying function  $f(x_y)$  (dashed gray line) and the noisy training examples. In all the three, the models are trained until the training loss goes below  $1e-6$ . **Good memorization (Left,  $\sigma_\epsilon = 1e-4$ ):** Model learns the true function  $f(x_y)$  well but slightly memorizes residual noise in the training data using the input noise  $\epsilon$ . This type of memorization is benign, as it does not compromise generalization. **Bad memorization (Middle,  $\sigma_\epsilon = 1e-3$ ):** The model relies more on noise features than learning the true function  $f(x_y)$ , leading to partial learning of  $f(x_y)$  and fitting of noise-dominated input features. This type of memorization impedes learning of generalizable patterns and is considered malign. **Ugly memorization (Right,  $\sigma_\epsilon = 0.0$ ):** Without input noise, the model overfits the training data, including label noise, resulting in a highly non-linear and complex model that fails to generalize to new data. This type is referred to as catastrophic overfitting.

396 In this context,  $x_y$  can be interpreted as the core feature (e.g., the object in an object classification  
 397 task),  $\epsilon$  as irrelevant random noise, and  $\epsilon$  as labeling noise or error. Now, consider training linear  
 398 regression models  $\hat{y} = g(x)$  on this dataset. Fixing  $\sigma_\epsilon$ , we train three models under three different  
 399 input noise levels:  $\sigma_\epsilon \in \{0, 0.01, 0.1\}$ . The results, summarized in Figure 3, showcases three types  
 400 of memorization:

401 **The Good.** At an intermediate level of input noise,  $\sigma_\epsilon = 1e-4$ , the model effectively captures  
 402 the true underlying function,  $f(x_y)$ . However, due to the label noise, the model cannot achieve a

403 zero training loss solely by learning  $f(x_y)$ . As a result, it begins to memorize the residual noise in  
404 the training data by using the input noise  $\epsilon$ . This is evidenced by sharp spikes at each training point,  
405 where the model,  $g(x)$ , precisely predicts the noisy label if given the exact same input as during  
406 training. Nevertheless, for a neighboring test example with no input noise, the model’s predictions  
407 align well with  $f(x_y)$ , demonstrating good generalization.

408 This phenomenon is often referred to as “benign overfitting” where a model can perfectly fit (overfit  
409 in fact) the training data while relying on noise and unreliable features, yet still generalize well to  
410 unseen data (Belkin et al., 2019a; Muthukumar et al., 2020; Bartlett et al., 2020). The key insight is  
411 that the overfitting in this case is “benign” because the model’s memorization by relying on noise  
412 features does not compromise the underlying structure of the true signal. Instead, the model retains  
413 a close approximation to the true function on test data, even though it memorizes specific noise in  
414 the training data. This has been shown to occur particularly in over-parameterized neural networks  
415 (Belkin et al., 2019b; Nakkiran et al., 2021).

416 **The Bad.** At a higher level of input noise,  $\sigma_\epsilon = 1e - 3$ , the model increasingly rely on the  
417 input noise features  $\epsilon$  rather than fully learning the true underlying function  $f(x_y)$ . In this case,  
418 memorization is more tempting for the model because the noise dominates the input, making it  
419 difficult to recover the true signal. As a result, the model  $g(x)$  might achieve zero training loss by  
420 only partially learning  $f(x_y)$  and instead relying heavily on the noise in the inputs to fit the remaining  
421 variance in the training data.

422 This is an instance of bad memorization as it hinders the learning of generalizable patterns. It becomes  
423 particularly problematic when the data contains spurious correlations. A model can achieve zero  
424 training loss by relying on a combination of spurious correlations and memorization of any errors  
425 that are not already satisfied by the spurious correlation. This phenomenon is referred to as "malign  
426 overfitting" in Wald et al. (2022), where a model perfectly fits the training data but in a way that  
427 compromises its ability to generalize, especially in situations where robustness, fairness, or invariance  
428 are critical.

429 It is important to note that both good and bad memorization stem from the same learning dynamics.  
430 ERM, and the SGD that drives it, do not differentiate between the types of correlations or features  
431 they are learning. Whether a features contributes to generalization or memorization is only revealed  
432 when the model is evaluated on held-out data. If the features learned are generalizable, the model  
433 will perform well on new data; if they are not, the model will struggle, revealing its reliance on  
434 memorized, non-generalizable patterns.

435 **The Ugly.** Finally, consider the case where there is no input noise,  $\sigma_\epsilon = 0.0$ . In this case, the  
436 model may initially capture the true function  $f(x_y)$ , but due to the presence of label noise, it cannot  
437 achieve zero training loss by learning only  $f(x_y)$ . Unlike the previous cases, the absence of input  
438 noise means the model has no additional features to leverage in explaining the residual error. As a  
439 result, the model is forced to learn a highly non-linear and complex function of the input  $x = x_y$  to  
440 fit the noisy labels.

441 In this situation, memorization is ugly: The model may achieve perfect predictions on the training  
442 data, but this comes at the cost of catastrophic overfitting— where the model overfits so severely that  
443 it not only memorizes every detail of the training data, including noise, but also loses its ability to  
444 generalize to new data (Mallinar et al., 2022).

445 These examples show that memorization is not always bad; its impact varies with the nature of  
446 the data. While MAT mitigates the negative effects of memorization in the presence of spurious  
447 correlations, there are cases where memorization can benefit generalization or even be essential  
448 (Feldman & Zhang, 2020). Future work could focus on distinguishing these scenarios and exploring  
449 the nuanced role of memorization in large language models (LLMs). Recent work (Carlini et al.,  
450 2022; Schwarzschild et al., 2024) have highlighted the importance of defining and understanding  
451 memorization in LLMs, as it can inform how these models balance between storing training data and  
452 learning generalizable patterns.

## 453 D Related Work

454 **Detecting Spurious Correlations.** Early methods for detecting spurious correlations rely on human  
455 annotations (Kim et al., 2019; Sagawa et al., 2019; Li & Vasconcelos, 2019), which are costly and  
456 susceptible to bias. Without explicit annotations, detecting spurious correlations requires assumptions.  
457 A common assumption is that spurious correlations are learned more quickly or are simpler to learn  
458 than core features (Geirhos et al., 2020; Arjovsky et al., 2019; Sagawa et al., 2020). Based on  
459 this, methods like Just Train Twice (JTT) (Liu et al., 2021), Environment Inference for Invariant  
460 Learning (EIIL) (Creager et al., 2021), Too-Good-To-Be-True Prior (Dagaev et al., 2023), and  
461 Correct-n-Contrast (CnC) (Zhang et al., 2022) train models with limited capacity to identify "hard"  
462 (minority) examples. Other methods such as Learning from Failure (LfF) (Nam et al., 2020) and  
463 Logit Correction (LC) (Liu et al., 2022) use generalized cross-entropy to bias classifiers toward  
464 spurious features. Closely related to this work is Cross-Risk Minimization (XRM) (Pezeshki et al.,  
465 2023), where uses the held-out mistakes as a signal for the spurious correlations.

466 **Mitigating Spurious Correlations.** Reweighting, resampling, and retraining techniques are widely  
467 used to enhance minority group performance by adjusting weights or sampling rates (Idrissi et al.,  
468 2022; Nagarajan et al., 2020; Ren et al., 2018). Methods like Deep Feature Reweighting (DFR)  
469 (Kirichenko et al., 2022) and Selective Last-Layer Finetuning (SELF) (LaBonte et al., 2024) retrain  
470 the last layer on balanced or selectively sampled data. GroupDRO (Sagawa et al., 2019) minimizes  
471 worst-case group loss, while approaches like LfF and JTT increase loss weights for likely minority  
472 examples. Data balancing can also be achieved through data synthesis, feature augmentation, or  
473 domain mixing (Hemmat et al., 2023; Yao et al., 2022; Han et al., 2022).

474 *Logit adjustment* methods are crucial for robust classification under biased training conditions.  
475 Menon et al. (2020) propose a method that corrects model predictions based on class frequencies,  
476 building on prior work in post-hoc adjustments (Collell et al., 2016; Kim & Kim, 2020; Kang et al.,  
477 2019). Other methods, such as Label-Distribution-Aware Margin (LDAM) loss (Cao et al., 2019),  
478 Balanced Softmax (Ren et al., 2020), Logit Correction (LC) (Liu et al., 2022), and Unsupervised  
479 Logit Adjustment (uLA) (Tsirigotis et al., 2024), adjust classifier margins to handle class or group  
480 imbalance effectively.

481 **Memorization and Spurious Correlations.** Research has shown that memorization in neural  
482 networks can significantly affect model robustness and generalization. Arpit et al. (2017); Maini et al.  
483 (2022); Stephenson et al. (2021); Maini et al. (2023); Krueger et al. (2017) explore memorization's  
484 impact on neural networks, examining aspects like loss sensitivity, curvature, and the layer where  
485 memorization occurs. Yang et al. (2022) investigate "rare spurious correlations," which are akin  
486 to example-specific noise features that models memorize. Bombari & Mondelli (2024) provide a  
487 theoretical framework quantifying the memorization of spurious features, differentiating between  
488 model stability with respect to individual samples and alignment with spurious patterns. Finally,  
489 Yang et al. (2024) propose Residual-Memorization (ResMem), which combines neural networks with  
490 k-nearest neighbor-based regression to fit residuals, enhancing test performance across benchmarks.

491 **E Proofs**

492 **Lemma E.1.** Let  $p_1(y = j | x) = \frac{e^{\phi_j(x)}}{\sum_{i=1}^k e^{\phi_i(x)}}$  be a softmax over the logits  $\phi_i(x)$  and define  
 493  $p_2(y = j | x)$  such that  $p_2(y = j | x) \propto w(j, x)p_1(y = j | x)$  for some weighting function  $w(j, x)$ .  
 494 Then,

$$p_2(y = j | x) = \frac{e^{\phi_j(x) + \log w(j, x)}}{\sum_{i=1}^k e^{\phi_i(x) + \log w(i, x)}}.$$

495 *Proof.* Given  $p_2(y = j | x) \propto w(j, x)p_1(y = j | x)$ , we substitute the expression for  $p_1(y = j | x)$ :

$$p_2(y = j | x) \propto w(j, x) \cdot \frac{e^{\phi_j(x)}}{\sum_{i=1}^k e^{\phi_i(x)}}.$$

496 Since  $w(j, x) \cdot e^{\phi_j(x)} = e^{\phi_j(x) + \log w(j, x)}$ , we have

$$p_2(y = j | x) \propto \frac{e^{\phi_j(x) + \log w(j, x)}}{\sum_{i=1}^k e^{\phi_i(x) + \log w(i, x)}}.$$

497 To ensure that  $p_2(y = j | x)$  is a valid probability distribution that sums to 1, we need a normalization  
 498 factor. Define the normalization constant  $Z$  as follows:

$$Z = \sum_{j=1}^k e^{\phi_j(x) + \log w(j, x)}.$$

499 Thus, the properly normalized form of  $p_2(y = j | x)$  is:

$$p_2(y = j | x) = \frac{e^{\phi_j(x) + \log w(j, x)}}{Z}.$$

500 Substituting back the expression for  $Z$ , we get

$$p_2(y = j | x) = \frac{e^{\phi_j(x) + \log w(j, x)}}{\sum_{i=1}^k e^{\phi_i(x) + \log w(i, x)}}.$$

501 This completes the proof. □

502 **F Proof of Theorem 2.2**

Setting derivatives of the objective equation ?? zero gives the normal equation

$$\frac{1}{n} \sum_{j=1}^n (s(x_j^\top w) - y_j)x_j + \lambda w = 0.$$

503 Solving for  $w$  then gives

$$\hat{w} = \sum_{i=1}^n \alpha_i x_i, \text{ with } \alpha_i := \frac{\pi_i - \hat{\pi}_i}{\eta}, \text{ with } \pi_i := 1_{\{y_i > 0\}}, \hat{\pi}_i := s(v_i), v_i := x_i^\top \hat{w}, \eta := n\lambda. \quad (9)$$

504 Note that the  $v_i$ 's correspond to logits, while the  $\alpha = (\alpha_1, \dots, \alpha_n) \in \mathbb{R}^n$  should be thought of as the  
 505 dual representation of the weights vector  $\hat{w}$ . Indeed, by construction, one has

$$\hat{w} = X^\top \alpha, \quad (10)$$

506 where  $X \in \mathbb{R}^{n \times d}$  is the design matrix.

507 Our mission is then to derive necessary and sufficient conditions for  $e > 0$ , where

$$e := \gamma \hat{w}_{\text{spure}} - \hat{w}_{\text{core}} = \sum_{i=1}^n (\gamma x_i^{(1)} - x_i^{(2)}) \alpha_i = \sum_{i=1}^n (\gamma^2 a_i - s_i) \alpha_i, \quad (11)$$

508 where  $s_j := 2y_j - 1$ .

509 **F.1 Fixed-Point Equations**

510 Define subsets  $I_{\pm}, S, L \subseteq [n]$  and integers  $m, k \in [n]$  by

$$I_{\pm} := \{i \in [n] \mid y_i = \pm 1\}, \quad S := \{i \in [n] \mid a_i = \gamma y_i\}, \quad L := \{i \in [n] \mid a_i = -\gamma y_i\}, \quad (12)$$

$$m := \mathbb{E}|S| = pn, \quad k := \mathbb{E}|L| = (1-p)n. \quad (13)$$

511 Thus,  $S$  (resp.  $L$ ) corresponds to the sample indices in the majority (resp. the minority) class.

512 One computes the logits as follows

$$v_i = x_i^{\top} \hat{w} = \sum_{j=1}^n \alpha_j x_j^{\top} x_i = \gamma^2 \sum_{j=1}^n \alpha_j z_j a_i + \sum_j \alpha_j s_j s_i + \sum_{j=1}^n \alpha_j \epsilon_j^{\top} \epsilon_i$$

$$= \begin{cases} a + \sum_{j=1}^n \alpha_j \epsilon_j^{\top} \epsilon_i, & \text{if } i \in S \cap I_+, \\ b + \sum_{j=1}^n \alpha_j \epsilon_j^{\top} \epsilon_i, & \text{if } i \in S \cap I_-, \\ c + \sum_{j=1}^n \alpha_j \epsilon_j^{\top} \epsilon_i, & \text{if } i \in L \cap I_+, \\ e + \sum_{j=1}^n \alpha_j \epsilon_j^{\top} \epsilon_i, & \text{if } i \in L \cap I_-, \end{cases} \quad (14)$$

513 where  $a, b, c, e \in \mathbb{R}$  are defined by

$$a := \gamma \hat{w}_{spu} + \hat{w}_{core} = \sum_{j=1}^n (\gamma^2 z_j + s_j) \alpha_j,$$

$$b := -\gamma \hat{w}_{spu} - \hat{w}_{core} = -\sum_{j=1}^n (\gamma^2 z_j + s_j) \alpha_j,$$

$$e := \gamma \hat{w}_{spu} - \hat{w}_{core} = \sum_{j=1}^n (\gamma^2 z_j - s_j) \alpha_j,$$

$$c := \hat{w}_{core} - \gamma \hat{w}_{spu} = \sum_{j=1}^n (-\gamma^2 z_j + s_j) \alpha_j. \quad (15)$$

514 Observe that

$$b = -a, \quad c = -e. \quad (16)$$

515 The following lemma will be crucial to our proof.

516 **Lemma F.1.** *If  $a \geq 0$  and  $e \geq 0$ , then part (B) of Theorem ?? holds. On the other hand, if  $a \geq 0$  and  $e \leq 0$ , then part (C) of Theorem ?? holds.*

518 *Proof.* Indeed, for a random test point  $(x, a, y)$ , we have

$$\begin{aligned} \mathbb{P}(C_{ERM}(x) = C_{spu}(x)) &= \mathbb{P}(x_{spu} \times x^{\top} \hat{w} \geq 0) = \mathbb{P}(\gamma^2 \hat{w}_{spu} + y x_{spu} \hat{w}_{core} + x_{spu} x_{\epsilon}^{\top} \hat{w}_{\epsilon} \geq 0) \\ &= \mathbb{P}(-x_{spu} x_{\epsilon}^{\top} \hat{w}_{\epsilon} \leq \gamma^2 \hat{w}_{spu} + y x_{spu} \hat{w}_{core}) \end{aligned}$$

519 Now, independent of  $y$ , the random variable  $-x_{spu} x_{\epsilon}^{\top} \hat{w}_{\epsilon}$  has distribution  $N(0, \sigma_{\epsilon}^2 \|\hat{w}_{\epsilon}\|^2)$ . Now,  
520 because  $\hat{w} = X^{\top} \alpha$  by construction, the variance can be written as  $\sigma_{\epsilon}^2 \|\hat{w}_{\epsilon}\|^2 = \sigma_{\epsilon}^2 \|X_{\epsilon}^{\top} \alpha\|^2$ , which is  
521 itself chi-squared random variable which concentrates around its mean  $\sigma_{\epsilon}^4 \|\alpha\|^2$ . Furthermore, thanks  
522 to equation 19,  $\|\alpha\|^2 \leq 1/(n\lambda^2)$ , which vanishes in the limit ??. We deduce that

$$\begin{aligned} \mathbb{P}(C_{ERM}(x) = C_{spu}(x)) &\rightarrow \mathbb{P}(\gamma^2 \hat{w}_{spu} + y x_{spu} \hat{w}_{core} \geq 0) \\ &= p \mathbb{1}_{\{\gamma \hat{w}_{spu} + \hat{w}_{core} \geq 0\}} + (1-p) \mathbb{1}_{\{\gamma \hat{w}_{spu} - \hat{w}_{core} \geq 0\}} \\ &= p \mathbb{1}_{\{a \geq 0\}} + (1-p) \mathbb{1}_{\{e \geq 0\}}. \end{aligned}$$

523 Thus, if  $a \geq 0$  and  $e \geq 0$ , we must have  $\mathbb{P}(C_{ERM}(x) = C_{spu}(x)) = p + 1 - p = 1$ , that is, part (B)  
524 of Theorem 2.2 holds.

525 On the other hand, one has

$$\begin{aligned} \mathbb{P}(C_{ERM}(x) = C_{core}(x)) &= \mathbb{P}(x_{core} \times x^{\top} \hat{w} \geq 0) = \mathbb{P}(\hat{w}_{core} + y x_{spu} \hat{w}_{spu} \geq 0) \\ &= p \mathbb{1}_{\{\hat{w}_{core} + \gamma \hat{w}_{spu} \geq 0\}} + (1-p) \mathbb{1}_{\{\hat{w}_{core} - \gamma \hat{w}_{spu} \geq 0\}} \\ &= q \mathbb{1}_{\{a \geq 0\}} + (1-q) \mathbb{1}_{\{e \leq 0\}}, \end{aligned}$$

526 where  $q := \mathbb{P}(a = y)$ . We deduce that if  $a \geq 0$  and  $e \leq 0$ , then  $\mathbb{P}(C_{ERM}(x) = C_{core}(x)) =$   
527  $q + 1 - q = 1$ , i.e part (C) of Theorem ?? holds.  $\square$



528 **F.2 Structure of the Dual Weights**

529 The following result shows that the dual weights  $\alpha_1, \dots, \alpha_n$  cluster into 4 lumps corresponding to  
 530 the following 4 sets of indices  $S \cap I_+$ ,  $S \cap I_-$ ,  $L \cap I_+$ , and  $L \cap I_-$ .

531 **Lemma F.2.** *There exist positive constants  $A, B, C, E > 0$  such that the following holds with large  
 532 probability uniformly over all indices  $i \in [n]$*

$$\alpha_i \simeq \begin{cases} A, & \text{if } i \in S \cap I_+, \\ -B, & \text{if } i \in S \cap I_-, \\ C, & \text{if } i \in L \cap I_+, \\ -E, & \text{if } i \in L \cap I_-. \end{cases} \quad (17)$$

533 Furthermore, the empirical probabilities predicted by ERM are given by

$$\hat{\pi}_i = y_i - \eta \alpha_i = \begin{cases} 1 - \eta A, & \text{if } i \in S \cap I_+, \\ \eta B, & \text{if } i \in S \cap I_-, \\ 1 - \eta C, & \text{if } i \in L \cap I_+, \\ \eta E, & \text{if } i \in L \cap I_-. \end{cases} \quad (18)$$

534 *Proof.* First observe that

$$\|\alpha\| \leq \frac{1}{\lambda \sqrt{n}}. \quad (19)$$

535 Indeed, one computes

$$\|\alpha\|^2 = \frac{1}{\eta^2} \sum_{i=1}^n (\pi_i - \hat{\pi}_i)^2 \leq \frac{1}{\eta^2} \sum_{i=1}^n 1 \leq \frac{n}{\eta^2} = \frac{1}{\lambda^2 n}$$

Next, observe that  $\sum_j \alpha_j \epsilon_j^\top \epsilon_i = \alpha_i \|\epsilon_i\|^2 + \sum_{j \neq i} \alpha_j \epsilon_j^\top \epsilon_i \simeq \sigma_\epsilon^2 \alpha_i d$ . This is because  $\alpha_i \|\epsilon_i\|^2$  concentrates around its mean which equals  $\sigma_\epsilon^2 \alpha_i d$ , while w.h.p,

$$\frac{1}{\sigma_\epsilon^2 d} \sup_{i \in [n]} \left| \sum_{j \neq i} \alpha_j \epsilon_j^\top \epsilon_i \right| \lesssim \|\alpha\| \sqrt{\frac{n \log n}{d}} = \sigma_\epsilon \|\alpha\| \sqrt{n} \cdot \sqrt{\frac{\log n}{d}} \leq \sigma_\epsilon \lambda \sqrt{\frac{\log n}{d}} = o(1).$$

536 The above is because  $\lambda \rightarrow 0$  and  $(\log n)/d \rightarrow 0$  by assumption. Henceforth we simply ignore the  
 537 contributions of the terms  $\sum_{j \neq i} \alpha_j \epsilon_j^\top \epsilon_i$ . We get the following equations in the limit equation ??

$$v_i = \begin{cases} \sigma_\epsilon^2 \alpha_i d + a, & \text{if } i \in S \cap I_+, \\ \sigma_\epsilon^2 \alpha_i d + b, & \text{if } i \in S \cap I_-, \\ \sigma_\epsilon^2 \alpha_i d + c, & \text{if } i \in L \cap I_+, \\ \sigma_\epsilon^2 \alpha_i d + e, & \text{if } i \in L \cap I_-, \end{cases} \quad (20)$$

$$\eta \alpha_i = y_i - s(v_i) = \begin{cases} 1 - s(\sigma_\epsilon^2 \alpha_i d + a), & \text{if } i \in S \cap I_+, \\ -s(\sigma_\epsilon^2 \alpha_i d + b), & \text{if } i \in S \cap I_-, \\ 1 - s(\sigma_\epsilon^2 \alpha_i d + c), & \text{if } i \in L \cap I_+, \\ -s(\sigma_\epsilon^2 \alpha_i d + e), & \text{if } i \in L \cap I_-. \end{cases}$$

538 Now, because of monotonicity of  $s$ , we can find  $A, B, C, E > 0$  such that

$$\alpha_i = \begin{cases} A, & \text{if } i \in S \cap I_+, \\ -B, & \text{if } i \in S \cap I_-, \\ C, & \text{if } i \in L \cap I_+, \\ -E, & \text{if } i \in L \cap I_-, \end{cases}$$

539 as claimed. □

540 We will make use of the following lemma.

541 **Lemma F.3.** *In the unregularized limit  $\lambda \rightarrow 0^+$ , it holds that  $\eta A, \eta B, \eta C, \eta E \in [0, 1/2]$ .*

542 *Proof.* Indeed, in that unregularized limit, ERM attains zero classification error on the training dataset  
 543 (first part of Theorem 2.2). This means mean that  $\hat{\pi}_i \geq 1/2$  iff  $y_i = 1$ , and the result follows. □

544 **F.3 Final Touch (Proof of Theorem ??)**

545 We resume the proof of Theorem 2.2. The scalars  $A, B, C, E$  must verify

$$\begin{aligned}
\eta A &= 1 - s(\sigma_\epsilon^2 Ad + a) = s(-\sigma_\epsilon^2 Ad - a), \\
\eta B &= s(-\sigma_\epsilon^2 Bd + b) = s(-\sigma_\epsilon^2 Bd - a) = 1 - s(\sigma_\epsilon^2 Bd + a), \\
\eta E &= s(-\sigma_\epsilon^2 Ed + e), \\
\eta C &= 1 - s(\sigma_\epsilon^2 Cd + c) = 1 - s(\sigma_\epsilon^2 Cd - e) = s(-\sigma_\epsilon^2 Cd + e).
\end{aligned} \tag{21}$$

546 We deduce that

$$A = B, \quad C = E, \tag{22}$$

$$\eta A = s(-\sigma_\epsilon^2 Ad - a), \quad \eta E = s(-\sigma_\epsilon^2 Ed + e). \tag{23}$$

547 **Proof of Part (C).** In particular, for the noiseless case where  $\sigma_\epsilon \rightarrow 0^+$ , we have  $\eta A \simeq s(-a)$  and  
548  $\eta E \simeq s(e)$ . We know from Lemma F.3 that  $\eta A, \eta E \leq 1/2$ . This implies  $a \geq 0$  and  $e \leq 0$ , and  
549 thanks to Lemma F.1, we deduce part (C) of Theorem 2.2.

550 **Proof of Part (B).** It remains to show that  $a \geq 0$  and  $e \geq 0$  in the noisy regime  $\sigma_\epsilon > 0$ , and then  
551 conclude via Lemma F.1.

552 Define  $N_1 := |S \cap I_+|$ ,  $N_2 := |S \cap I_-|$ ,  $N_3 := |L \cap I_+|$ ,  $N_4 := |L \cap I_-|$ . Note that from the  
553 definition of  $a, b, c, e$  in equation 15, one has

$$\begin{aligned}
a &= (\gamma^2 + 1)(N_1 + N_2)A - (\gamma^2 - 1)(N_3 + N_4)E, \\
e &= (\gamma^2 - 1)(N_1 + N_2)A - (\gamma^2 + 1)(N_3 + N_4)E, \\
b &= -a, \quad c = -e, \\
\eta A &= s(-\sigma_\epsilon^2 Ad - a), \quad \eta E = s(-\sigma_\epsilon^2 Ed + e), \\
B &= A, \quad C = E.
\end{aligned} \tag{24}$$

554 We now show that  $a \geq 0$  and  $e \geq 0$  under the conditions  $d \gg \log n$  and  $\gamma \gg \sigma_\epsilon \sqrt{d/m}$ .

555 Indeed, under the second condition, the following holds w.h.p

$$\begin{aligned}
\sigma_\epsilon^2 d + (\gamma^2 + 1)(N_1 + N_2) &= ((\gamma^2 + 1)(N_1 + N_2) + \sigma_\epsilon^2 d) \simeq ((\gamma^2 + 1)m + \sigma_\epsilon^2 d) \\
&\simeq (\gamma^2 + 1)m \simeq (\gamma^2 + 1)(N_1 + N_2),
\end{aligned}$$

556 where we have used the fact that  $N_1 + N_2$  concentrates around its mean  $m = pn$ . We deduce that

$$\begin{aligned}
\sigma_\epsilon^2 Ad + a &= (\sigma_\epsilon^2 d + (\gamma^2 + 1)(N_1 + N_2))A - (\gamma^2 - 1)(N_3 + N_4)E \\
&\simeq (\gamma^2 + 1)(N_1 + N_2)A - (\gamma^2 - 1)(N_3 + N_4)E \\
&\simeq a,
\end{aligned}$$

557 from which we get.

$$1/2 \geq \eta A \geq s(-\sigma_\epsilon^2 Ad - a) = s(-(1 + o(1))a) = s(-a) + o(1),$$

558 i.e  $s(-a) \geq 1/2 - o(1)$ . But this can only happen if  $a \geq 0$ .

559 Finally, the conditions  $d \gg \log n$  and  $\gamma \gg \sigma_\epsilon \sqrt{d/m}$  imply  $\gamma \gg K \sigma_\epsilon \sqrt{d/k}$  and  $g \geq K \log(3n)$  for  
560 any constant  $K > 0$ . Theorem 1 of Puli et al. (2023) <https://arxiv.org/abs/2308.12553>  
561 then gives  $e = \gamma \hat{w}_{spu} - \hat{w}_{core} > 0$ , and we are done.  $\square$