# An investigation of why overparameterization exacerbates spurious correlations

**Shiori Sagawa** [* 1]  **Aditi Raghunathan** [* 1]  **Pang Wei Koh** [* 1]  **Percy Liang** [1]
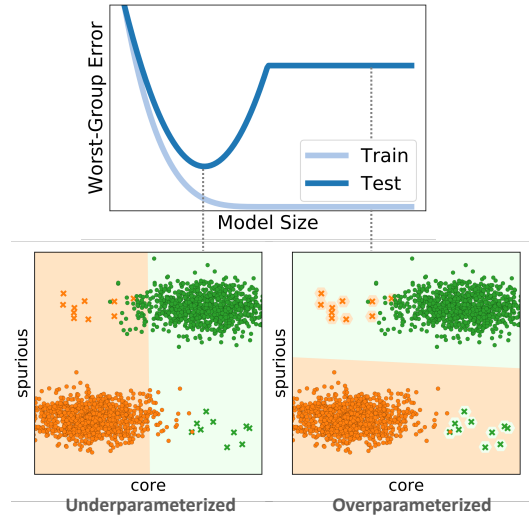
## Abstract

We study why overparameterization—increasing model size well beyond the point of zero training error—can hurt test error on minority groups despite improving average test error when there are spurious correlations in the data. Through simulations and experiments on two image datasets, we identify two key properties of the training data that drive this behavior: the proportions of majority versus minority groups, and the signal-to-noise ratio of the spurious correlations. We then analyze a linear setting and theoretically show how the inductive bias of models towards "memorizing" fewer examples can cause overparameterization to hurt. Our analysis leads to a counterintuitive approach of subsampling the majority group, which empirically achieves low minority error in the overparameterized regime, even though the standard approach of upweighting the minority fails. Overall, our results suggest a tension between using overparameterized models versus using all the training data for achieving low worst-group error.

## 1. Introduction

The typical goal in machine learning is to minimize the average error on a test set that is independent and identically distributed (i.i.d.) to the training set. A large body of prior work has shown that overparameterization—increasing model size beyond the point of zero training error—improves average test error in a variety of settings, both empirically (with neural networks, e.g., Nakkiran et al. (2019)) and theoretically (with linear and random projection models, e.g., Belkin et al. (2019); Mei & Montanari (2019)).

However, recent work has also demonstrated that models with low average error can still fail on particular groups of



*Figure 1.* **Top**: Overparameterization *hurts* test error on the worst group when models are trained with the reweighted objective that upweights minority groups (Equation 3). Without reweighting, models have poor worst-group error regardless of model size (Appendix A.1). **Bottom**: Consider data points $(x, y)$, where $x \in \mathbb{R}^2$ comprises a core feature $x_{\text{core}}$ (x-axis) and a spurious feature $x_{\text{spu}}$ (y-axis). The label $y$ is highly correlated with $x_{\text{spu}}$, except on two minority groups (crosses). Underparameterized models use the core feature (left), but overparameterized models use the spurious feature and memorize the minority points (right).

data points (Blodgett et al., 2016; Hashimoto et al., 2018; Buolamwini & Gebru, 2018). This problem of high worst-group error arises especially in the presence of spurious correlations, such as strong associations between label and background in image classification (McCoy et al., 2019; Sagawa et al., 2020). To mitigate this problem, common approaches reduce the worst-group training loss, e.g., through distributionally robust optimization (DRO) or simply upweighting the minority groups. Sagawa et al. (2020) showed these approaches improve worst-group error on strongly regularized neural networks but fail to help standard neural networks that can achieve zero training error, suggesting that increasing model capacity by reducing regularization—and perhaps by increasing overparameterization as well—can exacerbate spurious correlations.

In this paper, we investigate why overparameterization exacerbates spurious correlations under the above approach of upweighting minority groups. We first confirm on two
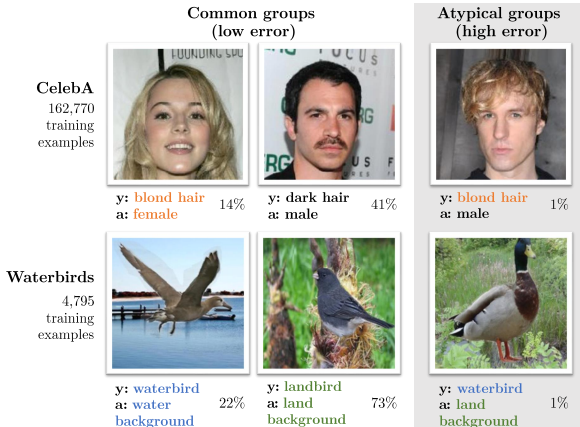
[*]Equal contribution [1]Stanford University. Correspondence to: Shiori Sagawa <ssagawa@cs.stanford.edu>, Aditi Raghunathan <aditir@stanford.edu>, Pang Wei Koh <pangwei@cs.stanford.edu>.

*Figure 2.* We consider two image datasets, CelebA and Waterbirds, where the label $y$ is correlated with a spurious attribute $a$ in a majority of the training data. The % beside each group shows its frequency in the training data. To measure how robust a model is to the spurious attribute, we divide the data into groups based on $(y, a)$ and record the highest error incurred by a group. Figure adapted from Sagawa et al. (2020).

image datasets (Figure 2) that directly increasing overparameterization (i.e., increasing model size) indeed hurts worst-group error, leading to models that are highly inaccurate on the minority groups where the spurious correlation does not hold (Section 3). In contrast, their underparameterized counterparts obtain much better worst-group error, but do worse on average. We also confirm that models trained via empirical risk minimization (i.e., without upweighting the minority) have poor worst-group test error regardless of whether they are under- or overparameterized. Through simulations on a synthetic setting, we further identify two properties of the training data that modulate the effect of overparameterization: (i) the relative sizes of the majority versus minority groups, and (ii) how informative the spurious features are relative to the core features (Section 4).

Why does overparameterization exacerbate spurious correlations? Underparameterized models do not rely on spurious features because that would incur high training error on the (upweighted) minority groups where the spurious correlation does not hold. In contrast, overparameterized models can always obtain zero training error by memorizing training examples, and instead rely on their inductive bias to pick a solution—which features to use and which examples to memorize—out of all solutions with zero training error. Our results suggest an intuitive story of why overparameterization can hurt: because overparameterized models can have an inductive bias towards "memorizing" fewer examples (Figure 1). If (i) the majority groups are sufficiently large and (ii) the spurious features are more informative than the core features for these groups, then overparameterized models could choose to use the spurious features because it entails less memorization, and therefore suffer high worst-

group test error. We test this intuition through simulations and formalize it in a theoretical analysis (Section 5).

Our analysis also leads to the counterintuitive result that on overparameterized models, subsampling the majority groups is much more effective at improving worst-group error than upweighting the minority groups. Indeed, an overparameterized model trained on a subset of $<5\%$ of the data performs similarly (on average and on the worst group) to an underparameterized model trained on all the data (Section 6). This suggests a possible tension between using overparameterized models and using all the data; average error benefits from both, but improving worst-group error seems to rely on using only one but not both.

## 2. Setup

**Spurious correlation setup.** We adopt the setting studied in Sagawa et al. (2020), where each example comprises the input features $x$, a label (core attribute) $y \in \mathcal{Y}$, and a spurious attribute $a \in \mathcal{A}$. Each example belongs to a group $g \in \mathcal{G} = \mathcal{Y} \times \mathcal{A}$, where $g = (y, a)$. Importantly, the spurious attribute $a$ is correlated with the label $y$ in the training set. We focus on the binary setting in which $\mathcal{Y} = \{1, -1\}$ and $\mathcal{A} = \{1, -1\}$.

**Applications.** We study two image classification tasks (Figure 2). In the first task, the label is spuriously correlated with demographics: specifically, we use the CelebA dataset (Liu et al., 2015) to classify hair color between the labels $\mathcal{Y} = \{\text{blonde, non-blonde}\}$, which are correlated with the gender $\mathcal{A} = \{\text{female, male}\}$. In the second task, the label is spuriously correlated with image background. We use the Waterbirds dataset (based on datasets from Wah et al. (2011); Zhou et al. (2017) and modified by Sagawa et al. (2020)) to classify between the labels $\mathcal{Y} = \{\text{waterbird, landbird}\}$, which are spuriously correlated with the image background $\mathcal{A} = \{\text{water background, land background}\}$. See Appendix A.5 for more dataset details.

**Objectives and metrics.** We evaluate a model $w$ by its *worst-group* error,

$$\text{Err}_{\text{wg}}(w) := \max_{g \in \mathcal{G}} \mathbb{E}_{x,y|g} \left[ \ell_{0-1}(w; (x,y)) \right], \quad (1)$$

where $\ell_{0-1}$ is the 0-1 loss. In other words, we measure the error (% of examples that are incorrectly labeled) in each group, and then record the highest error across all groups. The standard approach to training models is empirical risk minimization (ERM): given a loss function $\ell$, find the model $w$ that minimizes the average training loss

$$\hat{\mathcal{R}}_{\text{ERM}}(w) = \hat{\mathbb{E}}_{(x,y,g)} \left[ \ell(w; (x,y)) \right]. \quad (2)$$

However, in line with Sagawa et al. (2020), we find that models trained via ERM have poor worst-group test error

regardless of whether they are under- or overparameterized (Appendix A.1). To achieve low worst-group test error, prior work proposed modified objectives that focus on the worst-group loss, such as group distributionally robust optimization (group DRO) which directly optimizes for the worst-group training loss (Hu et al., 2018; Sagawa et al., 2020) or reweighting (Shimodaira, 2000; Byrd & Lipton, 2019). Sagawa et al. (2020) showed that both approaches can help worst-group loss, though group DRO is typically more effective. For simplicity, we focus on the well-studied reweighting approach, which optimizes

$$\hat{\mathcal{R}}_{\text{reweight}}(w) = \hat{\mathbb{E}}_{(x,y,g)} \left[ \frac{1}{\hat{p}_g} \ell(w; (x,y)) \right], \qquad (3)$$

where $\hat{p}_g$ is the fraction of training examples in group $g$. The intuition behind reweighting is that it makes each group contribute the same weight to the training objective: that is, minority groups are upweighted, while majority groups are downweighted. Note that this approach requires the groups $g$ to be specified at training time, though not at test time.

## 3. Overparameterization hurts worst-group error

Sagawa et al. (2020) observed that decreasing $L_2$ regularization hurts worst-group error. Though increasing overparameterization and reducing regularization can have different effects (Zhang et al., 2017; Mei & Montanari, 2019), this suggests that overparameterization might similarly exacerbate spurious correlations. Here, we show that directly increasing overparameterization (model size) indeed hurts worst-group error even though it improves average error.

**Models.** We study the CelebA and Waterbirds datasets described above. For CelebA, we train a ResNet10 model (He et al., 2016), varying model size by increasing the network width from 1 to 96, as in Nakkiran et al. (2019). For Waterbirds, we use logistic regression over random projections, as in Mei & Montanari (2019). Specifically, let $x \in \mathbb{R}^d$ denote the input features, which we obtain by passing the input image through a pre-trained, fixed ResNet-18 model. We train an unregularized logistic regression model over the feature representation $\text{ReLU}(Wx) \in \mathbb{R}^m$, where $W \in \mathbb{R}^{m \times d}$ is a random matrix with each row sampled uniformly from the unit sphere $\mathbb{S}^{d-1}$. We vary model size by increasing the number of projections $m$ from 1 to 10,000. We train each model by minimizing the reweighted objective (Equation (3)). For more details, see Appendix A.5.

**Results.** Overparameterization improves average test error across both datasets, in line with prior work (Belkin et al., 2019; Nakkiran et al., 2019) (Figure 3). However, in stark contrast, overparameterization *hurts* worst-group error: the best worst-group test error is achieved by an *underparameterized* model with non-zero training error. On CelebA,

the smallest model (width 1) has 12.4% worst-group training error but comparatively low worst-group test error of 25.6%. As width increases, training error goes to zero but worst-group test error gets worse, reaching >60% for overparameterized models with zero training error. Similarly, on Waterbirds, an underparameterized model with 90 random features and worst-group training error of 17.7% obtains the best worst-group test error of 26.6%, while overparameterized models with zero training error yield worst-group test error of 42.4% at best.

In Appendix A.2, we also confirm that stronger regularization improves worst-group error but hurts average error in overparameterized models, while it has little effect on both worst-group and average error in underparameterized models. However, we focus on understanding the effect of overparameterization in the remainder of the paper.

**Discussion.** Why does overparameterization hurt worst-group test error? We make two observations. First, in the overparameterized regime, the smallest groups incur the highest test error (blonde males in CelebA and waterbirds on land background in Waterbirds), despite having zero training error. In other words, overparameterized models perfectly fit the minority points at training time, but seem to do so by using patterns that do not generalize. We informally refer to this behavior as "memorizing" the minority points.

Second, underparameterized models do obtain low worst-group error by learning patterns that generalize to both majority and minority groups. Therefore, overparameterized models should also be able to learn these patterns while attaining zero training error (e.g., by memorizing the training points that the underparameterized model cannot fit). Despite this, overparameterized models seem to learn patterns that generalize well on the majority but do not work on the minority (such as the spurious attributes $a$ in Figure 2).

What makes overparameterized models memorize the minority instead of learning patterns that generalize well on both majority and minority groups? We study this question in the next two sections: in Section 4, we use simulations to understand properties of the data distribution that give rise to this trend, and in Section 5 we analyze a simplified linear setting and show how the inductive bias of models towards memorizing fewer points can lead to overparameterized models choosing to use spurious correlations.

## 4. Simulation studies

The discussion in Section 3 suggests two properties of the training distribution that modulate the effect of overparameterization on worst-group error. Intuitively, overparameterized models should be more incentivized to use the spurious features and memorize the minority groups if (i) the proportion of the majority group, $p_{\text{maj}}$, is higher, and (ii) the ratio
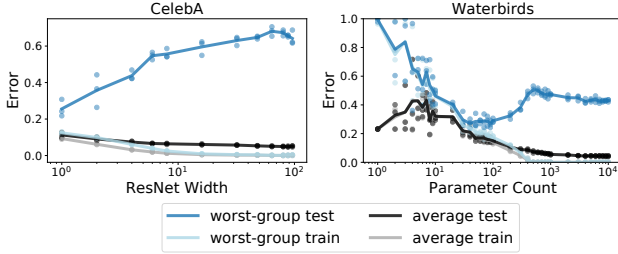
*Figure 3.* Increasing overparameterization (i.e., increasing model size) hurts the worst-group test error even though it improves the average test error. Here, we show results for models trained on the reweighted objective for CelebA (left) and Waterbirds (right).

of how informative the spurious features are relative to the core features, $r_{\text{s:c}}$, is higher. In this section, we use simulations to confirm these intuitions and probe how $p_{\text{maj}}$ and $r_{\text{s:c}}$ affect worst-group error in overparameterized models.

## 4.1. Synthetic experiment setup

**Data distribution.** We construct a synthetic dataset that replicates the empirical trends in Section 3. As in Section 2, the label $y \in \{1, -1\}$ is spuriously correlated with a spurious attribute $a \in \{1, -1\}$. We divide our training data into four groups accordingly: two majority groups with $a = y$, each of size $n_{\text{maj}}/2$, and two minority groups with $a = -y$, each of size $n_{\text{min}}/2$. We define $n = n_{\text{maj}} + n_{\text{min}}$ as the total number of training points, and $p_{\text{maj}} = n_{\text{maj}}/n$ as the fraction of majority examples. The higher $p_{\text{maj}}$ is, the more strongly $a$ is correlated with $y$ in the training data.

Each $(y, a)$ group has its own distribution over input features $x = [x_{\text{core}}, x_{\text{spu}}] \in \mathbb{R}^{2d}$ comprising core features $x_{\text{core}} \in \mathbb{R}^d$ generated from the label/core attribute $y$, and spurious features $x_{\text{spu}} \in \mathbb{R}^d$ generated from the spurious attribute $a$:

$$x_{\text{core}} \mid y \sim \mathcal{N}(y\mathbf{1}, \sigma_{\text{core}}^2 I_d)$$
$$x_{\text{spu}} \mid a \sim \mathcal{N}(a\mathbf{1}, \sigma_{\text{spu}}^2 I_d). \quad (4)$$

The core and spurious features are both noisy and encode their respective attributes at different signal-to-noise ratios. We define the *spurious-core information ratio* (SCR) as $r_{\text{s:c}} = \sigma_{\text{core}}^2 / \sigma_{\text{spu}}^2$. The higher the SCR, the more signal there is about the spurious attribute in the spurious features, relative to the signal about the label in the core features.

Compared to the image datasets we studied in Section 3, this synthetic dataset offers two key simplifications. First, the only differences between groups stem from their differences in $(y, a)$, which isolates the effect of flipping the spurious attribute $a$. In contrast, in real datasets, groups can differ in other ways, e.g., more label noise in one group. Second, the relative difficulty of estimating $y$ versus $a$ is completely governed by changing $\sigma_{\text{core}}^2$ and $\sigma_{\text{spu}}^2$. In contrast, real datasets have additional complications, e.g., estimating $y$ might involve a more complex function of the input $x$ than
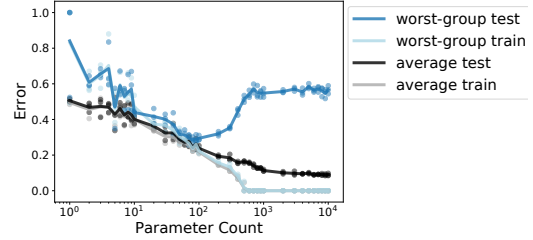


*Figure 4.* Overparameterization hurts worst-group test error but improves average test error on synthetic data, reproducing the trends we observe in real data.
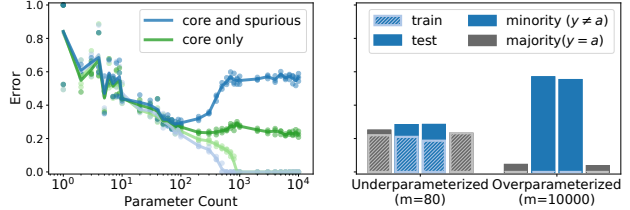


*Figure 5.* Overparameterized models have poor worst-group performance on the synthetic data because they rely on spurious features. **Left**: removing the spurious feature (green) eliminates the detrimental effect of overparameterization. **Right**: overparamerized models do well on the majority groups where the spurious features match the label, but poorly on the minority groups.

estimating $a$, and there might be an inductive bias towards learning a simpler model over a more complex one.

In all of the experiments below, we fix the total number of training points $n$ to 3000, and set $d = 100$ (so each input $x$ has $2d = 200$ dimensions). Unless otherwise specified, we set the majority fraction $p_{\text{maj}} = 0.9$ and the noise levels $\sigma_{\text{spu}}^2 = 1$ and $\sigma_{\text{core}}^2 = 100$ to encourage the model to use the spurious features over the core features.

**Model.** To avoid the complexities of optimizing neural networks, we follow the same random features setup we used for Waterbirds in Section 3: unregularized logistic regression using the reweighted objective on the random feature representation $\text{ReLU}(Wx) \in R^m$, where $W \in \mathbb{R}^{m \times d}$ is a random matrix (Mei & Montanari, 2019).

## 4.2. Observations on synthetic dataset

**The synthetic dataset replicates the trends we observe on real datasets.** Figure 4 shows how average and worst-group error change with the number of parameters/random projections $m$. This matches the trends we obtained on CelebA and Waterbirds in Section 3. The best worst-group test error of 28.5% is achieved by an underparameterized model, whereas highly overparameterized models achieve high worst-group test error that plateaus at around 55%. In contrast, the average test error is better for overparameterized models than for underparameterized models.

**Overparameterized models use spurious features.** Fig-

ure 5-Right shows that overparameterized models have high test error on minority groups ($a = -y$) despite zero training error, but perform very well on the majority groups ($a = y$). Since the only difference between the minority and majority groups in the synthetic dataset is the relative signs of the core and spurious attributes, this suggests overparameterized models are using spurious features and simply memorizing the minority groups to get zero training error, consistent with our discussion in Section 3. In contrast, the underparameterized model has low training and test errors across all groups, suggesting that it relies mainly on core features.

These results imply that the degradation in the worst-group test error is due to the spurious features. We confirm that overparameterization no longer hurts when we "remove" the spurious features by replacing them with noise centered around zero (i.e., we replace the mean of $x_{\mathsf{spu}}$ by 0). In this case, the best worst-group test error is now obtained by an overparameterized model, as shown in Figure 5-Left.

### 4.3. Distributional properties

What properties of the training data make overparameterization hurt worst-group error? We study (i) $p_{\mathsf{maj}}$, which controls the relative size of majority to minority groups, and (ii) $r_{\mathsf{s:c}}$, the relative informativeness of spurious to core features. In the synthetic dataset, overparameterization hurts worst-group test error only when both are sufficiently high. In contrast, overparameterization helps average test error regardless; see Appendix A.3.

**Effect of the majority fraction $p_{\mathsf{maj}}$.** We observe that increasing $p_{\mathsf{maj}} = n_{\mathsf{maj}}/n$, which controls the relative size of the majority versus minority groups, makes overparameterization hurt worst-group error more (Figure 6). When the groups are perfectly balanced with $p_{\mathsf{maj}} = 0.5$, overparameterization no longer hurts the worst-group test error, with overparameterized models achieving better worst-group test error than all underparameterized models. This suggests that group imbalance can be a key factor inducing the detrimental effect of overparameterization.

**Effect of the spurious-core information ratio $r_{\mathsf{s:c}}$.** Next, we characterize the effect of $r_{\mathsf{s:c}} = \sigma_{\mathsf{core}}^2/\sigma_{\mathsf{spu}}^2$, which measures the relative informativeness of the spurious versus core features. A high $r_{\mathsf{s:c}}$ means that the spurious features are more informative. We vary $r_{\mathsf{s:c}}$ by changing $\sigma_{\mathsf{spu}}^2$ while keeping $\sigma_{\mathsf{core}}^2 = 100$ fixed, since this does not change the best possible worst-group test error (with a model that uses only the core features $x_{\mathsf{core}}$). Figure 6 shows that the higher $r_{\mathsf{s:c}}$ is, the more overparameterization hurts. As $r_{\mathsf{s:c}}$ increases, the spurious features become more informative, and overparameterized models rely more on them than the core features; underparameterized models outperform overparameterized models only for sufficiently large $r_{\mathsf{s:c}} \geq 1$. Note that increasing $r_{\mathsf{s:c}}$ does not significantly affect the worst-group
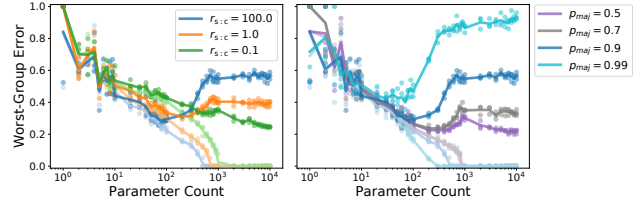


*Figure 6.* The higher the majority fraction $p_{\mathsf{maj}}$ and the spurious-core information ratio $r_{\mathsf{s:c}}$, the more overparameterization hurts the worst-group test error. With sufficiently low $p_{\mathsf{maj}}$ and $r_{\mathsf{s:c}}$, overparameterization switches to helping worst-group test error.

test error in the underparameterized regime, since the core features $x_{\mathsf{core}}$ are unaffected. In contrast, increasing the majority fraction $p_{\mathsf{maj}}$ hurts the worst-group test error in both underparameterized and overparameterized models.

### 4.4. An intuitive story

We return to the question of what makes overparameterized models memorize the minority instead of learning patterns that generalize on both majority and minority groups. The simulation results above show that of all overparameterized models that achieve zero training error, the inductive bias of the model class and training algorithm favors models that use spurious features which generalize only for the majority groups, instead of learning to use core features that also generalize well on the minority groups.

What is the nature of this inductive bias? Consider a model that predicts the label $y$ by returning its estimate of the spurious attribute $a$ from $x_{\mathsf{spu}}$, taking advantage of the fact that $y$ and $a$ are correlated in the training data. To get achieve zero training error, it will need to memorize the points in the minority group, e.g., by exploiting variations due to noise in the features $x$. On the other hand, consider a model that predicts $y$ by returning a direct estimate of $y$ based on the core features $x_{\mathsf{core}}$. Because $x_{\mathsf{core}}$ provides a noisier estimate of $y$ than $x_{\mathsf{spu}}$ does for $a$, this model will need to memorize all points for which $x_{\mathsf{core}}$ gives an inaccurate prediction of $y$ due to noise. Since the estimators of the core and spurious attributes are equally easy to learn, the main difference between these two models is the number of examples to be memorized.

We therefore hypothesize that *the inductive bias favors memorizing as few points as possible*. This is consistent with the results above: the model uses $x_{\mathsf{spu}}$ and memorizes the minority points only when the fraction of minority points is small (high majority fraction $p_{\mathsf{maj}}$). Similarly, the model uses $x_{\mathsf{spu}}$ over $x_{\mathsf{core}}$ to fit the majority points only when the spurious features are less noisy (high $r_{\mathsf{s:c}}$) and therefore require less memorization to obtain zero training error than the core features. In the next section, we make this intuition formal by analyzing a related but simpler linear setting.

## 5. Theoretical analysis

In this section, we show how the inductive bias against memorization leads to overparameterization exacerbating spurious correlations. Our analysis explicates the effect of the inductive bias and the importance of the data parameters $p_{\mathsf{maj}}$ and $r_{\mathsf{s:c}}$ discussed in Section 4.

The synthetic setting discussed in Section 4 is difficult to analyze because of the non-linear random projections, so we introduce a linear *explicit-memorization* setting that allows us to precisely define the concept of memorization. For clarity, we refer to the previous synthetic setting in Section 4 as the *implicit-memorization* setting. In Appendix A.4, we show empirically that models in these two settings behave similarly in the overparameterized regime, though they differ in the underparameterized regime.

In the previous implicit-memorization setting, we varied model size and memorization capacity by varying the number of random projections of the input. In the new explicit-memorization setting, we instead use linear models that act directly on the input and introduce explicit "noise features" that can be used to memorize. We vary the memorization capacity by varying the number of explicit noise features.

### 5.1. Explicit-memorization setup

**Training data.** We consider input features $x = [x_{\mathsf{core}}, x_{\mathsf{spu}}, x_{\mathsf{noise}}]$, where the core feature $x_{\mathsf{core}} \in \mathbb{R}$ and the spurious feature $x_{\mathsf{spu}} \in \mathbb{R}$ are scalars. As in the implicit-memorization setup, they are generated based on the label and the spurious attribute, respectively:

$$x_{\mathsf{core}} \mid y \sim \mathcal{N}(y, \sigma_{\mathsf{core}}^2), \quad x_{\mathsf{spu}} \mid a \sim \mathcal{N}(a, \sigma_{\mathsf{spu}}^2).$$

The "noise" features $x_{\mathsf{noise}} \in \mathbb{R}^N$ are generated as

$$x_{\mathsf{noise}} \sim \mathcal{N}\left(0, \frac{\sigma_{\mathsf{noise}}^2}{N} I_N\right),$$

where $\sigma_{\mathsf{noise}}^2$ is a constant. The scaling by $1/N$ ensures that for large $N$, the norm of the noise vectors $\|x_{\mathsf{noise}}\|_2^2 \approx \sigma_{\mathsf{noise}}^2$ is approximately constant with high probability. Intuitively, when $N$ is large, overparameterized models can use $x_{\mathsf{noise}}$ to fit a training point $x$ without affecting its predictions on other points, thereby memorizing $x$. We formalize this notion of memorization later in Section 5.2.

As before, the training data is composed of four groups, each corresponding to a combination of the label $y \in \{-1, 1\}$ and the spurious attribute $a \in \{-1, 1\}$: two majority groups with $a = y$, each of size $n_{\mathsf{maj}}/2$, and two minority groups with $a = -y$, each of size $n_{\mathsf{min}}/2$. Combined, there are $n$ training examples $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$.

**Model.** We study unregularized logistic regression on the input features $x \in \mathbb{R}^{N+2}$. As before, we consider the reweighted estimator $\hat{w}^{\mathsf{rw}}$. When the training data is linearly separable, the minimizer of the unregularized logistic loss on the training data is not well-defined. We therefore define $\hat{w}^{\mathsf{rw}}$ in terms of the sequence of $L_2$-regularized models $\hat{w}_\lambda^{\mathsf{rw}}$:

$$\hat{w}_\lambda^{\mathsf{rw}} \overset{\text{def}}{=} \underset{w \in \mathbb{R}^{N+2}}{\arg\min} \ \hat{\mathbb{E}}_{(x,y,g)}\left[\frac{1}{\hat{p}_g}\ell(w; (x, y))\right] + \frac{\lambda}{2}\|w\|_2^2,$$

where $\ell$ is the logistic loss and $\hat{p}_g$ is the fraction of training examples in group $g$. Since scaling a model does not affect its 0-1 error, we define $\hat{w}^{\mathsf{rw}}$ as the limit of this sequence, scaled to unit norm, as the regularization strength $\lambda \to 0^+$:

$$\hat{w}^{\mathsf{rw}} \overset{\text{def}}{=} \lim_{\lambda \to 0^+} \frac{\hat{w}_\lambda^{\mathsf{rw}}}{\|\hat{w}_\lambda^{\mathsf{rw}}\|_2}. \tag{5}$$

In the underparameterized regime, the training data is not linearly separable and we simply have $\hat{w}^{\mathsf{rw}} = \hat{w}_0^{\mathsf{rw}}/\|\hat{w}_0^{\mathsf{rw}}\|_2$. In the overparameterized regime where $N \gg n$, the training data is linearly separable, and Rosset et al. (2004) showed that $\hat{w}^{\mathsf{rw}} = \hat{w}^{\mathsf{mm}}$, where $\hat{w}^{\mathsf{mm}}$ is the max-margin classifier

$$\hat{w}^{\mathsf{mm}} \overset{\text{def}}{=} \underset{\|w\|_2=1}{\arg\max} \ \min_i y^{(i)}(w \cdot x^{(i)}). \tag{6}$$

The equivalence $\hat{w}^{\mathsf{rw}} = \hat{w}^{\mathsf{mm}}$ holds regardless of the reweighting by $1/\hat{p}_g$: if we define the ERM estimator $\hat{w}^{\mathsf{erm}}$ analogously to (5) without the reweighting, it is also equal to $\hat{w}^{\mathsf{mm}}$. We will therefore analyze $\hat{w}^{\mathsf{mm}}$ in the overparameterized regime since it subsumes both $\hat{w}^{\mathsf{rw}}$ and $\hat{w}^{\mathsf{erm}}$.

We also note that if we use gradient descent to directly optimize the unregularized logistic regression objective (either reweighted or not), the resulting solution after scaling to unit norm also converges to $\hat{w}^{\mathsf{mm}}$ as the number of gradient steps goes to infinity (Soudry et al., 2018).

### 5.2. Analysis of worst-group error

We now state our main analytical result: in the explicit-memorization setting, the worst-group test error of a sufficiently overparameterized model is greater than $1/2$ (worse than random) under certain settings of $\sigma_{\mathsf{spu}}^2, \sigma_{\mathsf{core}}^2, n_{\mathsf{maj}}, n_{\mathsf{min}}$. In contrast, underparameterized models attain reasonable worst-group error even under such a setting.

**Theorem 1.** *For any $p_{\mathsf{maj}} \geq \left(1 - \frac{1}{2001}\right)$, $\sigma_{\mathsf{core}}^2 \geq 1$, $\sigma_{\mathsf{spu}}^2 \leq \frac{1}{16 \log 100 n_{\mathsf{maj}}}$, $\sigma_{\mathsf{noise}}^2 \leq \frac{n_{\mathsf{maj}}}{600^2}$ and $n_{\mathsf{min}} \geq 100$, there exists $N_0$ such that for all $N > N_0$ (overparameterized regime), with high probability over draws of the data,*

$$Err_{\mathsf{wg}}(\hat{w}^{\mathsf{mm}}) \geq 2/3, \tag{7}$$

*where $\hat{w}^{\mathsf{mm}}$ is the max-margin classifier.*

*However, for $N = 0$ (underparameterized regime), with $p_{\mathsf{maj}} = \left(1 - \frac{1}{2001}\right)$, $\sigma_{\mathsf{core}}^2 = 1$, and $\sigma_{\mathsf{spu}}^2 = 0$, and in the asymptotic regime with $n_{\mathsf{maj}}, n_{\mathsf{min}} \to \infty$, we have*

$$Err_{\mathsf{wg}}(\hat{w}^{\mathsf{rw}}) < 1/4, \tag{8}$$

*where $\hat{w}^{\mathsf{rw}}$ minimizes the reweighted logistic loss.*

The result in the overparameterized regime applies to the max-margin classifier $\hat{w}^{\mathsf{mm}}$, which as discussed above subsumes both $\hat{w}^{\mathsf{rw}}$ and $\hat{w}^{\mathsf{erm}}$ when the data is linearly separable. The proof of Theorem 1 appears in Appendix B.

The conditions on $\sigma^2_{\mathsf{spu}}$ and $\sigma^2_{\mathsf{core}}$ in Theorem 1 above imply high spurious-core information ratio $r_{\mathsf{s:c}}$. Theorem 1 therefore provides a setting where high $p_{\mathsf{maj}}$ and high $r_{\mathsf{s:c}}$ provably make overparameterized models obtain high worst-group error, matching the trends we observed upon varying $p_{\mathsf{maj}}$ and $r_{\mathsf{s:c}}$ in the implicit-memorization setting (Figure 6). Furthermore, underparameterized models obtain reasonable worst-group error despite these conditions, mirroring the observations in earlier sections.

## 5.3. Overparameterization and memorization

We now sketch the key ideas in the proof of Theorem 1 (full proof in Appendix B), focusing first on the overparameterized regime. We start by establishing an inductive bias towards learning the minimum-norm model that fits the training data. We then define memorization and show how the minimum-norm inductive bias translates into a bias against memorization. Finally, we illustrate how the bias against memorization leads to learning the spurious feature and suffering high worst-group error.

**Minimum-norm inductive bias.** Define a *separator* as any model that correctly classifies all of the training points $(x, y)$ with margin $yw \cdot x \geq 1$. Then from standard duality arguments, $\hat{w}^{\mathsf{mm}}$ can be rewritten as $\hat{w}^{\mathsf{minnorm}}/\|\hat{w}^{\mathsf{minnorm}}\|$, the scaled version of the *minimum-norm separator* $\hat{w}^{\mathsf{minnorm}}$

$$\hat{w}^{\mathsf{minnorm}} \stackrel{\text{def}}{=} \underset{w \in \mathbb{R}^{N+2}}{\arg\min} \|w\|_2^2 \text{ s.t. } y^{(i)}(w \cdot x^{(i)}) \geq 1 \ \forall i. \quad (9)$$

Since scaling does not affect the 0-1 test error, it suffices to analyze $\hat{w}^{\mathsf{minnorm}}$. Equation (9) shows that out of the set of all separators (which all perfectly fit the training data), the inductive bias favors the separator with the minimum norm. We now discuss how this minimum-norm inductive bias favors less memorization.

**Memorization.** For convenience, we denote the three components of a model $w$ as

$$w = [w_{\mathsf{core}}, w_{\mathsf{spu}}, w_{\mathsf{noise}}], \quad (10)$$

where $w_{\mathsf{core}} \in \mathbb{R}$, $w_{\mathsf{spu}} \in \mathbb{R}$, and $w_{\mathsf{noise}} \in \mathbb{R}^N$. By the representer theorem, we can decompose $w_{\mathsf{noise}}$ as follows:

$$w_{\mathsf{noise}} = \sum_i \alpha^{(i)} x^{(i)}_{\mathsf{noise}}. \quad (11)$$

In the overparameterized regime when $N \gg n$, a model can "memorize" a training point $x^{(i)}$ via $w_{\mathsf{noise}}$, in particular by putting a large weight $\alpha^{(i)}$ in the direction of $x^{(i)}$ (Equation (11)):

**Definition 1** ($\gamma$-memorization). *A model $w$ memorizes a point $x^{(i)}$ if $|\alpha^{(i)}| \geq \gamma^2/\sigma^2_{\mathsf{noise}}$ for some constant $\gamma \in \mathbb{R}$.*

Because the noise vectors of the training points (high-dimensional Gaussians) are nearly orthogonal for large $N$, the component $\alpha^{(i)} x^{(i)}_{\mathsf{noise}}$ affects the prediction on $x^{(i)}$, but not on any other training or test points.

This ability to memorize plays a crucial role in making overparameterized models obtain high worst-group error. Intuitively, the minimum-norm inductive bias favors less memorization in overparameterized models. Roughly speaking, models that memorize more have larger weights $|\alpha^{(i)}|$ on the noise vectors $x^{(i)}_{\mathsf{noise}}$. Since these noise vectors are nearly orthogonal and have similar norm, this translates into a larger norm $\|w_{\mathsf{noise}}\|_2^2$.

**Comparing using $x_{\mathsf{core}}$ versus using $x_{\mathsf{spu}}$.** To illustrate how the inductive bias against memorization leads to high worst-group error, we consider two extreme sets of separators: (i) ones that use the spurious feature but not the core feature, denoted by $\mathcal{W}^{\mathsf{use-spu}}$ (ii) ones that use the core feature but not the spurious feature, denoted by $\mathcal{W}^{\mathsf{use-core}}$.

$$\mathcal{W}^{\mathsf{use-spu}} \stackrel{\text{def}}{=} \{w \in \mathbb{R}^{N+2} : w \text{ is a separator}, w_{\mathsf{core}} = 0\}$$
$$\mathcal{W}^{\mathsf{use-core}} \stackrel{\text{def}}{=} \{w \in \mathbb{R}^{N+2} : w \text{ is a separator}, w_{\mathsf{spu}} = 0\}. \quad (12)$$

In scenario (i), using the spurious feature $x_{\mathsf{spu}}$ alone allows models to fit the majority groups very well. Thus, models that use $x_{\mathsf{spu}}$ only need to memorize the minority points. In Proposition 1, we construct a separator $w^{\mathsf{use-spu}} \in \mathcal{W}^{\mathsf{use-spu}}$ and show that its norm *only* scales with the number of minority points $n_{\mathsf{min}}$.

Conversely, in scenario (ii), using the core feature $x_{\mathsf{core}}$ alone allows models to fit all groups equally well. However, when $r_{\mathsf{s:c}}$ is high, $x_{\mathsf{core}}$ is noisier than $x_{\mathsf{spu}}$, so models that use $x_{\mathsf{core}}$ still need to memorize a constant fraction of *all* the training points. In Proposition 2, we show that norms of all separators $w^{\mathsf{use-core}} \in \mathcal{W}^{\mathsf{use-core}}$ are lower bounded by a quantity linear in the total number of training points $n$.

When the majority fraction $p_{\mathsf{maj}}$ is sufficiently large such that $n_{\mathsf{min}} \ll n$, the separator $w^{\mathsf{use-spu}}$ that uses $x_{\mathsf{spu}}$ will have a lower norm than any separator $w^{\mathsf{use-core}} \in \mathcal{W}^{\mathsf{use-core}}$ that uses $x_{\mathsf{core}}$. Since the inductive bias favors the minimum-norm separator, it prefers a separator $w^{\mathsf{use-spu}}$ that memorizes the minority points and suffers high worst-group error over any $w^{\mathsf{use-core}} \in \mathcal{W}^{\mathsf{use-core}}$.

**Proposition 1** (Norm of models using the spurious feature). *When $\sigma^2_{\mathsf{core}}, \sigma^2_{\mathsf{spu}}$ satisfy the conditions in Theorem 1, there exists $N_0$ such that for all $N > N_0$, with high probability,*

*there exists a separator $w^{\text{use-spu}} \in \mathcal{W}^{\text{use-spu}}$ such that*

$$\|w^{\text{use-spu}}\|_2^2 \leq \gamma_1^2 + \left(\frac{\gamma_2 n_{\min}}{\sigma_{\text{noise}}^2}\right),$$

*for some constants $\gamma_1, \gamma_2 > 0$.*

*Proof sketch.* To simplify exposition in this sketch, suppose that the noise vectors $x_{\text{noise}}^{(i)}$ are orthogonal and have constant norm $\|x_{\text{noise}}^{(i)}\|_2^2 = \sigma_{\text{noise}}^2$. We construct a separator $w^{\text{use-spu}} \in \mathcal{W}^{\text{use-spu}}$ that does not use the core feature $x_{\text{core}}$ as follows. Set $w_{\text{spu}}^{\text{use-spu}} = \gamma_1$ for some large enough constant $\gamma_1 > 0$. This is sufficient to satisfy the margin condition on the majority points: since $\sigma_{\text{spu}}^2$ is very small, w.h.p. all majority training points satisfy $y^{(i)}(x_{\text{spu}}^{(i)}\gamma_1) \geq 1$.

However, for the minority training points, the spurious attribute $a$ does not match the label $y$, and in order to satisfy the margin condition with a positive $w_{\text{spu}}^{\text{use-spu}}$, these $n_{\min}$ minority points have to be memorized. Since $\sigma_{\text{spu}}^2$ is very small, the decrease in the margin due to $w_{\text{spu}}^{\text{use-spu}} = \gamma_1$ is at most $-\rho\gamma_1$ w.h.p. for some constant $\rho$ that depends on $\sigma_{\text{spu}}^2$. To satisfy the margin condition, it thus suffices to set $\alpha_{\text{use-spu}}^{(i)} = y^{(i)}(1 + \rho\gamma_1)/\sigma_{\text{noise}}^2$, and the bound on the norm follows. The full proof appears in Section B.2.6. $\square$

**Proposition 2** (Norm of models using the core feature). *When $\sigma_{\text{core}}^2, \sigma_{\text{spu}}^2$ satisfy the conditions in Theorem 1 and $n_{\min} \geq 100$, there exists $N_0$ such that for all $N > N_0$, with high probability, all separators $w^{\text{use-core}} \in \mathcal{W}^{\text{use-core}}$ satisfy*

$$\|w^{\text{use-core}}\|_2^2 \geq \frac{\gamma_3 n}{\sigma_{\text{noise}}^2},$$

*for some constant $\gamma_3 > 0$.*

*Proof sketch.* Any model $w^{\text{use-core}} \in \mathcal{W}^{\text{use-core}}$ has $w_{\text{spu}}^{\text{use-core}} = 0$ by definition. We show that a constant fraction of training points have to be $\gamma$-memorized in order to satisfy the margin condition. We do so by first showing that the probability that a training point $x$ satisfies the margin condition *without* being $\gamma$-memorized cannot be too large. For simplicity, suppose again that the noise vectors $x_{\text{noise}}^{(i)}$ are orthogonal and have constant norm $\|x_{\text{noise}}^{(i)}\|_2^2 = \sigma_{\text{noise}}^2$. Then this probability is $\mathbb{P}(x_{\text{core}}w_{\text{core}}^{\text{use-core}} \leq 1 - \gamma^2) \geq \Phi(-1/\sigma_{\text{core}})$ for small $\gamma$, where $\Phi$ is the Gaussian CDF. Hence, in expectation, at least a constant fraction of points from the training distribution need to be memorized in order for $w^{\text{use-core}}$ to satisfy the margin condition. With high probability, this is also true on the training set consisting of $n$ points (via the DKW inequality) and the bound on the norm follows. The full proof appears in Section B.2.7. $\square$

In the full proof of Theorem 1 in Appendix B, we generalize the above ideas to consider all separators in $\mathbb{R}^{N+2}$ instead of just the separators in $\mathcal{W}^{\text{use-spu}} \bigcup \mathcal{W}^{\text{use-core}}$. Note the importance of both $r_{\text{s:c}}$ and $p_{\text{maj}}$: when $r_{\text{s:c}}$ is high, models that use $x_{\text{spu}}$ only need to memorize the minority groups (Proposition 1), and when $p_{\text{maj}}$ is also high, these models end up memorizing fewer points than models that use $x_{\text{core}}$ and have to memorize a constant fraction of the entire training set (Proposition 2).

## 6. Subsampling

Our results above highlight the role of the majority fraction $p_{\text{maj}}$ in determining if overparameterization hurts worst-group test error. When $p_{\text{maj}}$ is large, the inductive bias favors using spurious features because it entails memorizing only a relatively small number of minority points, while the alternative of using core features requires memorizing a large number of majority points. This suggests that reducing the memorization cost of using core features by directly removing some majority points could induce overparameterized models to obtain low worst-group error.

Here, we show that this approach of *subsampling* the majority group achieves good worst-group test error on the datasets studied above. Subsampling creates a new *group-balanced* dataset by randomly removing training points in all other groups to match the number of points from the smallest group (Japkowicz & Stephen, 2002; Haixiang et al., 2017; Buda et al., 2018). We then train a model to minimize the average loss on this subsampled dataset. For a precise description, see Appendix A.6.

Figure 7 shows that overparameterized models trained via subsampling (Equation 15) obtain low worst-group error on the CelebA, Waterbirds, and synthetic (implicit-memorization) datasets. Across all three datasets, training via subsampling makes increasing overparameterization help *both* average and worst-group test error. Moreover, overparameterized models trained on subsampled data are comparable to or better than the best models trained on the full dataset (i.e., underparameterized models trained with reweighting).
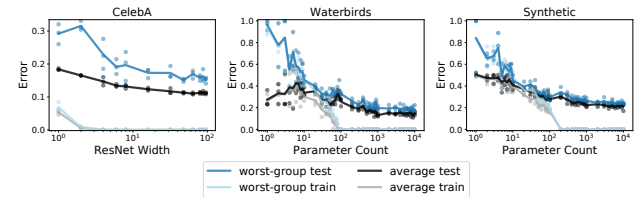


*Figure 7.* Overparameterization helps worst-group test error when training via subsampling, which involves creating a group-balanced dataset by reducing the number of majority points and minimizing average training loss on the new dataset.

Subsampling seems wasteful since it throws away a large fraction of the training data: we only use 3.4% of the full training data for CelebA, 4.6% for Waterbirds, and 10% for the synthetic dataset. However, the results above show that subsampling in overparameterized models matches or outperforms reweighting with underparameterized models. For example, on CelebA, an overparameterized model trained via subsampling obtains 11.1% average test and 15.1% worst-group test error, whereas an underparameterized model trained with reweighting obtains 11.3% average and 25.6% worst-group test error.

**Subsampling vs. reweighting.** Both subsampling and reweighting artificially balance the groups in the training data, and previous work on imbalanced datasets has concluded that reweighting is typically at least as effective as subsampling (Buda et al., 2018). However, we find a clear difference between subsampling and reweighting in the overparameterized regime: increasing overparameterization with reweighting increases worst-group error, while doing so with subsampling decreases worst-group error. The intuition developed in Sections 4 and 5 shed some light on this difference. Consider an overparameterized model: as in Section 5.1, reweighting does not change the learned model which is the max-margin classifier. However, subsampling reduces $p_{\mathsf{maj}}$. Recall that the inductive bias favors spurious features when the alternative of using core features requires memorizing a large number of training points. By reducing $p_{\mathsf{maj}}$, we reduce this memorization cost associated with core features, thereby inducing the model to use core features and achieve low worst-group test error.

## 7. Related work

**The effect of overparameterization.** The effect of overparameterization on average test error has been widely studied. In what is commonly referred to as "double descent", increasing model size beyond zero training error decreases test error, despite conventional wisdom that overfitting should increase test error. This behavior has been observed empirically (Belkin et al., 2019; Opper, 1995; Advani & Saxe, 2017; Nakkiran et al., 2019) and shown analytically in high-dimensional regression (Hastie et al., 2019; Bartlett et al., 2019; Mei & Montanari, 2019). These works focus on average test error and are consistent with our findings there. However, our focus is on worst-group test error, particularly when the groups are defined based on spurious attributes, and in this paper we establish that worst-group test error can behave quite differently from average test error.

Increasing overparameterization can actually improve model robustness to some types of distributional shifts (Hendrycks et al., 2019; Hendrycks & Dietterich, 2019; Yang et al., 2020). In this light, our results show that the effect of overparameterization on model robustness can depend heavily on the dataset (e.g., properties like $p_{\mathsf{maj}}$ and $r_{\mathsf{s:c}}$), type of distributional shift, and training procedure.

**Worst-group error.** Prior work on improving worst-group error focused on the underparameterized regime, with methods based on weighting/sampling (Shimodaira, 2000; Japkowicz & Stephen, 2002; Buda et al., 2018; Cui et al., 2019), distributionally robust optimization (DRO) (Ben-Tal et al., 2013; Namkoong & Duchi, 2017; Oren et al., 2019), and fair algorithms (Dwork et al., 2012; Hardt et al., 2016; Kleinberg et al., 2017). Our focus is on the overparameterized, zero-training-error regime; here, previous methods based on reweighting and DRO are ineffective (Wen et al., 2014; Byrd & Lipton, 2019; Sagawa et al., 2020). As mentioned in Section 1, Sagawa et al. (2020) demonstrated that stronger $L_2$-regularization can improve worst-group error on neural networks (when coupled with reweighting or group DRO). Similarly Cao et al. (2019) show that data-dependent regularization can improve error on rare labels. While their work focuses on developing methods to improve worst-group error, our focus is on understanding the mechanisms by which overparameterization hurts worst-group error.

## 8. Discussion

Our work shows that overparameterization hurts worst-group error on real datasets that contain spurious correlations. We studied the implicit- and explicit-memorization settings to provide a potential story for why this might occur: there can be an inductive bias towards solutions that do not need to memorize as many training points, and this can favor models that exploit the spurious correlations.

However, our synthetic settings make several simplifying assumptions, e.g., they suppose that the model prefers the spurious feature because it is less noisy than the core feature. This assumption need not always apply, and different assumptions might also lead to overparameterization exacerbating spurious correlations. For example, there might exist a true classifier based on the core features which has high accuracy but which is relatively more complex (e.g., high parameter norm) and therefore not favored by the training procedure. Studying the effect of overparameterization in settings such as those is important future work.

We also observed that subsampling allows overparameterized models to achieve low average and worst-group test error, despite eliminating a large fraction of training examples. In contrast, when using the full training data, only underparameterized models attain low worst-group test error under our current training methods. These observations call for future work to develop methods that can exploit both the statistical information in the full training data as well as the expressivity of overparameterized models, so as to attain good worst-group and average test error.

## Reproducibiltity

Code is available at https://github.com/ssagawa/overparam_spur_corr.
All code, data, and experiments are available on the Codalab platform at https://worksheets.codalab.org/worksheets/0x1db77e603a8d48c8abebd67fce39cf8b.

## References

Advani, M. S. and Saxe, A. M. High-dimensional dynamics of generalization error in neural networks. *arXiv preprint arXiv:1710.03667*, 2017.

Bartlett, P. L., Long, P. M., Lugosi, G., and Tsigler, A. Benign overfitting in linear regression. *arXiv*, 2019.

Belkin, M., Hsu, D., Ma, S., and Mandal, S. Reconciling modern machine-learning practice and the classical bias–variance trade-off. *Science*, 116(32), 2019.

Ben-Tal, A., den Hertog, D., Waegenaere, A. D., Melenberg, B., and Rennen, G. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59:341–357, 2013.

Blodgett, S. L., Green, L., and O'Connor, B. Demographic dialectal variation in social media: A case study of African-American English. In *Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1119–1130, 2016.

Buda, M., Maki, A., and Mazurowski, M. A. A systematic study of the class imbalance problem in convolutional neural networks. *Neural Networks*, 106:249–259, 2018.

Buolamwini, J. and Gebru, T. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on Fairness, Accountability and Transparency*, pp. 77–91, 2018.

Byrd, J. and Lipton, Z. What is the effect of importance weighting in deep learning? In *International Conference on Machine Learning (ICML)*, pp. 872–881, 2019.

Cao, K., Wei, C., Gaidon, A., Arechiga, N., and Ma, T. Learning imbalanced datasets with label-distribution-aware margin loss. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.

Cui, Y., Jia, M., Lin, T., Song, Y., and Belongie, S. Class-balanced loss based on effective number of samples. In *Computer Vision and Pattern Recognition (CVPR)*, pp. 9268–9277, 2019.

Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. Fairness through awareness. In *Innovations in Theoretical Computer Science (ITCS)*, pp. 214–226, 2012.

Haixiang, G., Yijing, L., Shang, J., Mingyun, G., Yuanyue, H., and Bing, G. Learning from class-imbalanced data: Review of methods and applications. *Expert Systems with Applications*, 73:220–239, 2017.

Hardt, M., Price, E., and Srebo, N. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pp. 3315–3323, 2016.

Hashimoto, T. B., Srivastava, M., Namkoong, H., and Liang, P. Fairness without demographics in repeated loss minimization. In *International Conference on Machine Learning (ICML)*, 2018.

Hastie, T., Montanari, A., Rosset, S., and Tibshirani, R. J. Surprises in high-dimensional ridgeless least squares interpolation. *arXiv preprint arXiv:1903.08560*, 2019.

He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *Computer Vision and Pattern Recognition (CVPR)*, 2016.

Hendrycks, D. and Dietterich, T. Benchmarking neural network robustness to common corruptions and perturbations. *arXiv preprint arXiv:1903.12261*, 2019.

Hendrycks, D., Zhao, K., Basart, S., Steinhardt, J., and Song, D. Natural adversarial examples. *arXiv preprint arXiv:1907.07174*, 2019.

Hu, W., Niu, G., Sato, I., and Sugiyama, M. Does distributionally robust supervised learning give robust classifiers? In *International Conference on Machine Learning (ICML)*, 2018.

Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International Conference on Machine Learning (ICML)*, pp. 448–456, 2015.

Japkowicz, N. and Stephen, S. The class imbalance problem: A systematic study. *Intelligent data analysis*, 6(5):429–449, 2002.

Kleinberg, J., Mullainathan, S., and Raghavan, M. Inherent trade-offs in the fair determination of risk scores. In *Innovations in Theoretical Computer Science (ITCS)*, 2017.

Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning face attributes in the wild. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3730–3738, 2015.

McCoy, R. T., Pavlick, E., and Linzen, T. Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference. In *Association for Computational Linguistics (ACL)*, 2019.

Mei, S. and Montanari, A. The generalization error of random features regression: Precise asymptotics and double descent curve. *arXiv preprint arXiv:1908.05355*, 2019.

Nakkiran, P., Kaplun, G., Bansal, Y., Yang, T., Barak, B., and Sutskever, I. Deep double descent: Where bigger models and more data hurt. *arXiv preprint arXiv:1912.02292*, 2019.

Namkoong, H. and Duchi, J. Variance regularization with convex objectives. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

Opper, M. Statistical mechanics of learning: Generalization. *The Handbook of Brain Theory and Neural Networks,*, pp. 922–925, 1995.

Oren, Y., Sagawa, S., Hashimoto, T., and Liang, P. Distributionally robust language modeling. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2019.

Rosset, S., Zhu, J., and Hastie, T. J. Margin maximizing loss functions. In *Advances in neural information processing systems*, pp. 1237–1244, 2004.

Sagawa, S., Koh, P. W., Hashimoto, T. B., and Liang, P. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. In *International Conference on Learning Representations (ICLR)*, 2020.

Shimodaira, H. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*, 90:227–244, 2000.

Soudry, D., Hoffer, E., Nacson, M. S., Gunasekar, S., and Srebro, N. The implicit bias of gradient descent on separable data. *Journal of Machine Learning Research (JMLR)*, 19(1):2822–2878, 2018.

Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research (JMLR)*, 15(1):1929–1958, 2014.

Wah, C., Branson, S., Welinder, P., Perona, P., and Belongie, S. The Caltech-UCSD Birds-200-2011 dataset. Technical report, California Institute of Technology, 2011.

Wen, J., Yu, C., and Greiner, R. Robust learning under uncertain test distributions: Relating covariate shift to model misspecification. In *International Conference on Machine Learning (ICML)*, pp. 631–639, 2014.

Yang, Z., Yu, Y., You, C., Steinhardt, J., and Ma, Y. Rethinking bias-variance trade-off for generalization of neural networks. *arXiv preprint arXiv:2002.11328*, 2020.

Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals, O. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations (ICLR)*, 2017.

Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., and Torralba, A. Places: A 10 million image database for scene recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(6):1452–1464, 2017.

# A. Supplemental experiments

## A.1. ERM models have poor worst-group error regardless of the degree of overparameterization

In the main text, we focused on reweighted models, trained with the reweighted objective on the full data (Sections 3-5), as well as subsampled models, trained on subsampled data with the ERM objective (Section 6). Here, we study the effect of overparameterization on ERM models, trained with the ERM objective on the full data. Consistent with prior work, we observe that ERM models obtain poor worst-group error (near or worse than random), regardless of whether the model is underparameterized or overparameterized (Sagawa et al., 2020). We also confirm that overparameterization helps average test error (see, e.g., Nakkiran et al. (2019); Belkin et al. (2019); Mei & Montanari (2019)).

**Empirical results.** We first consider the CelebA and Waterbirds dataset, following the experimental set-up of Section 3 but now training with the standard ERM objective (Equation (2)) instead of the reweighted objective (Equation (3)).

On these datasets, overparameterization helps the average test error (Figure 8). As model size increases past the point of zero training error, the average test error decreases. The best average test error is obtained by highly overparameterized models with zero training error—4.6% for CelebA at width 96, and 4.2% for Waterbirds at 6,000 random features.

In contrast, the worst-group error is consistently high across model sizes: it is consistently worse than random ($>50\%$) for CelebA and nearly random (44%) for Waterbirds (Figure 8). These worst-group errors are much worse than those obtained by reweighted, underparameterized models (25.6% for CelebA and 26.6% for Waterbirds; see Section 3). Thus, while overparameterization helps ERM models achieve better test error, these models all fail to yield good worst-group error regardless of the degree of overparameterization.
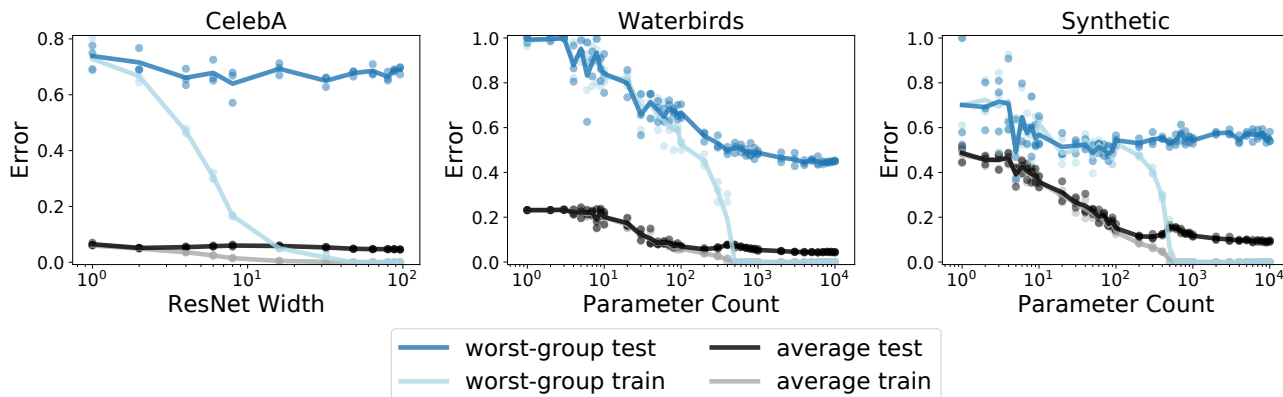


*Figure 8.* The effect of overparameterization on the average and worst-group error of an ERM model. Increasing model size helps average test error, but worst-group error remains poor across model sizes.

**Simulation results.** We also evaluate the effect of overparameterization on ERM models on the synthetic dataset introduced in Section 4. As above, ERM models fail to achieve reasonable worst-group test error across model sizes, but improve in average test error as model size increases (Figure 8). The best average test error is obtained by a highly overparameterized model with zero training error—9.0% error at 9,000 random features—while the worst-group test error is nearly random or worse ($> 48\%$) across model sizes.

## A.2. Stronger $L_2$ regularization improves worst-group error in overparameterized reweighted models

In the main text, we studied models with default/weak or no $L_2$ regularization. In this section, we study the role of $L_2$ regularization in modulating the effect of overparameterization on worst-group error by changing the hyperparameter $\lambda$ that controls $L_2$ regularization strength. Overall, we find that increasing $L_2$ regularization (to the point where models do not have zero training error) improves worst-group error but hurts average error in overparameterized reweighted models. In contrast, $L_2$ regularization has little effect on both worst-group and average error in the underparameterized regime.

**Strong $L_2$ regularization improves worst-group error in overparameterized reweighted models.** In the main text, we trained ResNet10 models with default, weak regularization ($\lambda = 0.0001$) on the CelebA dataset, and unregularized logistic regression on the Waterbirds and synthetic datasets. Here, we consider strongly-regularized models with $\lambda = 0.1$ for both types of models; unlike before, these models no longer achieve zero training error even when overparameterized. Figure 9 shows the results of varying model size on strongly-regularized ERM, reweighted, and subsampled models on the three datasets.

On all three datasets, with strong regularization, ERM models continue to yield poor worst-group test error across model sizes, with similar or worse worst-group test error compared to with weak/ no regularization. Conversely, strongly-regularized subsampled models continue to achieve low worst-group test error across model sizes.

Where strong regularization has a large effect is on reweighted models. With reweighting, we find that strong regularization improves worst-group error in overparameterized models: across all three datasets, the worst-group test error in the overparameterized regime is much lower for the strongly-regularized models than their weakly regularized or unregularized counterparts (Figure 3). These results are consistent with similar observations made in Sagawa et al. (2020). However, even though strongly-regularized overparameterized models outperform weakly-regularized overparameterized models, overparameterization can still hurt the worst-group error in strongly-regularized reweighted models. On the CelebA and synthetic datasets, with $\lambda = 0.1$, the best worst-group error is still obtained by an underparameterized model for the CelebA and synthetic datasets, though overparameterization seems to help worst-group error on the Waterbirds dataset at least in the range of model sizes studied.



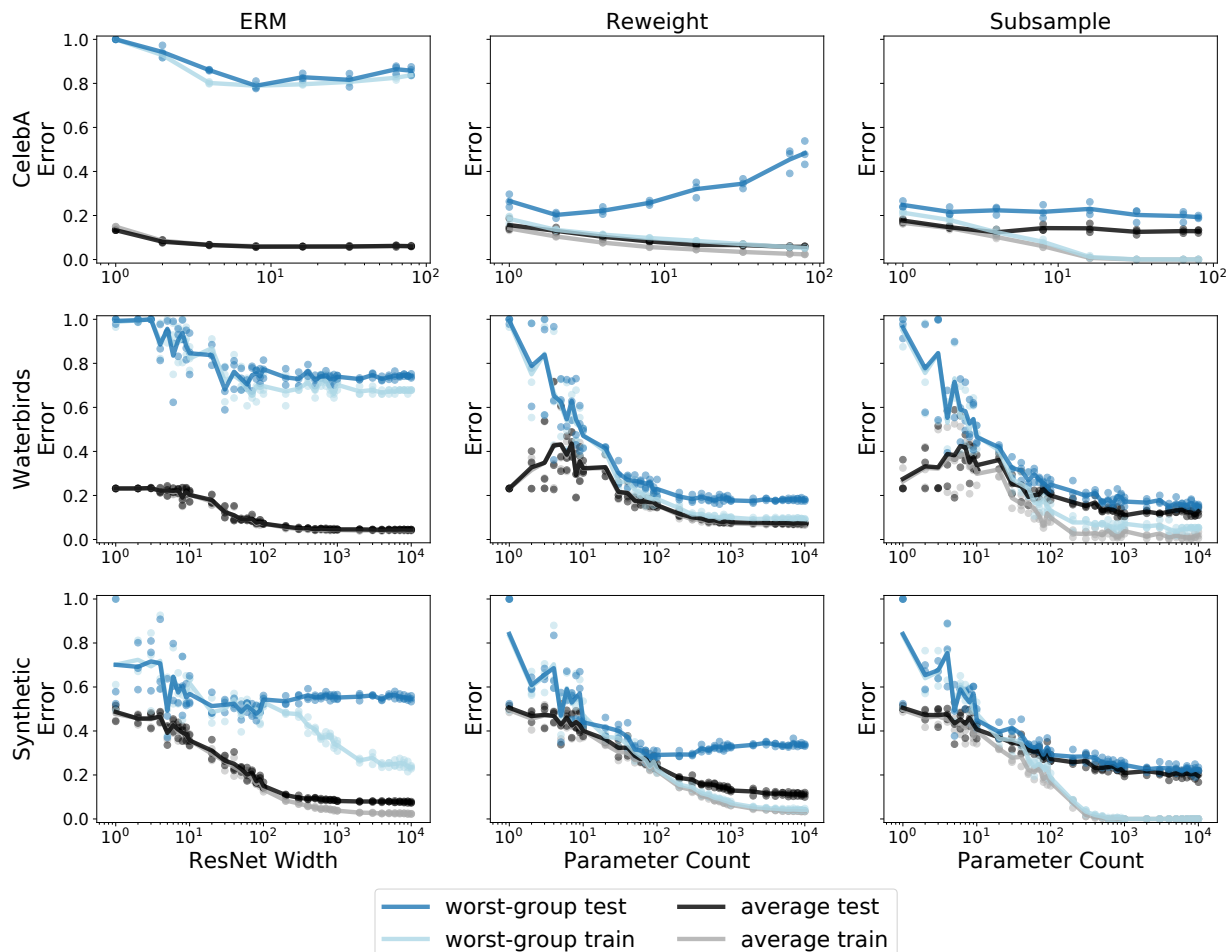*Figure 9.* Strongly-regularized models have lower worst-group error than their weakly-regularized counterparts in the overparameterized regime (Figure 3). Even under strong regularization, increasing model size can hurt the worst-group error on the CelebA (top) and synthetic (bottom) datasets, although overparameterization seems to improve worst-group error in the Waterbirds datase (middle) for the range of model sizes studied.

**Overparameterized models require strong regularization for worst-group test error but not average test error.** Given a fixed overparameterized model size, how does its performance change with the $L_2$ regularization strength $\lambda$? We study this with the logistic regression model on the Waterbirds and synthetic datasets, using a model size of $m = 10,000$ random features and varying the $L_2$ regularization strength from $\lambda = 10^{-9}$ to $\lambda = 10^2$. [1]

Results are in Figure 10. As before, ERM models obtain poor worst-group error regardless of the regularization strength, and subsampled models are relatively insensitive to regularization, achieving reasonable worst-group error at most settings of $\lambda$.

For reweighted models, however, having the right level of regularization is critical for obtaining good worst-group test error. On both datasets, the best worst-group test error is obtained by strongly-regularized models that do not achieve zero training error. In contrast, increasing regularization strength hurts average error, with the best average test error attained by models with nearly zero regularization.
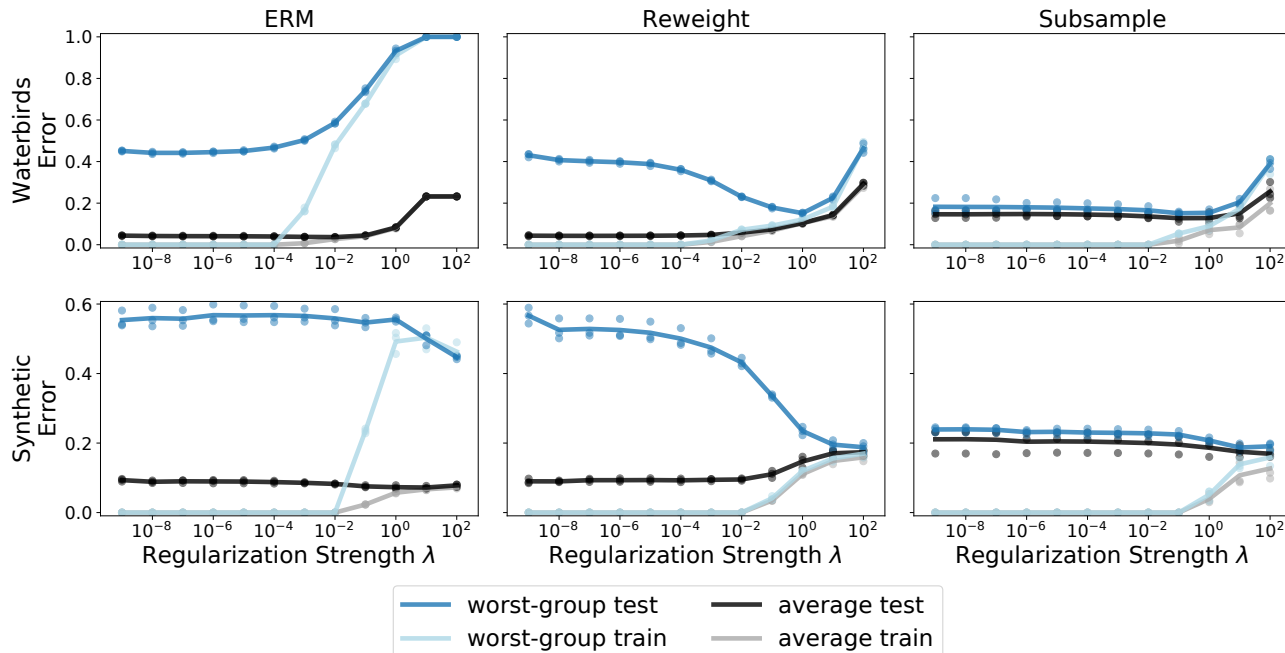


*Figure 10.* The effect of regularization on overparameterized random features logistic regression models ($m = 10,000$). ERM models (left) do consistently poorly while subsampled models (right) do consistently well on worst-group error. For reweighted models (middle), the best worst-group error is obtained by a strongly-regularized model that does not achieve zero training error.

$L_2$ **regularization affects where worst-group test error plateaus as model size increases.** In the above experiments, we kept either model size or regularization strength fixed, and varied the other. Here, we vary both: we consider $L_2$ regularization strengths $\lambda \in \{10^{-9}, 10^{-6}, 0.001, 0.1, 10\}$ and investigate the effect of increasing model size for each $\lambda$. We plot the results for Waterbirds and the synthetic dataset in Figure 11 and Figure 12 respectively.

For reweighted models, the results match what we observed above. Strengthening $L_2$ regularization reduces the detrimental effect of overparameterization on worst-group error. For any fixed model size in the overparameterized regime, the worst-group test error improves as $\lambda$ increases up to a certain value. Worst-group test error seems to plateau at different values as model size increases, depending on the regularization strength, though we note that it is possible that further increasing model size beyond the range we studied might lead models with different regularization strengths to eventually converge. Further empirical studies as well as theoretical characterization of the interaction between regularization and overparameterization are needed to confirm this phenomenon.

Given sufficiently large $\lambda$ (e.g., $\lambda = 10$ for both Waterbirds and synthetic datasets), overparameterized models seem to

---

[1]We did not run this experiment on the CelebA dataset for computational reasons, as doing so would have required tuning a different learning rate for each choice of regularization strength.

outperform underparameterized models, at least for the range of model sizes studied. However, we caution that this trend does not seem to hold on the CelebA dataset (Figure 9).

Finally, in contrast with its effects on overparameterized models, regularization seems to only have a modest effect on worst-group test error in the underparameterized regime.



*Figure 11.* The effect of overparameterization on models with different $L_2$ regularization strengths $\lambda$ on the Waterbirds dataset. Different regularization strengths are shown in different colors, with training and test errors plotted in light and dark colors, respectively.



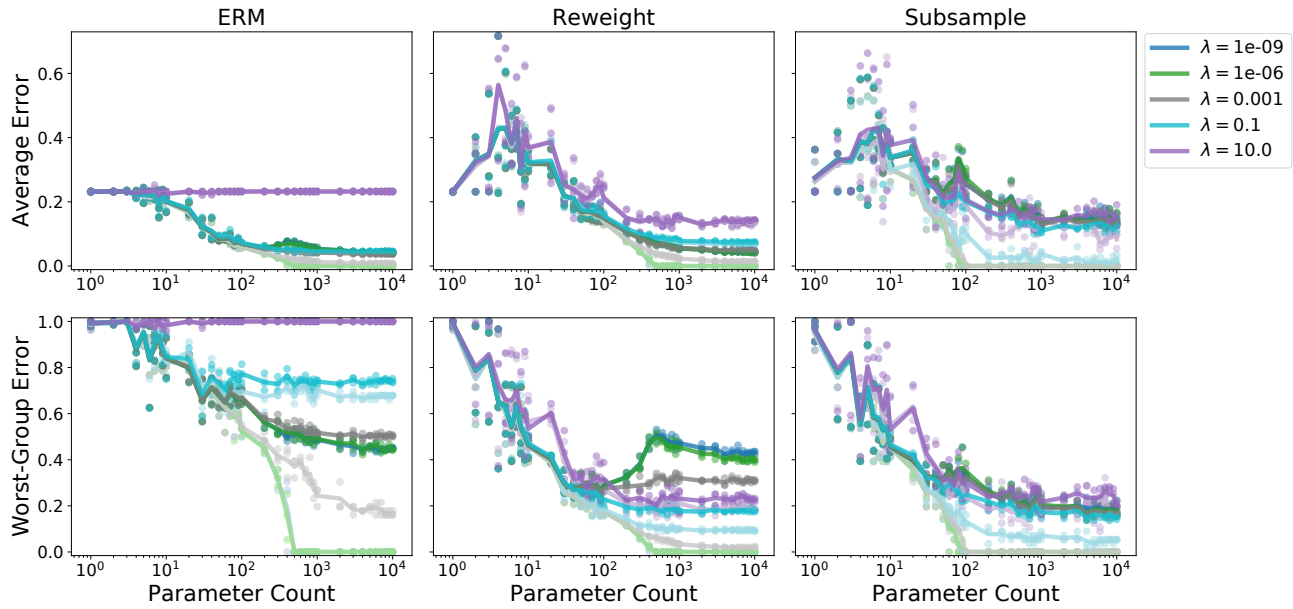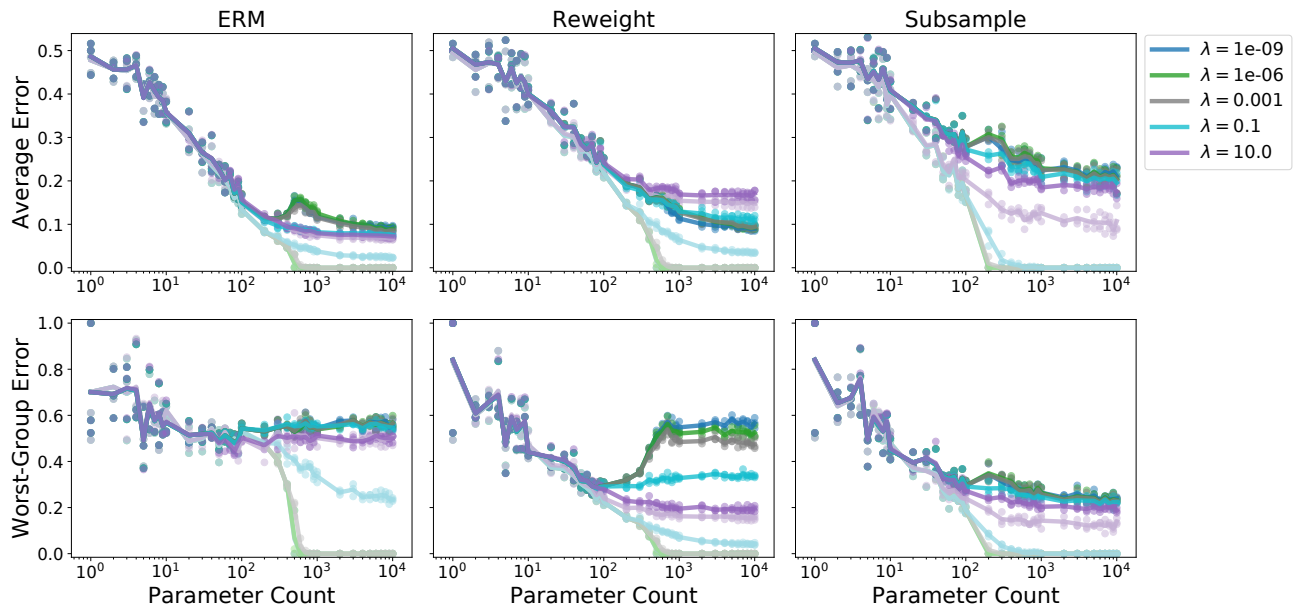*Figure 12.* The effect of overparameterization on models with different $L_2$ regularization strengths $\lambda$ on the synthetic dataset. The plotting scheme follows that of Figure 11.

### A.3. Overparameterization helps average test error on the synthetic data regardless of $p_{\mathsf{maj}}$ and $r_{\mathsf{s:c}}$

Figure 13 shows how the average test error changes as a function of model size under different settings of the majority fraction $p_{\mathsf{maj}}$ and the spurious-core ratio $r_{\mathsf{s:c}}$ on the synthetic dataset introduced in Section 4. As expected, overparameterization helps the average test error regardless of SCR and the majority fraction.
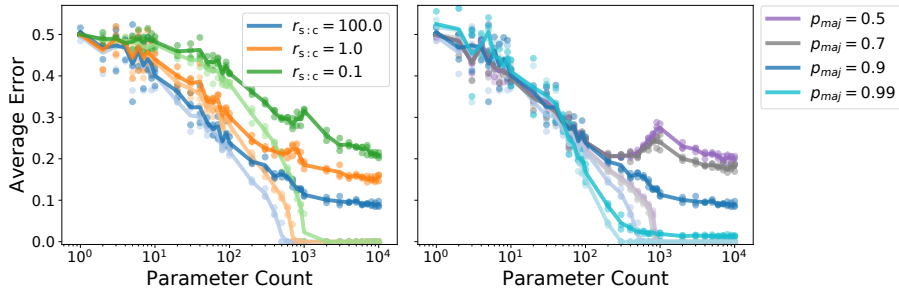


*Figure 13.* The effect of overparameterization on average error of a reweighted model on synthetic data. Different values of $p_{\mathsf{maj}}$ and $r_{\mathsf{s:c}}$ are plotted in different colors, with training and test errors plotted in light and dark colors, respectively. Across all values of $p_{\mathsf{maj}}$ and $r_{\mathsf{s:c}}$, overparameterization helps the average test error.

### A.4. Comparison between implicit and explicit implicit memorization

To motivate the explicit-memorization setting, we ran some brief experiments to show that in the overparameterized regime, linear models in the explicit-memorization setting behave similarly to random projection (RP) models in the implicit-memorization setting, with $\sigma^2_{\mathsf{core}}$ and $\sigma^2_{\mathsf{spu}}$ in the latter scaled up by a factor of $d$ (Figure 14). Recall that in the latter, $x_{\mathsf{core}} \in \mathbb{R}^d$ is distributed as $x_{\mathsf{core}}|y \sim \mathcal{N}(y, \sigma^2_{\mathsf{core}}I_d)$. Roughly speaking, all the information about $y$ is contained in the mean $\bar{x}_{\mathsf{core}} = \frac{1}{d}\sum_j x_{\mathsf{core},j}$, which is distributed as $\mathcal{N}(y, \sigma^2_{\mathsf{core}}I_d/d)$. In the explicit-memorization setting, we can view $x_{\mathsf{core}} \in \mathbb{R}$ as equivalent to $\bar{x}_{\mathsf{core}}$ in the implicit-memorization setting (and similarly for $x_{\mathsf{spu}}$), explaining the quantitative fit observed in Figure 14.

However, in the highly underparameterized regime, the RP models do poorly because of model misspecification (owing to a small number of random projections), whereas the linear models can still learn to use $x_{\mathsf{core}}$ and therefore do well.



*Figure 14.* The effect of overparameterization on the worst-group test error for linear models in the explicit-memorization setting ($\sigma^2_{\mathsf{core}} = 1, \sigma^2_{\mathsf{spu}} = 0.01, \sigma^2_{\mathsf{noise}} = 1$) and random projection models in the implicit-memorization setting ($\sigma^2_{\mathsf{core}} = 100, \sigma^2_{\mathsf{spu}} = 1, d = 100$). The models agree in the overparameterized regime.

### A.5. Experimental details

**Waterbirds and CelebA datasets.** For the CelebA dataset, we use the official train-val-test split from Liu et al. (2015), with the *Blond_Hair* attribute as the target $y$ and the *Male* as the spurious association $a$.

For the Waterbirds dataset, we follow the setup in Sagawa et al. (2020); for convenience, we reproduce some details of how it was constructed here. This dataset was obtained by combining bird images from the CUB dataset (Wah et al., 2011) with backgrounds from the Places dataset (Zhou et al., 2017). The CUB dataset comes with annotations of bird species. For the Waterbirds dataset, each bird was labeled was a waterbird if it was a seabird or waterfowl in the CUB dataset; otherwise, it was labeled as a landbird. Bird images were cropped using the provided segmentation masks and placed on either a land (bamboo forest or broadleaf forest) or water (ocean or natural lake) background obtained from the Places dataset.

For Waterbirds, we follow the same train-val-test split as in Sagawa et al. (2020). Note that in these validation and test sets,

landbirds and waterbirds are uniformly distributed on land and water backgrounds so that accuracy on the rare groups can be more accurately estimated. When calculating average test accuracy, we therefore first compute the average test accuracy over each group and then report a weighted average, with weights corresponding to the relative proportion of each group in the skewed training dataset.

We post-process Waterbirds by extracting feature representations taken from the last layer of a ResNet18 model pre-trained on ImageNet. We use the Pytorch `torchvision` implementation of the ResNet18 model for this. All models on the Waterbirds dataset in our paper are logistic regression models trained on top of this (fixed) feature representation.

**ResNet.** We used a modified ResNet10 with variable widths, following the approach in Nakkiran et al. (2019) and extending the `torchvision` implementation. We trained all ResNet10 models with stochastic gradient descent with momentum of 0.9 and a batch size of 128, with the $L_2$ regularization parameter $\lambda$ was passed in to the optimizer as the weight decay parameter. In the experiments in the main text, we used the default setting of $\lambda = 10^{-4}$. We used a fixed learning rate instead of a learning rate schedule and selected the largest learning rate for which optimization was stable, following Sagawa et al. (2020). This resulted in learning rates of 0.01 and 0.0001 for $\lambda = 10^{-4}$ and $\lambda = 0.1$, respectively, across all training procedures. As in the original ResNet paper (He et al., 2016), we used batch normalization (Ioffe & Szegedy, 2015) and no dropout (Srivastava et al., 2014), and for simplicity, we trained all models without data augmentation.

We trained for 50 epochs for ERM and reweighted models and 500 epochs for subsampled models (due to smaller number of examples per epoch). We found that worst-group error can be unstable across epochs due to the small sample size and relatively large learning rate, so in our results we report the error averaged over the last 10 epochs.

**Logistic regression.** We used the logistic regression implementation from `scikit-learn`, training with the L-BFGS solver until convergence with tolerance 0.0001, and setting the regularization parameter as $C = 1/(n\lambda)$. For unregularized models, we set $\lambda = 10^{-9}$ for numerical stability.

### A.6. Subsampling

Formally, given a set of groups $\mathcal{G}$ and a dataset D comprising a set of $n$ training points with their group identities $\{(x^{(i)}, y^{(i)}, g^{(i)})\}$, the subsampling procedure involves two steps. First, we group training points based on group identities:

$$\mathrm{D}_g \overset{\text{def}}{=} \{(x^{(i)}, y^{(i)}) \mid g^{(i)} = g\} \text{ for each } g \in \mathcal{G}. \tag{13}$$

For each group $g$, we select a subset $\mathrm{D}_g^{\text{ss}} \subseteq \mathrm{D}_g$ uniformly at random from $\mathrm{D}_g$ such that each subset has the same number of points as the smallest group in the training set. We form a new dataset $\mathrm{D}^{\text{ss}}$ by combining these subsets:

$$\mathrm{D}^{\text{ss}} = \bigcup_{g \in \mathcal{G}} \mathrm{D}_g^{\text{ss}}, \text{ where} \tag{14}$$

$$\mathrm{D}_g^{\text{ss}} \subseteq \mathrm{D}_g \ \text{ and } \ |\mathrm{D}_g^{\text{ss}}| = \min_{g \in \mathcal{G}} |\mathrm{D}_g|$$

Note that $\mathrm{D}^{\text{ss}}$ is group-balanced, with $p_{\text{maj}} = 0.5$. We then train a model by minimizing the average loss on $\mathrm{D}^{\text{ss}}$,

$$\hat{\mathcal{R}}_{\text{subsample}}(w) \overset{\text{def}}{=} \frac{1}{|\mathrm{D}^{\text{ss}}|} \sum_{(x,y) \in \mathrm{D}^{\text{ss}}} \ell(w; (x, y)). \tag{15}$$

Since $\mathrm{D}^{\text{ss}}$ is group-balanced, the reweighted training loss (Equation 3) has the same weight on all training points and minimizing the reweighted objective on $\mathrm{D}^{\text{ss}}$ is equivalent to minimizing the average loss objective above.

## B. Proof of Theorem 1

Here, we detail the proof of Theorem 1 presented in Section 5. We structure the proof by splitting Theorem 1 into two smaller theorems: one for the overparameterized regime (Appendix B.2), and another for the underparameterized regime (Appendix B.3).

## B.1. Notation and definitions.

We denote the separate components of the weight vector $\hat{w}_{\text{core}} \in \mathbb{R}, \hat{w}_{\text{spu}} \in \mathbb{R}, \hat{w}_{\text{noise}} \in \mathbb{R}^N$ such that

$$\hat{w} = [\hat{w}_{\text{core}}, \hat{w}_{\text{spu}}, \hat{w}_{\text{noise}}]. \tag{16}$$

Further, by the representer theorem, we decompose $\hat{w}_{\text{noise}}$ as

$$\hat{w}_{\text{noise}} = \sum_{i=1}^{n} \alpha^{(i)}(\hat{w}) x_{\text{noise}}^{(i)}. \tag{17}$$

Note that $\alpha^{(i)}(w)$ is equivalent to the $\alpha^{(i)}$ referred to in the main text. Recall that we define memorization of each training point $x^{(i)}$ by the weight $\alpha^{(i)}$ as follows.

**Definition 2** ($\gamma$-memorization). *Consider a separator $\hat{w}$ on training data $\{(x^{(i)}, y^{(i)})\}_{i=1}^{n}$. For some constant $\gamma \in \mathbb{R}$, we say that a model $\gamma$-memorizes a training point if*

$$\left| \alpha^{(i)}(\hat{w}) \right| > \frac{\gamma^2}{\sigma_{\text{noise}}^2}. \tag{18}$$

The component $\alpha^{(i)}(\hat{w}) x_{\text{noise}}^{(i)}$ serves to "memorize" $x^{(i)}$ when $N$ is sufficiently large, as it affects the prediction on $x^{(i)}$ but not on any other training or test points (because noise vectors are nearly orthogonal when $N$ is large). In the proof, we set the constant $\gamma^2$ appropriately (based on other parameter settings in Theorem 1) to get the required result.

Finally, let $G_{\text{maj}}, G_{\text{min}}$ denote the indices of training points in the majority and minority group respectively.

## B.2. Overparameterized regime

In our explicit-memorization set-up, sufficiently overparameterized models provably have high worst-group error under certain settings of $\sigma_{\text{spu}}^2, \sigma_{\text{core}}^2, n_{\text{maj}}, n_{\text{min}}$ as stated in Theorem 1 (restated below as Theorem 2).

**Theorem 2.** *For any $p_{\text{maj}} \geq \left(1 - \frac{1}{2001}\right)$, $\sigma_{\text{core}}^2 \geq 1$, $\sigma_{\text{spu}}^2 \leq \frac{1}{16 \log 100 n_{\text{maj}}}$, $\sigma_{\text{noise}}^2 \leq \frac{n_{\text{maj}}}{600^2}$ and $n_{\text{min}} \geq 100$, there exists $N_0$ such that for all $N > N_0$ (overparametrized regime), with high probability over draws of the data,*

$$Err_{\text{wg}}(\hat{w}^{\text{mm}}) \geq 2/3, \tag{19}$$

*where $\hat{w}^{\text{mm}}$ is the max-margin classifier.*

In Section 5, we sketched key ideas in the proof by considering special families of separators: because the minimum-norm inductive bias favors less memorization, models can prefer to learn the spurious feature and memorize the minority examples (entailing high worst-group error), instead of learning the core feature and memorizing some fraction of all training points (possibly attaining reasonable worst-group error). We now provide the full proof of Theorem 2, generalizing the above key concepts by considering *all* separators.

*Proof.* Recall from Section 5 that we consider the maximum-margin classifier $\hat{w}^{\text{minnorm}}$:

$$\hat{w}^{\text{minnorm}} = \arg\min \|w\|_2^2 \text{ s.t. } y^{(i)}(w \cdot x^{(i)}) \geq 1, \forall i. \tag{20}$$

In other words, $\hat{w}^{\text{minnorm}}$ is the minimum-norm separator, where separator is a classifier with zero training error and required margins, satisfying $y^{(i)}(w \cdot x^{(i)}) \geq 1$ for all $i$. We analyze the worst-group error of the minimum-norm separator $\hat{w}^{\text{minnorm}}$ as outlined below:

1. We first upper bound the fraction of *majority* examples memorized by the minimum-norm separator $\hat{w}^{\text{minnorm}}$. We show that there exists a separator that can use spurious features and needs to memorize only the minority points (Lemma 1) for the parameter settings in Theorem 2 where $\sigma_{\text{spu}}$ is sufficiently small. Since the norm of a separator is roughly scales with the number of points memorized ($|\alpha^{(i)}(\hat{w})| \geq \gamma^2/\sigma_{\text{noise}}^2$), we have an upper bound on the number of training points memorized by $\hat{w}^{\text{minnorm}}$. Since the number of majority points is much larger than the number of minority points, this says that only a small fraction of majority points could be memorized by $\hat{w}^{\text{minnorm}}$.

2. Next, we observe that since the core feature is noisy as per the parameter setting in Theorem 2, if we do not use the spurious feature, a constant fraction of majority points have to be memorized if spurious features are not used. Conversely, if less than this fraction of majority points can be memorized, the separator must use spurious features. Since using spurious features leads to higher worst-group test error, this reveals a trade-off between the worst-group test error of a separator and the fraction of *majority points* that it memorizes at training time. Succinctly, smaller fraction memorized implies the use of spurious features which in turn implies higher worst-group test error. Smaller worst-group test error requires eliminating the use of spurious features which would lead to a large fraction of majority points requiring memorization in order for a classifier to be a separator. We formalize the above trade-off between the worst-group test error and fraction of majority examples to be memorized in Proposition 3.

Combining the two steps together, since $\hat{w}^{\mathsf{minnorm}}$ memorizes only a small fraction of majority points by virtue of being the minimum norm separator, $\hat{w}^{\mathsf{minnorm}}$ suffers high worst-group test error.

We now formally prove Theorem 2, invoking propositions that we prove in subsequent sections.

### B.2.1. BOUNDING THE FRACTION OF MEMORIZED EXAMPLES IN THE MAJORITY GROUPS.

In the first part of the proof, we show that the minimum-norm separator $\hat{w}^{\mathsf{minnorm}}$ "memorizes" a small fraction of the majority examples. Formally, we study the quantity $\delta_{\mathsf{maj\text{-}train}}\left(\hat{w}, \gamma^2\right)$ defined as follows.

**Definition 3.** *Consider a separator $\hat{w}$ on training data $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$. Let $\delta_{maj\text{-}train}\left(\hat{w}, \gamma^2\right)$ be the fraction of training examples that $\hat{w}$ $\gamma$-memorizes in the majority groups:*

$$\delta_{maj\text{-}train}\left(\hat{w}, \gamma^2\right) \stackrel{\text{def}}{=} \frac{1}{n_{\mathsf{maj}}} \sum_{i \in G_{\mathsf{maj}}} \mathbb{I}\left[\left|\alpha^{(i)}(\hat{w})\right| > \frac{\gamma^2}{\sigma_{\mathsf{noise}}^2}\right] \tag{21}$$

We provide an upper bound on $\delta_{\mathsf{maj\text{-}train}}\left(\hat{w}^{\mathsf{minnorm}}, \gamma^2\right)$ (Lemma 4) by first bounding $\|\hat{w}^{\mathsf{minnorm}}\|$ and then bounding $\delta_{\mathsf{maj\text{-}train}}\left(\hat{w}^{\mathsf{minnorm}}, \gamma^2\right)$ in terms of $\|\hat{w}^{\mathsf{minnorm}}\|$.

<u>**Bounding $\|\hat{w}^{\mathsf{minnorm}}\|$**</u>

**Lemma 1.** *There exists a separator $w^{\mathsf{use-spu}}$ that satisfies $y^{(i)}(w^{\mathsf{use-spu}} \cdot x^{(i)}) \geq 1, \forall i \in G_{\mathsf{maj}}, G_{\mathsf{min}}$. The norm of this separator gives a bound on $\|\hat{w}^{\mathsf{minnorm}}\|$ as follows. For the parameter settings under Theorem 2, with high probability, we have*

$$\|\hat{w}^{\mathsf{minnorm}}\|_2^2 \leq \|w^{\mathsf{use-spu}}\|_2^2 \leq u^2 + s^2\sigma_{\mathsf{noise}}^2(1+c_1)n_{\mathsf{min}} + \frac{s^2\sigma_{\mathsf{noise}}^2}{n^4}, \tag{22}$$

*for constants $u = 1.3125, s = \frac{2.61}{\sigma_{\mathsf{noise}}^2}$.*

*Proof.* In order to get an upper bound on $\|\hat{w}^{\mathsf{minnorm}}\|$, we compute the norm of a particular separator. Concretely, we consider a separator $w^{\mathsf{use-spu}}$ of the following form:

$$w_{\mathsf{core}}^{\mathsf{use-spu}} = 0$$
$$w_{\mathsf{spu}}^{\mathsf{use-spu}} = u$$
$$w_{\mathsf{noise}}^{\mathsf{use-spu}} = \sum_i \alpha^{(i)}(w^{\mathsf{use-spu}})x_{\mathsf{noise}}^{(i)}$$
$$\alpha^{(i)}(w^{\mathsf{use-spu}}) = 0 \text{ for } i \in G_{\mathsf{maj}}$$
$$\alpha^{(i)}(w^{\mathsf{use-spu}}) = y^{(i)}s \text{ for } i \in G_{\mathsf{min}}$$

First, because we are interested in $w^{\mathsf{use-spu}}$ that does not use the core feature and relies on the spurious feature instead, we let $w_{\mathsf{core}}^{\mathsf{use-spu}} = 0$ and $w_{\mathsf{spu}}^{\mathsf{use-spu}} = u, u \in \mathbb{R}$. We set the value $u$ appropriately so that none of the majority points are memorized (corresponding to $\alpha^{(i)}(w^{\mathsf{use-spu}}) = 0$ for all $i \in G_{\mathsf{maj}}$). However since the spurious correlations are reversed in the minority

points and $w^{\text{use}-\text{spu}}_{\text{core}} = 0$, the minority points have to be memorized. For simplicity, we set $\alpha^{(i)}(w^{\text{use}-\text{spu}}) = y^{(i)}s$ for all $i \in G_{\text{min}}$.

Now it remains to select appropriate values of constants $u$ and $s$ such that $y^{(i)}(w^{\text{use}-\text{spu}} \cdot x^{(i)}) \geq 1$ is satisfied for all training examples.

For majority points, this involves setting $u$ large enough such that the less noisy spurious feature can be used to obtain the required margin. Without loss of generality, assume $y^{(i)} = 1$. Formally, for $i \in G_{\text{maj}}$,

$$
\begin{aligned}
w^{\text{use}-\text{spu}} \cdot x^{(i)} &\geq x^{(i)}_{\text{spu}}u + \sum_{j \in G_{\text{min}}} sx^{(i)}_{\text{noise}} \cdot x^{(j)}_{\text{noise}} \\
&\geq 4/5u + \sum_{j \in G_{\text{min}}} sx^{(i)}_{\text{noise}} \cdot x^{(j)}_{\text{noise}}, \text{ w.h.p. from Lemma 5 with } a = y = 1 \\
&\geq 4/5u - \frac{s\sigma^2_{\text{noise}}}{n^5}, \text{ w.h.p. from Lemma 8.} \\
&\geq 4/5u - \frac{s\sigma^2_{\text{noise}}}{100}.
\end{aligned}
$$

The first inequality follows from the fact that $\sigma_{\text{spu}}$ is small enough under the parameter settings of Theorem 2 to allow a uniform bound on $x^{(i)}_{\text{spu}}$ (Lemma 5). The second inequality follows from setting the number of random features $N$ to be large enough so that the noise features are near orthogonal (Lemma 8). Conversely, we have

$$
4/5u - \frac{s\sigma^2_{\text{noise}}}{100} \geq 1 \implies w^{\text{use}-\text{spu}} \text{ is a separator on the majority points w.h.p.} \tag{23}
$$

Notice that the condition in Equation 23 requires that $u$ be greater than $0$. Since the minority points have spurious attribute $a = -y$, we need to set $s$ to be large enough so that $w^{\text{use}-\text{spu}}$ as defined above separates the minority points. Just as before, we set $y = 1$ WLOG. For $i \in G_{\text{min}}$, we have

$$
\begin{aligned}
w^{\text{use}-\text{spu}} \cdot x^{(i)} &\geq x^{(i)}_{\text{spu}}u + \sum_{j \in G_{\text{min}}} sx^{(i)}_{\text{noise}} \cdot x^{(j)}_{\text{noise}} \\
&\geq -6/5u + \sum_{j \in G_{\text{min}}} sx^{(i)}_{\text{noise}} \cdot x^{(j)}_{\text{noise}}, \text{ From Lemma 5 with } a = -y = -1 \\
&\geq -6/5u + s(1 - c_1)\sigma^2_{\text{noise}} - \frac{s\sigma^2_{\text{noise}}}{n^5}, \text{ w.h.p from Lemma 8 and Lemma 9} \\
&\geq -6/5u + s(1 - c_1)\sigma^2_{\text{noise}} - \frac{s\sigma^2_{\text{noise}}}{100}.
\end{aligned}
$$

The steps are similar to the condition for majority points, with the key difference that the contribution from the noise term involves $s\|x^{(i)}_{\text{noise}}\|^2_2$ (Lemma 9).

Conversely, we have

$$
-6/5u + s(1 - c_1)\sigma^2_{\text{noise}} - \frac{s\sigma^2_{\text{noise}}}{100} \geq 1 \implies w^{\text{use}-\text{spu}} \text{ is a separator on the minority points w.h.p..} \tag{24}
$$

A set of parameters that satisfies both conditions above Equation 24 and Equation 23 is the following:

$$
u = 1.3125, s\sigma^2_{\text{noise}} = 2.61.
$$

We use the fact that $c_1 < 1/2000$ (From Lemma 9).

Finally, we have w.h.p,

$$
\|w^{\text{use}-\text{spu}}\|^2_2 \leq u^2 + s^2\sigma^2_{\text{noise}}(1 + c_1)n_{\text{min}} + \frac{s^2\sigma^2_{\text{noise}}}{n^4}. \tag{25}
$$

This follows from bounds on $\|x^{(i)}_{\text{noise}}\|^2_2$ (Lemma 9) and sum of less than $n^2$ terms involving $s^2x^{(i)}_{\text{noise}} \cdot x^{(j)}_{\text{noise}}$ (using Lemma 8). $\qquad \square$

<u>**Bounding $\delta_{\text{maj-train}}\left(\hat{w}, \gamma^2\right)$ in terms of $\|\hat{w}\|$**</u>

**Lemma 2.** *For a separator $\hat{w}$ with bounded $\alpha^{(i)}(\hat{w})^2 \leq \frac{10n}{\sigma_{\text{noise}}^2}$ for all $i = 1, \ldots, n$, its norm can be bounded with high probability as*

$$\|\hat{w}\|_2^2 \geq \frac{\gamma^4(1-c_1)}{\sigma_{\text{noise}}^2} \delta_{\text{maj-train}}\left(\hat{w}, \gamma^2\right) n_{\text{maj}} - \frac{10}{\sigma_{\text{noise}}^2 n^3} \tag{26}$$

*under the parameter settings of Theorem 2.*

*Proof.* The result follows bounded norms (Lemma 9), bounded dot products (Lemma 8), and the definition of $\delta_{\text{maj-train}}\left(\hat{w}, \gamma^2\right)$ (Definition 3).

$$\|\hat{w}\|_2^2 \geq \sum_{i \in G_{\text{maj}}} \alpha^{(i)}(\hat{w})^2 \|x_{\text{noise}}^{(i)}\|_2^2 + \sum_{j \neq k} \alpha^{(j)}(\hat{w}) \alpha^{(k)}(\hat{w}) x_{\text{noise}}^{(j)} \cdot x_{\text{noise}}^{(k)} \tag{27}$$

$$\geq \underbrace{\left(\frac{\gamma^4(1-c_1)}{\sigma_{\text{noise}}^2}\right) \delta_{\text{maj-train}}\left(\hat{w}, \gamma^2\right) n_{\text{maj}}}_{\text{Choosing only points with } \alpha^{(i)}(\hat{w}) \geq \gamma^2/\sigma_{\text{noise}}^2} - \underbrace{\frac{M^2}{\sigma_{\text{noise}}^2 n^4}}_{\max \alpha^{(i)}(\hat{w}) = M/\sigma_{\text{noise}}^2} \quad , \text{ w.h.p.} \tag{28}$$

$$\geq \frac{\gamma^4(1-c_1)}{\sigma_{\text{noise}}^2} \delta_{\text{maj-train}}\left(\hat{w}, \gamma^2\right) n_{\text{maj}} - \frac{10}{\sigma_{\text{noise}}^2 n^3} \tag{29}$$

$\square$

<u>**Bounding $\delta_{\text{maj-train}}\left(\hat{w}^{\text{minnorm}}, \gamma^2\right)$**</u>

We now apply Lemma 1 and Lemma 2 in order to bound $\delta_{\text{maj-train}}\left(\hat{w}^{\text{minnorm}}, \gamma^2\right)$, showing that the fraction of majority points that are memorized is small for appropriate choice of $\gamma$.

To invoke Lemma 2, we first show that the coefficient $\alpha^{(i)}(\hat{w}^{\text{minnorm}})$ is bounded above with high probabiltity.

**Lemma 3.** *Under the parameter settings of Theorem 2, with high probability, $\alpha^{(i)}(\hat{w}^{\text{minnorm}})$ is bounded above for $i = 1, \ldots, n$ as*

$$\alpha^{(i)}(\hat{w}^{\text{minnorm}})^2 \leq \frac{10n}{\sigma_{\text{noise}}^4}. \tag{30}$$

*Proof.* Let $\max_i \alpha^{(i)}(\hat{w}^{\text{minnorm}}) = \frac{M}{\sigma_{\text{noise}}^2}$.

$$\|\hat{w}^{\text{minnorm}}\|_2^2 \geq \|\hat{w}_{\text{noise}}^{\text{minnorm}}\|_2^2 \tag{31}$$

$$= \sum_{i \in G_{\text{min}} G_{\text{maj}}} \alpha^{(i)}(\hat{w}^{\text{minnorm}})^2 \|x_{\text{noise}}^{(i)}\|_2^2 + \sum_{i,j} \alpha^{(i)}(\hat{w}^{\text{minnorm}}) \alpha^{(j)}(\hat{w}^{\text{minnorm}}) x_{\text{noise}}^{(i)} \cdot x_{\text{noise}}^{(j)} \tag{32}$$

$$\geq \frac{M^2(1-c_1)}{\sigma_{\text{noise}}^2} - \frac{M^2}{\sigma_{\text{noise}}^2 n^6} n^2 \tag{33}$$

$$\geq \frac{M^2(1-c_1)}{\sigma_{\text{noise}}^2} - \frac{M^2}{\sigma_{\text{noise}}^2 n^4}. \tag{34}$$

From the upper bound on $\|\hat{w}^{\mathsf{minnorm}}\|_2^2$ (Lemma 1), we have

$$\frac{M^2(1 - c_1)}{\sigma_{\mathsf{noise}}^2} - \frac{M^2}{\sigma_{\mathsf{noise}}^2 n^4} \leq u^2 + s^2 \sigma_{\mathsf{noise}}^2 (1 + c_1) n_{\mathsf{min}} + \frac{s^2 \sigma_{\mathsf{noise}}^2}{n^4} \tag{35}$$

$$\implies M^2 \left( 1 - c_1 - \frac{1}{n^4} \right) \leq u^2 \sigma_{\mathsf{noise}}^2 + (s\sigma_{\mathsf{noise}}^2)^2 \left( (1 + c_1) n_{\mathsf{min}} + \frac{1}{n^4} \right) \tag{36}$$

$$\implies M^2 \left( 1 - c_1 - \frac{1}{n^4} \right) \leq u^2 \frac{n_{\mathsf{maj}}}{360000} + (s\sigma_{\mathsf{noise}}^2)^2 \left( (1 + c_1) n_{\mathsf{min}} + \frac{1}{n^4} \right), \tag{37}$$

$$\text{From a bound on } \sigma_{\mathsf{noise}}^2 \text{ in the parameter settings.} \tag{38}$$

Since $c_1 < 1/2000$, and $n \geq 2000$, setting $u = 1.3125$, $s\sigma_{\mathsf{noise}}^2 = 2.61$, we get $M^2 \leq 10n$. $\qquad\square$

Now, we are ready to show that $\delta_{\mathsf{maj\text{-}train}}\left( \hat{w}^{\mathsf{minnorm}}, \gamma^2 \right)$ is small.

**Lemma 4.** *Under the parameter settings of Theorem 2, the following is true with high probability.*

$$\delta_{\mathit{maj\text{-}train}}\left( \hat{w}^{\mathsf{minnorm}}, \frac{9}{10} \right) \leq 1/200, \tag{39}$$

*Proof.* Applying Lemma 2 to $\hat{w}^{\mathsf{minnorm}}$ by invoking the bounds on $\alpha^{(i)}(\hat{w}^{\mathsf{minnorm}})$ (Lemma 3),

$$\|\hat{w}^{\mathsf{minnorm}}\|_2^2 \geq \frac{\gamma^4(1 - c_1)}{\sigma_{\mathsf{noise}}^2} \delta_{\mathsf{maj\text{-}train}}\left( \hat{w}^{\mathsf{minnorm}}, \gamma^2 \right) n_{\mathsf{maj}} - \frac{10}{\sigma_{\mathsf{noise}}^2 n^3} \tag{40}$$

with high probability. Putting this together with Lemma 1, we have

$$\frac{\gamma^4(1 - c_1)}{\sigma_{\mathsf{noise}}^2} \delta_{\mathsf{maj\text{-}train}}\left( \hat{w}^{\mathsf{minnorm}}, \gamma^2 \right) n_{\mathsf{maj}} - \frac{10}{\sigma_{\mathsf{noise}}^2 n^3} \leq u^2 + s^2 \sigma_{\mathsf{noise}}^2 (1 + c_1) n_{\mathsf{min}} + \frac{s^2 \sigma_{\mathsf{noise}}^2}{n^4}$$

$$\implies \delta_{\mathsf{maj\text{-}train}}\left( \hat{w}^{\mathsf{minnorm}}, \gamma^2 \right) \leq \underbrace{\frac{u^2 \sigma_{\mathsf{noise}}^2}{\gamma^4 n_{\mathsf{maj}}(1 - c_1)}}_{\text{Very small}} + \underbrace{\left( \frac{(s\sigma_{\mathsf{noise}}^2)^2(1 + c_1)}{\gamma^4(1 - c_1)} \right) \frac{n_{\mathsf{min}}}{n_{\mathsf{maj}}}}_{\approx 0.0042} + \underbrace{\frac{(s\sigma_{\mathsf{noise}}^2)^2}{n^4 n_{\mathsf{maj}}}}_{\text{Very small}} + \underbrace{\frac{10}{\gamma^4(1 - c_1)n^3}}_{\text{Very small}}$$

$$\implies \delta_{\mathsf{maj\text{-}train}}\left( \hat{w}^{\mathsf{minnorm}}, \frac{9}{10} \right) \leq 1/200, \text{ w.h.p,}$$

where in the last step we substitute the constants $\gamma^2 = 9/10$, $u = 1.3125$, $s\sigma_{\mathsf{noise}}^2 = 2.61$, $n_{\mathsf{maj}}/n_{\mathsf{min}} \leq 1/2000$ and $\sigma_{\mathsf{noise}}^2 \leq n_{\mathsf{maj}}/360000$. $\qquad\square$

### B.2.2. CONCENTRATION INEQUALITIES

**Lemma 5.** *With probability $> 1 - 1/100$, if $\sigma_{\mathsf{spu}} \leq \frac{1}{4\sqrt{\log 100n}}$,*

$$a - 1/5 \leq x_{\mathsf{spu}}^{(i)} \leq a + 1/5, \ \forall i = 1, \dots n, \tag{41}$$

*where $a$ is the spurious attribute.*

This follows from standard subgaussian concentration and union bound over $n = n_{\mathsf{maj}} + n_{\mathsf{min}}$ points.

**Lemma 6.** *For a vector $z \in \mathbb{R}^N$ such that $z \in \mathcal{N}(0, \sigma^2 I)$,*

$$\mathbb{P}(|\|z\|^2 - \sigma^2 N| \geq \sigma^2 t) \leq 2 \exp\left( \frac{-Nt^2}{8} \right). \tag{42}$$

**Lemma 7.** *For two vectors $z_i, z_j \in \mathbb{R}^N$ such that $z_i, z_j \sim \mathcal{N}(0, \sigma^2 I)$, by Hoeffding's inequality, we have*

$$\mathbb{P}(|z_i \cdot z_j| \geq \sigma^2 t) \leq 2 \exp\left( -\frac{t^2}{2\|z_i\|^2} \right). \tag{43}$$

**Corollary 1.** *Combining Lemma 6 and Lemma 7, we get*

$$\mathbb{P}(|z_i \cdot z_j| \geq \sigma^2 t) \leq 2\exp\left(\frac{-N^3}{8}\right) + 2\exp\left(-\frac{t^2}{8N}\right). \tag{44}$$

**Lemma 8.** *For $N = \Omega(poly(n))$, with probability greater than $1 - 1/2000$,*

$$|x_{\text{noise}}^{(i)} \cdot x_{\text{noise}}^{(j)}| \leq \frac{\sigma_{\text{noise}}^2}{n^6} \quad \forall x_{\text{noise}}^{(i)}, x_{\text{noise}}^{(j)}. \tag{45}$$

This follows from Corollary 1 and union bound over $n^2$ pairs of training points.

**Lemma 9.** *For $N = \Omega(poly(n))$, with probability greater than $1 - 1/2000$,*

$$(1 - c_1)\sigma^2 \leq \|x_{\text{noise}}^{(i)}\|^2 \leq (1 + c_1)\sigma^2, \forall i. \tag{46}$$

*This follows from Lemma 6 and union bound over $n$ training points. In particular, we can set $c_1 < 1/2000$ for large enough $N$.*

### B.2.3. SMALL $\delta_{\text{MAJ-TRAIN}}\left(\hat{w}^{\text{minnorm}}, \gamma^2\right)$ IMPLIES HIGH WORST-GROUP ERROR

In the previous section, we proved that $\delta_{\text{maj-train}}\left(\hat{w}^{\text{minnorm}}, \gamma^2\right)$, the fraction of majority training samples that can have coefficient on the noise vectors greater than $\gamma^2/\sigma_{\text{noise}}^2$ in the max margin separator $\hat{w}^{\text{minnorm}}$ is bounded for suitable value of $\gamma$. We showed this using the fact that the norm of $\hat{w}^{\text{minnorm}}$ is the smallest among all separators and the observation that the squared norm of a separator roughlty scales proportional the number of training points that have large coefficient along the noise vectors.

What does small $\delta_{\text{maj-train}}\left(\hat{w}^{\text{minnorm}}, \gamma^2\right)$ imply? We now show that the bound on $\delta_{\text{maj-train}}\left(\hat{w}^{\text{minnorm}}, \gamma^2\right)$ has an important consequence on the worst-group error $\text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}})$; low $\delta_{\text{maj-train}}\left(\hat{w}^{\text{minnorm}}, \gamma\right)$ would imply high worst-group error $\text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}})$. We show that there is a trade-off between the worst-group test error of a separator and the fraction of *majority points* that it "memorizes" at training time. If a model that has low worst-group test error must use the core feature and not the spurious feature, and to obtain zero training error such a model would memorize a potentially large fraction of majority and minority points. In contrast, if the model instead uses only the spurious feature, then the worst-group test error would be high, but it would memorize only a small fraction of majority examples at training time; because we assume that the spurious feature is much less noisy than the core feature ($\sigma_{\text{core}} \gg \sigma_{\text{spu}}$), much fewer majority examples would need to be memorized. To summarize, *a large $\hat{w}_{\text{spu}}$ would require smaller fraction of majority points to be memorized $\delta_{\text{maj-train}}\left(\hat{w}, \gamma^2\right)$ but increase the worst-group test error $\text{Err}_{\text{wg}}(\hat{w})$*. We formalize the above trade-off between the worst-group error and fraction of majority examples to be memorized in Proposition 3.

**Proposition 3.** *For the minimum norm separator $\hat{w}^{\text{minnorm}}$, under the parameter settings of Theorem 2, with high probability,*

$$Err_{\text{wg}}(\hat{w}^{\text{minnorm}}) \geq \Phi\left(\frac{-c_3 + \hat{w}_{\text{spu}}^{\text{minnorm}} - \hat{w}_{\text{core}}^{\text{minnorm}}}{\sqrt{\hat{w}_{\text{core}}^{\text{minnorm}\,2}\sigma_{\text{core}}^2 + \hat{w}_{\text{spu}}^{\text{minnorm}\,2}\sigma_{\text{spu}}^2}}\right) - c_4, \tag{47}$$

*for some constants $c_3, c_4 < 1/1000$ and $\Phi$ the Gaussian CDF.*

*For any separator $\hat{w}$ that spans the training points and satisfies*

$$\alpha^{(i)}(\hat{w})^2 \leq \frac{10n}{\sigma_{\text{noise}}^4}, \tag{48}$$

*under the parameter settings of Theorem 2, with high probability,*

$$\delta_{\text{maj-train}}\left(\hat{w}, \gamma^2\right) \geq \Phi\left(\frac{1 - (1 + c_1)\gamma^2 - c_5 - \hat{w}_{\text{spu}} - \hat{w}_{\text{core}}}{\sqrt{\hat{w}_{\text{core}}^2\sigma_{\text{core}}^2 + \hat{w}_{\text{spu}}^2\sigma_{\text{spu}}^2}}\right) - c_6, \tag{49}$$

*for some constants $c_1 < 1/2000; c_5, c_6 < 1/1000$ and $\Phi$ the Gaussian CDF.*

We prove Proposition 3 in Section B.2.5.

As mentioned before, we see that the spurious component weight $\hat{w}_{\text{spu}}^{\text{minnorm}}$ has opposite effects on the two quantities; $\text{Err}_{\text{wg}}(\hat{w})$ increases with increase $\hat{w}_{\text{spu}}$, but $\delta_{\text{maj-train}}(\hat{w}, \gamma)$ decreases with increase in $\hat{w}_{\text{spu}}$. This dependence can be exploited to relate the two quantities to each other as follows.

$$\Phi^{-1}\left(\delta_{\text{maj-train}}\left(\hat{w}^{\text{minnorm}}, \gamma\right) + c_6\right) + \Phi^{-1}(\text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}}) + c_4) \geq \frac{1 - c_3 - c_5 - (1 + c_1)\gamma^2 - 2\hat{w}_{\text{core}}^{\text{minnorm}}}{\sqrt{\hat{w}_{\text{core}}^2 \sigma_{\text{core}}^2 + \hat{w}_{\text{spu}}^2 \sigma_{\text{spu}}^2}}. \tag{50}$$

In other words, if the $\delta_{\text{maj-train}}\left(\hat{w}^{\text{minnorm}}, \gamma\right)$ is low, then $\text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}})$ would need to be high.

### B.2.4. WORST-GROUP ERROR IS HIGH

Recall from part 1 that $\delta_{\text{maj-train}}\left(\hat{w}^{\text{minnorm}}, \gamma\right) < 1/200$ for appropriate choice of $\gamma$, and from part 2 the trade-off between $\delta_{\text{maj-train}}\left(\hat{w}^{\text{minnorm}}, \gamma\right)$ and $\text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}})$ (Equation (50)). As a final step, we need to bound the quantities on the RHS of Equation (50). All the constants are small, and $\gamma^2 = 9/10$, $\delta_{\text{maj-train}}\left(\hat{w}^{\text{minnorm}}, 9/10\right) \leq 1/200$ (Lemma 4) which allows us to write

$$\Phi^{-1}(0.006) + \Phi^{-1}(\text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}}) + c_4) \geq \frac{-2\hat{w}_{\text{core}}^{\text{minnorm}}}{\sqrt{\hat{w}_{\text{core}}^{\text{minnorm}\,2}\sigma_{\text{core}}^2 + \hat{w}_{\text{spu}}^{\text{minnorm}\,2}\sigma_{\text{spu}}^2}} \geq \frac{-2}{\sigma_{\text{core}}} \tag{51}$$

$$\implies \Phi^{-1}(\text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}}) + c_4) \geq 0.512 \tag{52}$$

$$\implies \text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}}) \geq 0.67 \tag{53}$$

We have hence proved that the minimum-norm separator $\hat{w}^{\text{minnorm}}$ incurs high worst-group error with high probability under the specified conditions.

### B.2.5. PROOF OF PROPOSITION 3

**Proposition 3.** *For the minimum norm separator $\hat{w}^{\text{minnorm}}$, under the parameter settings of Theorem 2, with high probability,*

$$Err_{\text{wg}}(\hat{w}^{\text{minnorm}}) \geq \Phi\left(\frac{-c_3 + \hat{w}_{\text{spu}}^{\text{minnorm}} - \hat{w}_{\text{core}}^{\text{minnorm}}}{\sqrt{\hat{w}_{\text{core}}^{\text{minnorm}\,2}\sigma_{\text{core}}^2 + \hat{w}_{\text{spu}}^{\text{minnorm}\,2}\sigma_{\text{spu}}^2}}\right) - c_4, \tag{47}$$

*for some constants $c_3, c_4 < 1/1000$ and $\Phi$ the Gaussian CDF.*

*For any separator $\hat{w}$ that spans the training points and satisfies*

$$\alpha^{(i)}(\hat{w})^2 \leq \frac{10n}{\sigma_{\text{noise}}^4}, \tag{48}$$

*under the parameter settings of Theorem 2, with high probability,*

$$\delta_{\text{maj-train}}\left(\hat{w}, \gamma^2\right) \geq \Phi\left(\frac{1 - (1 + c_1)\gamma^2 - c_5 - \hat{w}_{\text{spu}} - \hat{w}_{\text{core}}}{\sqrt{\hat{w}_{\text{core}}^2 \sigma_{\text{core}}^2 + \hat{w}_{\text{spu}}^2 \sigma_{\text{spu}}^2}}\right) - c_6, \tag{49}$$

*for some constants $c_1 < 1/2000$; $c_5, c_6 < 1/1000$ and $\Phi$ the Gaussian CDF.*

*Proof.* We derive the two bounds below.

#### Worst-group test error

We bound the expected worst-group error $\text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}})$, which is the expected worst-group loss over the data distribution. Below, we lower bound the worst-group error $\text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}})$ by bounding the error on a particular group: minority positive

points which have label $y = 1$ and spurious attribute $a = -1$. The test error is the probability that a test example $x$ from this group gets misclassified, i.e. $\hat{w}^{\text{minnorm}} \cdot x < 0$.

$$\text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}}) \geq \mathbb{P}\left(\hat{w}^{\text{minnorm}} \cdot x < 0 \mid y = 1, a = -1\right) \tag{54}$$

$$= \mathbb{P}\left(\hat{w}_{\text{core}}^{\text{minnorm}} x_{\text{core}} + \hat{w}_{\text{spu}}^{\text{minnorm}} x_{\text{spu}} + \hat{w}_{\text{noise}}^{\text{minnorm}} \cdot x_{\text{noise}} < 0 \mid y = 1, a = -1\right) \tag{55}$$

$$= \mathbb{P}\left(\hat{w}_{\text{core}}^{\text{minnorm}}(1 + \sigma_{\text{core}} z_1) + \hat{w}_{\text{spu}}^{\text{minnorm}}(-1 + \sigma_{\text{spu}} z_2) + \hat{w}_{\text{noise}}^{\text{minnorm}} \cdot x_{\text{noise}} < 0\right) \tag{56}$$

In the last step, we rewrite for convenience $x_{\text{core}} = y + \sigma_{\text{core}} z_1$ and $x_{\text{spu}} = a + \sigma_{\text{spu}} z_2$, where $z_1, z_2 \sim \mathcal{N}(0, 1)$.

We use the properties of high-dimensional Gaussian random vectors to bound the quantity $\hat{w}_{\text{noise}}^{\text{minnorm}} \cdot x_{\text{noise}}$. Recall that $\hat{w}_{\text{noise}}^{\text{minnorm}}$ can be written as

$$\hat{w}_{\text{noise}}^{\text{minnorm}} = \sum_{i \in G_{\text{maj}}, G_{\text{min}}} \alpha^{(i)}(\hat{w}^{\text{minnorm}}) x_{\text{noise}}^{(i)}. \tag{57}$$

From Lemma 3, we know that $\max_i \alpha^{(i)}(\hat{w}^{\text{minnorm}})^2 < \frac{10n}{\sigma_{\text{noise}}^4}$. This, along with Lemma 7 gives $|x_{\text{noise}} \cdot \hat{w}_{\text{noise}}^{\text{minnorm}}| \leq c_3$ with probability $1 - c_4$ for some small constants $c_3, c_4 < 1/1000$. Let $B$ denote the event that this high probability event where the dot product $|x_{\text{noise}} \cdot \hat{w}_{\text{noise}}^{\text{minnorm}}| \leq c_3$. Using the fact that $\mathbb{P}(A) \geq \mathbb{P}(A \mid B) - \mathbb{P}(\neg B)$ which follows from simple algebra, we have

$$\text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}}) \geq \mathbb{P}\left(\hat{w}_{\text{core}}^{\text{minnorm}}(1 + \sigma_{\text{core}} z_1) + \hat{w}_{\text{spu}}^{\text{minnorm}}(-1 + \sigma_{\text{spu}} z_2) + \hat{w}_{\text{noise}}^{\text{minnorm}} \cdot x_{\text{noise}} < 0\right) \tag{58}$$

$$\geq \mathbb{P}\left(\hat{w}_{\text{core}}^{\text{minnorm}}(1 + \sigma_{\text{core}} z_1) + \hat{w}_{\text{spu}}^{\text{minnorm}}(1 - \sigma_{\text{spu}} z_2) < -c_3\right) - c_4 \tag{59}$$

$$= \mathbb{P}\left(\hat{w}_{\text{core}}^{\text{minnorm}} \sigma_{\text{core}} z_1 + \hat{w}_{\text{spu}}^{\text{minnorm}} \sigma_{\text{spu}} z_2 < -c_3 + \hat{w}_{\text{spu}}^{\text{minnorm}} - \hat{w}_{\text{core}}^{\text{minnorm}}\right) - c_4 \tag{60}$$

$$= \Phi\left(\frac{-c_3 + \hat{w}_{\text{spu}}^{\text{minnorm}} - \hat{w}_{\text{core}}^{\text{minnorm}}}{\sqrt{\hat{w}_{\text{core}}^{\text{minnorm}\,2} \sigma_{\text{core}}^2 + \hat{w}_{\text{spu}}^{\text{minnorm}\,2} \sigma_{\text{spu}}^2}}\right) - c_4. \tag{61}$$

From the expression above, we see that $\text{Err}_{\text{wg}}(\hat{w}^{\text{minnorm}})$ increases as the spurious component $\hat{w}_{\text{spu}}^{\text{minnorm}}$ increases. This is because in the minority group, the spurious feature is negatively correlated with the label.

**Fraction of memorized training examples in majority groups**

We now compute a lower bound on $\delta_{\text{maj-train}}\left(\hat{w}^{\text{minnorm}}, \gamma^2\right)$, which is the number of majority points (where $a = y$) that are "memorized." Intuitively, we want to show that the fraction depends on $\hat{w}_{\text{spu}} - \hat{w}_{\text{core}}$. The more the core feature is used relative to the spurious feature, the larger fraction of points need to be memorized because the core feature is more noisy.

First, consider a separator $\hat{w}$ with some core and spurious components $\hat{w}_{\text{core}}$ and $\hat{w}_{\text{spu}}$. Recall that $\hat{w}_{\text{noise}} = \sum_i \alpha^{(i)}(\hat{w}) x_{\text{noise}}^{(i)}$ and $y^{(i)}(\hat{w} \cdot x^{(i)}) \geq 1$ by the definition of separators. For a given $\hat{w}_{\text{core}}$ and $\hat{w}_{\text{spu}}$, we want to bound the fraction of majority points ($a = y$) which can have $\alpha^{(i)}(\hat{w}) < \frac{\gamma^2}{\sigma_{\text{noise}}^2}$. We focus only on separators with bounded memorization, i.e. those that satisfy $\alpha^{(i)}(\hat{w})^2 \leq \frac{10n}{\sigma_{\text{noise}}^4}$. Note that from Lemma 3, w.h.p., the mininum-norm separator $\hat{w}^{\text{minnorm}}$ satifies this condition.

We bound the above by bounding a related quantity: the fraction of points that are memorized in the training distribution in expectation. We then use concentration to relate it to the fraction of the training set.

Formally, we have fixed quantities $\hat{w}_{\text{core}}$ and $\hat{w}_{\text{spu}}$. The training set is generated as per the usual data generating distribution. As before, we are interested in separators on the training set. For any majority training point, the coefficient $\alpha^{(i)}(\hat{w})$ in a separator is a random variable. Since training point $i$ is separated, we have

$$\hat{w}_{\text{core}}(1 + \sigma_{\text{core}} z_1) + \hat{w}_{\text{spu}}(1 + \sigma_{\text{spu}} z_2) + \left(\sum_i \alpha^{(i)}(\hat{w}) x_{\text{noise}}^{(i)}\right)^\top x_{\text{noise}}^{(i)} \geq 1.$$

From Lemma 8, Lemma 6, and the condition on $\alpha^{(i)}(\hat{w})$, this implies with high probability that

$$\hat{w}_{\text{core}}(1 + \sigma_{\text{core}} z_1) + \hat{w}_{\text{spu}}(1 + \sigma_{\text{spu}} z_2) \geq 1 - (1 + c_1)\sigma_{\text{noise}}^2 \alpha^{(i)}(\hat{w}) - c_5,$$

for some constant $c_5 < 1/1000$. Conditioning on the high probability event just as before ($\mathbb{P}(A) \leq \mathbb{P}(A \mid B) + \mathbb{P}(\neg B)$), we get

$$\mathbb{P}(\alpha^{(i)}(\hat{w}) \leq \frac{\gamma^2}{\sigma_{\text{noise}}^2}) \leq \mathbb{P}\left(\hat{w}_{\text{core}}\sigma_{\text{core}}z_1 + \hat{w}_{\text{spu}}\sigma_{\text{spu}}z_2 \leq -1 + (1+c_1)\gamma^2 + c_5 + \hat{w}_{\text{core}} + \hat{w}_{\text{spu}}\right) + \delta \tag{62}$$

$$= \Phi\left(\frac{-1 + (1+c_1)\gamma^2 + c_5 + \hat{w}_{\text{spu}} + \hat{w}_{\text{core}}}{\sqrt{\hat{w}_{\text{core}}^2\sigma_{\text{core}}^2 + \hat{w}_{\text{spu}}^2\sigma_{\text{spu}}^2}}\right) + \delta \tag{63}$$

$$\implies \mathbb{P}(\alpha^{(i)}(\hat{w}) \geq \frac{\gamma^2}{\sigma_{\text{noise}}^2}) \geq \Phi\left(\frac{1 - (1+c_1)\gamma^2 - c_5 - \hat{w}_{\text{spu}} - \hat{w}_{\text{core}}}{\sqrt{\hat{w}_{\text{core}}^2\sigma_{\text{core}}^2 + \hat{w}_{\text{spu}}^2\sigma_{\text{spu}}^2}}\right) - \delta, \tag{64}$$

for some $\delta < 1/2000$. Finally, we connect to $\delta_{\text{maj-train}}(\hat{w})(\gamma^2)$ which is the finite sample version of the quantity $\mathbb{P}(\alpha^{(i)}(\hat{w}) \leq \frac{\gamma^2}{\sigma_{\text{noise}}^2})$. By DKW, we know that the empirical CDF converges to the population CDF. Under the conditions of Theorem 2, which lower bounds the number of majority elements, we have with high probability,

$$\delta_{\text{maj-train}}(\hat{w})(\gamma^2) \geq \Phi\left(\frac{1 - (1+c_1)\gamma^2 - c_5 - \hat{w}_{\text{spu}} - \hat{w}_{\text{core}}}{\sqrt{\hat{w}_{\text{core}}^2\sigma_{\text{core}}^2 + \hat{w}_{\text{spu}}^2\sigma_{\text{spu}}^2}}\right) - c_6, \tag{65}$$

for constants $c_5, c_6 < 1/1000$.

$\square$

$\square$

### B.2.6. PROOF OF PROPOSITION 1

**Proposition 1** (Norm of models using the spurious feature). *When $\sigma_{\text{core}}^2, \sigma_{\text{spu}}^2$ satisfy the conditions in Theorem 1, there exists $N_0$ such that for all $N > N_0$, with high probability, there exists a separator $w^{\text{use}-\text{spu}} \in \mathcal{W}^{\text{use}-\text{spu}}$ such that*

$$\|w^{\text{use}-\text{spu}}\|_2^2 \leq \gamma_1^2 + \left(\frac{\gamma_2 n_{\min}}{\sigma_{\text{noise}}^2}\right),$$

*for some constants $\gamma_1, \gamma_2 > 0$.*

*Proof.* The proposition follows directly from Lemma 1.

$$\|w^{\text{use}-\text{spu}}\|_2^2 \leq u^2 + s^2\sigma_{\text{noise}}^2(1+c_1)n_{\min} + \frac{s^2\sigma_{\text{noise}}^2}{n^4}$$
$$\leq u^2 + s^2\sigma_{\text{noise}}^2(2+c_1)n_{\min}.$$

The constant $\gamma_1 = u = 1.3125$ and $\gamma_2 = s\sigma_{\text{noise}}^2(2+c_1) = 2.61(2+c_1)$ for $c_1 < 1/2000$. $\square$

### B.2.7. PROOF OF PROPOSITION 2

**Proposition 2** (Norm of models using the core feature). *When $\sigma_{\text{core}}^2, \sigma_{\text{spu}}^2$ satisfy the conditions in Theorem 1 and $n_{\min} \geq 100$, there exists $N_0$ such that for all $N > N_0$, with high probability, all separators $w^{\text{use}-\text{core}} \in \mathcal{W}^{\text{use}-\text{core}}$ satisfy*

$$\|w^{\text{use}-\text{core}}\|_2^2 \geq \frac{\gamma_3 n}{\sigma_{\text{noise}}^2},$$

*for some constant $\gamma_3 > 0$.*

*Proof.* To bound the norm for all $w^{\mathsf{use-core}} \in \mathcal{W}^{\mathsf{use-core}}$, we provide a lower bound on the norm of the minimum-norm separator in the set $\mathcal{W}^{\mathsf{use-core}}$:

$$\bar{w}^{\mathsf{use-core}} \stackrel{\text{def}}{=} \arg \min_{w \in \mathcal{W}^{\mathsf{use-core}}} \|w\|^2. \tag{66}$$

We bound the $\|\bar{w}^{\mathsf{use-core}}\|$ in two steps:

1. We first provide a lower bound for $\|\bar{w}^{\mathsf{use-core}}\|$ in terms of the fraction of training points memorized $\delta_{\text{train}}\left(\bar{w}^{\mathsf{use-core}}, \gamma^2\right)$ (defined formally below) in Corollary 2.

2. We then provide a lower bound for $\delta_{\text{train}}\left(\bar{w}^{\mathsf{use-core}}, \gamma^2\right)$ in Corollary 3.

We first formally define $\delta_{\text{train}}\left(\hat{w}, \gamma^2\right)$.

**Definition 4.** *For a separator $\hat{w}$ on training data $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$, let $\delta_{\text{train}}\left(\hat{w}, \gamma^2\right)$ be the fraction of training examples that $\hat{w}$ $\gamma$-memorizes:*

$$\delta_{\text{train}}\left(\hat{w}, \gamma^2\right) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n \mathbb{I}\left[\left|\alpha^{(i)}(\hat{w})\right| > \frac{\gamma^2}{\sigma_{\mathsf{noise}}^2}\right] \tag{67}$$

## Bounding $\|\bar{w}^{\mathsf{use-core}}\|$ by $\delta_{\text{train}}\left(\bar{w}^{\mathsf{use-core}}, \gamma^2\right)$

**Lemma 10.** *For a separator $\hat{w}$ with bounded $\alpha^{(i)}(\hat{w})^2 \leq \frac{10n}{\sigma_{\mathsf{noise}}^2}$ for all $i = 1, \ldots, n$, its norm can be bounded with high probability as*

$$\|\hat{w}\|_2^2 \geq \frac{\gamma^4(1 - c_1)}{\sigma_{\mathsf{noise}}^2} \delta_{\text{train}}\left(\hat{w}, \gamma^2\right) n - \frac{10}{\sigma_{\mathsf{noise}}^2 n^3} \tag{68}$$

*Proof.* Similarly to the proof of Lemma 2, the result follows bounded norms (Lemma 9), bounded dot products (Lemma 8), and the definition of $\delta_{\text{train}}\left(\hat{w}, \gamma^2\right)$ (Definition 4).

$$\|\hat{w}\|_2^2 \geq \sum_{i \in G_{\mathsf{maj}}} \alpha^{(i)}(\hat{w})^2 \|x_{\mathsf{noise}}^{(i)}\|_2^2 + \sum_{j \neq k} \alpha^{(j)}(\hat{w})\alpha^{(k)}(\hat{w}) x_{\mathsf{noise}}^{(j)} \cdot x_{\mathsf{noise}}^{(k)} \tag{69}$$

$$\geq \underbrace{\left(\frac{\gamma^4(1 - c_1)}{\sigma_{\mathsf{noise}}^2}\right) \delta_{\text{train}}\left(\hat{w}, \gamma^2\right) n}_{\text{Choosing only points with } \alpha^{(i)}(\hat{w}) \geq \gamma^2/\sigma_{\mathsf{noise}}^2} - \underbrace{\frac{M^2}{\sigma_{\mathsf{noise}}^2 n^4}}_{\max \alpha^{(i)}(\hat{w}) = M/\sigma_{\mathsf{noise}}^2}, \text{ w.h.p.} \tag{70}$$

$$\geq \frac{\gamma^4(1 - c_1)}{\sigma_{\mathsf{noise}}^2} \delta_{\text{train}}\left(\hat{w}, \gamma^2\right) n - \frac{10}{\sigma_{\mathsf{noise}}^2 n^3} \tag{71}$$

$\square$

**Corollary 2.** *With high probability,*

$$\|\bar{w}^{\mathsf{use-core}}\|_2^2 \geq \frac{\gamma^4(1 - c_1)}{\sigma_{\mathsf{noise}}^2} \delta_{\text{maj-train}}\left(\bar{w}^{\mathsf{use-core}}, \gamma^2\right) n_{\mathsf{maj}} - \frac{10}{\sigma_{\mathsf{noise}}^2 n^3} \tag{72}$$

*Proof.* The result follows from applying Lemma 10 to $\bar{w}^{\mathsf{use-core}}$, invoking the bounds on any individual component $\alpha^{(i)}(\bar{w}^{\mathsf{use-core}})$ obtained below in Lemma 11. $\square$

Below, we bound $\alpha^{(i)}(\bar{w}^{\mathsf{use-core}})$, where $\alpha^{(i)}(\bar{w}^{\mathsf{use-core}})$ is the component of training point $i$ to the classifier $\bar{w}^{\mathsf{use-core}}$ via the representer theorem.

**Lemma 11.** *With high probability, $i = 1, \ldots, n$, $\alpha^{(i)}(\bar{w}^{\mathsf{use-core}})$ can be bounded as follows.*

$$\alpha^{(i)}(\bar{w}^{\mathsf{use-core}})^2 \leq \frac{10n}{\sigma_{\mathsf{noise}}^4}. \tag{73}$$

*Proof.* As a first step, we upper bound the norm of $\bar{w}^{\mathsf{use-core}}$ by the norm of another separator $w^{\mathsf{use-core}} \in \mathcal{W}^{\mathsf{use-core}}$, using the fact that $\bar{w}^{\mathsf{use-core}}$ is the minimum-norm separator in $\mathcal{W}^{\mathsf{use-core}}$. In particular, we construct a separator $w^{\mathsf{use-core}} \in \mathcal{W}^{\mathsf{use-core}}$ that "memorizes" all training points, of the following form:

$$w_{\mathsf{core}}^{\mathsf{use-core}} = 0$$
$$w_{\mathsf{spu}}^{\mathsf{use-core}} = 0$$
$$\alpha^{(i)}(w^{\mathsf{use-core}}) = y^{(i)}\alpha \text{ for all } i = 1, \ldots, n.$$

This is analogous to the construction of $w^{\mathsf{use-spu}} \in \mathcal{W}^{\mathsf{use-spu}}$ (Lemma 1), and similar calculations can be used to obtain a suitable value $\alpha$ to ensure that $w^{\mathsf{use-core}}$ is a separator with high probability. We provide it below for completeness. We show that the following condition is sufficient to satisfy the margin constraints $y^{(i)}w^{\mathsf{use-core}} \cdot x^{(i)} \geq 1$ for all $i = 1, \ldots, n$ with high probability:

$$\alpha\sigma_{\mathsf{noise}}^2 \geq \frac{1}{1 - c_1 - 1/n^5}. \tag{74}$$

for $c_1 < 1/2000$. We obtain the above condition by applying Lemma 8 and Lemma 9 to the margin condition.

$$w^{\mathsf{use-core}} \cdot x^{(i)} \geq 1 \tag{75}$$

$$\implies \alpha\|x_{\mathsf{noise}}^{(i)}\|^2 - \alpha\sum_{j \neq i}\left|x_{\mathsf{noise}}^{(i)} \cdot x_{\mathsf{noise}}^{(j)}\right| \geq 1 \tag{76}$$

$$\implies \alpha\sigma_{\mathsf{noise}}^2(1 - c_1) - \frac{\alpha\sigma_{\mathsf{noise}}^2}{n^5} \geq 1 \text{ with high probability} \tag{77}$$

Thus, we can construct $w^{\mathsf{use-core}}$ by setting some constant $\alpha\sigma_{\mathsf{noise}}^2 \leq 2$.

Now that we have constructed $w^{\mathsf{use-core}}$, we can bound the norm of the minimum norm separator $\bar{w}^{\mathsf{use-core}}$ by the norm of $w^{\mathsf{use-core}}$. The following is true with high probability,

$$\|\bar{w}^{\mathsf{use-core}}\|^2 \leq \|w_{\mathsf{noise}}^{\mathsf{use-core}}\|^2 \tag{78}$$

$$= \sum_{i=1}^{n}\alpha^2\|x_{\mathsf{noise}}^{(i)}\|^2 + \sum_{i \neq j}\alpha^2 x_{\mathsf{noise}}^{(i)} \cdot x_{\mathsf{noise}}^{(j)} \tag{79}$$

$$\leq \alpha^2\sigma_{\mathsf{noise}}^2(1 + c_1)n + \frac{\alpha^2\sigma_{\mathsf{noise}}^2}{n^4} \tag{80}$$

Finally, we bound $\alpha^{(i)}(\bar{w}^{\mathsf{use-core}})$ for all $i$ by bounding $\max_i \alpha^{(i)}(\bar{w}^{\mathsf{use-core}}) = \frac{M}{\sigma_{\mathsf{noise}}^2}$. As we showed in the proof of Lemma 3, following is true with high probability:

$$\|\bar{w}^{\mathsf{use-core}}\|_2^2 \geq \frac{M^2(1 - c_1)}{\sigma_{\mathsf{noise}}^2} - \frac{M^2}{\sigma_{\mathsf{noise}}^2 n^4}. \tag{81}$$

Combined with the upper bound on $\|\bar{w}^{\mathsf{use-core}}\|_2^2$ (Equation (80)), we have

$$\frac{M^2(1 - c_1)}{\sigma_{\mathsf{noise}}^2} - \frac{M^2}{\sigma_{\mathsf{noise}}^2 n^4} \leq \|\bar{w}^{\mathsf{use-core}}\| \leq \alpha^2\sigma_{\mathsf{noise}}^2(1 + c_1)n + \frac{\alpha^2\sigma_{\mathsf{noise}}^2}{n^4} \tag{82}$$

$$\implies M^2\left(1 - c_1 - \frac{1}{n^4}\right) \leq (\alpha\sigma_{\mathsf{noise}}^2)^2\left((1 + c_1)n + \frac{1}{n^4}\right). \tag{83}$$

Since $c_1 < 1/2000$, and $n \geq 2000$, setting $\alpha\sigma_{\mathsf{noise}}^2 = 2$ yields $M^2 \leq 10n$ with high probability. $\square$

<u>**Bounding $\delta_{\text{train}}\left(\bar{w}^{\text{use}-\text{core}}, \gamma^2\right)$**</u>

**Corollary 3.** *Under the parameter settings of Theorem 2, with high probability,*

$$\delta_{train}\left(\bar{w}^{\text{use}-\text{core}}, \gamma^2\right) \geq \Phi\left(\frac{1 - (1 + c_1)\gamma^2 - c_5 - \bar{w}_{\text{core}}^{\text{use}-\text{core}}}{\left|\bar{w}_{\text{core}}^{\text{use}-\text{core}}\sigma_{\text{core}}\right|}\right) - c_6, \tag{84}$$

*for some constants $c_1 < 1/2000$; $c_5, c_6 < 1/1000$ where $\Phi$ is the Gaussian CDF.*

*Proof.* The result follows from applying Proposition 3 (which computes a bound on the majority fraction of points that is $\gamma-$memorized) to $\bar{w}^{\text{use}-\text{core}}$, invoking Lemma 11, and plugging in $\bar{w}_{\text{spu}}^{\text{use}-\text{core}} = 0$. Note that when $\bar{w}_{\text{spu}}^{\text{use}-\text{core}} = 0$, $\delta_{\text{train}}\left(\bar{w}^{\text{use}-\text{core}}, \gamma^2\right) = \delta_{\text{maj-train}}\left(\bar{w}^{\text{use}-\text{core}}, \gamma^2\right)$. $\square$

Finally, the above bound on $\delta_{\text{train}}\left(\bar{w}^{\text{use}-\text{core}}, \gamma^2\right)$ translates to a bound on the norm $\|\bar{w}^{\text{use}-\text{core}}\|$ via simple algebra. For $\gamma$ that satisfies $1 - (1 + c_1)\gamma^2 - c_5 > 0$:

$$\delta_{\text{train}}\left(\bar{w}^{\text{use}-\text{core}}, \gamma^2\right) \geq \Phi\left(\frac{-1}{\sigma_{\text{core}}} + \frac{1 - (1 + c_1)\gamma^2 - c_5}{\left|\bar{w}_{\text{core}}^{\text{use}-\text{core}}\sigma_{\text{core}}\right|}\right) - c_6 \tag{85}$$

$$\geq \Phi\left(\frac{-1}{\sigma_{\text{core}}}\right) - c_6. \tag{86}$$

Plugging the above lower bound into the bound on $\|\bar{w}^{\text{use}-\text{core}}\|$ from Corollary 2, we have

$$\|\bar{w}^{\text{use}-\text{core}}\|_2^2 \geq \frac{\gamma^4(1 - c_1)}{\sigma_{\text{noise}}^2}\delta_{\text{train}}\left(\bar{w}^{\text{use}-\text{core}}, \gamma^2\right)n_{\text{maj}} - \frac{10}{\sigma_{\text{noise}}^2 n^3} \tag{87}$$

$$\geq \frac{n}{\sigma_{\text{noise}}^2}\left(\Phi\left(\frac{-1}{\sigma_{\text{core}}}\right) - c_6\right)\gamma^4(1 - c_1) - \frac{10}{\sigma_{\text{noise}}^2 n^3} \tag{88}$$

$$\geq \frac{n}{\sigma_{\text{noise}}^2}\underbrace{\left[\left(\Phi\left(\frac{-1}{\sigma_{\text{core}}}\right) - c_6\right)\gamma^4(1 - c_1) - c_7\right]}_{\text{set to } \gamma_3} \tag{89}$$

for some $c_7 < 1/1000$. $\square$

## B.3. Underparameterized regime

So far, we have studied the overparameterized regime for the data distribution described in Section 5. In the overparameterized setting, where the dimension of noise features $N$ is very large, logistic regression (both ERM and reweighted) leads to max-margin classifiers. We showed that for some setting of parameters $n_{\text{maj}}, n_{\text{min}}, \sigma_{\text{spu}}, \sigma_{\text{core}}$, the robust error of such max-margin classifiers can be $> 2/3$, worse than random guessing. How does the same reweighted logistic regression perform in the underparameterized regime? We focus on the setting where $N = 0$. In this setting, the data is two-dimensional, and w.h.p., the training data is not linearly separable unless $\sigma_{\text{core}} = 0$. Consequently, the learned model $\hat{w}^{\text{rw}}\mathbb{R}^2$ that minimizes the reweighted training loss is not generally a max-margin separator.

For intuition, consider the following two sets of models, which are analogous to what we considered in Equation 12 in the main text for the overparameterized regime:

$$\mathcal{W}^{\text{use}-\text{spu}} \overset{\text{def}}{=} \{w \in \mathbb{R}^2 \text{ such that } w_{\text{core}} = 0\}$$

$$\mathcal{W}^{\text{use}-\text{core}} \overset{\text{def}}{=} \{w \in \mathbb{R}^2 \text{ such that } w_{\text{spu}} = 0\}. \tag{90}$$

The first set $\mathcal{W}^{\text{use}-\text{spu}}$ comprises models that use the spurious feature but not the core feature, and the second set $\mathcal{W}^{\text{use}-\text{core}}$ comprises models that use the core feature but not the spurious feature. Models in $\mathcal{W}^{\text{use}-\text{spu}}$ that exclusively use $x_{\text{spu}}$ will have high training loss on the minorities since the minority points cannot be memorized. Due to upweighting the minorities, these models will have high reweighted training loss. On the other hand, models in $\mathcal{W}^{\text{use}-\text{core}}$ exclusively use the core

features that are informative for the label $y$ across all groups. Hence they obtain reasonable loss across all groups and have smaller reweighted training loss than models in $\mathcal{W}^{\text{use}-\text{spu}}$.

We will show in this section that the population minimizer of the reweighted loss is indeed in $\mathcal{W}^{\text{use}-\text{core}}$ and bound the asymptotic variance of the reweighted estimator, leading to the final result in Theorem 1. Our approach is to study the asypmtotic behavior of the reweighted estimator when the number of data points $n \gg d$.

**Data distribution.** We first recap the data generating distribution (described in Section 5). $x = [x_{\text{core}}, x_{\text{spu}}]$ where,

$$x_{\text{core}} \mid y \sim \mathcal{N}(y, \sigma_{\text{core}}^2), \quad x_{\text{spu}} \mid a \sim \mathcal{N}(a, \sigma_{\text{spu}}^2),$$

For $p_{\text{maj}}$ fraction of points, we have $a = y$ (majority points) and for $1 - p_{\text{maj}}$ fraction of points, we have $a = -y$ (minority points).

**Reweighted logistic loss.** Let $p_{\text{maj}}$ be the fraction of the majority group points and $(1 - p_{\text{maj}})$ be the fraction of minority points. In order to use standard results from the asymptotics of M-estimators, we rewrite the reweighted estimator (defined in Section 2) as the minimizer of the following loss over $n$ training points $[x_i, y_i]_{i=1}^n$.

$$\hat{w}^{\text{rw}} = \arg\min \frac{1}{n} \sum_{i=1}^n \ell_{\text{rw}}(x_i, y_i, w) \tag{91}$$

$$\ell_{\text{rw}}(x, y, w) = \frac{-1}{p_{\text{maj}}} \log\left(\frac{1}{1 + \exp(-yw^\top x)}\right), \quad \text{For } (x, y) \text{ from majority group} \tag{92}$$

$$\ell_{\text{rw}}(x, y, w) = \frac{-1}{1 - p_{\text{maj}}} \log\left(\frac{1}{1 + \exp(-yw^\top x)}\right), \quad \text{For } (x, y) \text{ from minority group.} \tag{93}$$

We follow the standard steps of asymptotic analysis where we:

1. Compute the population minimizer $w^\star$ that satisfies $\nabla L_{\text{rw}}(w^\star) = 0$, where $L_{\text{rw}}(w^\star) = \mathbb{E}[\ell_{\text{rw}}(x, y, w^\star)]$.

2. Bound the asymptotic variance $\nabla^2 L_{\text{rw}}(w^\star)^{-1} \text{Cov}[\nabla \ell_{\text{rw}}(x, y, w^\star)] \nabla^2 L_{\text{rw}}(w^\star)^{-1}$.

**Proposition 4.** *For the data distribution under study, the population minimizer $w^\star$ that satisfies $\nabla L_{\text{rw}}(w^\star) = 0$ is the following.*

$$w^\star = \left[\frac{2}{\sigma_{\text{core}}^2}, 0\right]. \tag{94}$$

This is a very important property in the underparameterized regime: the population minimizer has the best possible worst-group error by only using the core feature and not the spurious feature.

**Proposition 5.** *The asymptotic distribution of the reweighted logistic regression estimator is as follows.*

$$\sqrt{n}(\hat{w} - w^\star) \to^d \mathcal{N}(0, V), \tag{95}$$

$$V \preceq \text{diag}\left(\frac{16 \exp\left(\frac{8}{(\sigma_{\text{core}}^2 + 8)\sigma_{\text{core}}^2}\right)(\sigma_{\text{core}}^2 + 1)(1 + 8/\sigma_{\text{core}}^2)^3}{p_{\text{maj}}(1 - p_{\text{maj}})(\sigma_{\text{core}}^2 + 9)^2}, \frac{16 \exp\left(\frac{8}{(\sigma_{\text{core}}^2 + 8)\sigma_{\text{core}}^2}\right)(1 + 8/\sigma_{\text{core}}^2)}{p_{\text{maj}}(1 - p_{\text{maj}})(\sigma_{\text{spu}}^2 + 1)}\right). \tag{96}$$

*For $\sigma_{\text{core}} \geq 1$, we have*

$$V \preceq \text{diag}\left(\frac{C_1}{p_{\text{maj}}(1 - p_{\text{maj}})}, \frac{C_2}{p_{\text{maj}}(1 - p_{\text{maj}})}\right), \tag{97}$$

*for some constants $C_1, C_2$.*

We see that the asymptotic variance increases as $p_{\text{maj}}$ increases. This is expected because the reweighted estimator upweights the minority points by inverse of group size. As these weights increase, the variance also increases. However, as we noted before, since the population minimizer has small worst-group error, for large enough training set size, we get small worst-group error since the asymptotic variance is finite (for fixed $p_{\text{maj}}$) and the estimator approaches the population minimizer.

We now prove Theorem 1 for the underparameterized regime, restated as Theorem 3 below.

**Theorem 3.** *In the underparameterized regime with $N = 0$, for $p_{\text{maj}} = \left(1 - \frac{1}{2001}\right)$, $\sigma^2_{\text{core}} = 1$, and $\sigma^2_{\text{spu}} = 0$, in the asymptotic regime with $n_{\text{maj}}, n_{\text{min}} \to \infty$, we have*

$$Err_{\text{wg}}(\hat{w}^{\text{rw}}) < 1/4. \tag{98}$$

*Proof.* We now put the two Propositions 5 and 4 together. We have $\hat{w}^{\text{rw}}_{\text{core}} \geq 2 - \epsilon_1$ and $|\hat{w}^{\text{rw}}_{\text{spu}}| \leq \epsilon_2$ for $\epsilon_1, \epsilon_2 < 1/10$, i.e the estimator is very close to the population minimizer. This follows from setting $\sigma_{\text{core}}, \sigma_{\text{spu}}, p_{\text{maj}} = \frac{n_{\text{maj}}}{n_{\text{maj}} + n_{\text{min}}}$ to their corresponding values and setting $n = n_{\text{maj}} + n_{\text{min}}$ to be large enough. In order to compute the worst-group error, WLOG consider points with label $y = 1$ (labels are balanced in the population). For a point from the majority group, the probability of misclassification is as follows.

$$\Pr[\hat{w}^{\text{rw}}_{\text{core}} x_{\text{core}} + \hat{w}^{\text{rw}}_{\text{spu}} x_{\text{spu}} \geq 0] = \Pr[z \geq \frac{\hat{w}^{\text{rw}}_{\text{core}} + \hat{w}^{\text{rw}}_{\text{spu}}}{\sigma^2_{\text{core}} \hat{w}^{\text{rw}}_{\text{core}}{}^2 + \sigma^2_{\text{spu}} \hat{w}^{\text{rw}}_{\text{spu}}{}^2}], \tag{99}$$

where $z \sim \mathcal{N}(0,1)$.

Similarly, for the minority group, the probability of misclassification is

$$\Pr[z \geq \frac{\hat{w}^{\text{rw}}_{\text{core}} - \hat{w}^{\text{rw}}_{\text{spu}}}{\sigma^2_{\text{core}} \hat{w}^{\text{rw}}_{\text{core}}{}^2 + \sigma^2_{\text{spu}} \hat{w}^{\text{rw}}_{\text{spu}}{}^2}], \text{ where } z \sim \mathcal{N}(0,1). \tag{100}$$

Therefore, the worst-group error of $\hat{w}^{\text{rw}}$ can be bounded as.

$$\text{Err}_{\text{wg}}(\hat{w}^{\text{rw}}) \leq 1 - \Phi\left(\frac{\hat{w}^{\text{rw}}_{\text{core}} - |\hat{w}^{\text{rw}}_{\text{spu}}|}{\sigma^2_{\text{core}} \hat{w}^{\text{rw}}_{\text{core}}{}^2 + \sigma^2_{\text{spu}} \hat{w}^{\text{rw}}_{\text{spu}}{}^2}\right), \tag{101}$$

where $\Phi$ is the Gaussian CDF. Substituting $\sigma_{\text{core}} = 1, \sigma_{\text{spu}} = 0, \hat{w}^{\text{rw}}_{\text{core}} \geq 2 - \epsilon_1, |\hat{w}^{\text{rw}}_{\text{spu}}| \leq \epsilon_2$ gives the required result that $\text{Err}_{\text{wg}}(\hat{w}^{\text{rw}}) < 1/4$. In contrast, in the overparameterized regime where $N \gg n$, even for very large $n$, the reweighted estimator has high worst-group error, as shown in Theorem 1. $\qquad \square$

### B.3.1. COMPLETE PROOFS

We now provide the proofs for Proposition 4 and Proposition 5 which mostly follow from straightforward algebra.

**Proposition 4.** *For the data distribution under study, the population minimizer $w^\star$ that satisfies $\nabla L_{\text{rw}}(w^\star) = 0$ is the following.*

$$w^\star = \left[\frac{2}{\sigma^2_{\text{core}}}, 0\right]. \tag{94}$$

*Proof.* For convenience, we compute expectations over the majority and minority groups separately and express the population loss $L_{\text{rw}}$ as the weighted sum of the two terms. Recall that we denote $x = [x_{\text{core}}, x_{\text{spu}}]$.

$$L_{\text{rw}}(w) = p_{\text{maj}} L_{\text{rw-maj}} + (1 - p_{\text{maj}}) L_{\text{rw-min}} \tag{102}$$

$$L_{\text{rw-maj}}(w) = \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma^2_{\text{core}})} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(y, \sigma^2_{\text{spu}})}[\ell_{\text{rw}}(x, y, w)]. \tag{103}$$

$$L_{\text{rw-min}}(w) = \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma^2_{\text{core}})} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(-y, \sigma^2_{\text{spu}})}[\ell_{\text{rw}}(x, y, w)]. \tag{104}$$

We use the following expression for computing the population gradient.

$$\nabla \log\left(\frac{1}{1 + \exp(-yw^\top x)}\right) = \left(\frac{-y \exp(-yw^\top x)}{1 + \exp(-yw^\top x)}\right)x. \tag{105}$$

Combining the definition of the reweighted loss and population losses (Equation 91 and Equation 102) with the gradient expression above gives the following.

$$\nabla \mathsf{L}_{\text{rw-maj}}(w) = \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma_{\text{core}}^2)} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(y, \sigma_{\text{spu}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{-y \exp(-y w^\top x)}{1 + \exp(-y w^\top x)} \right) x \right]. \tag{106}$$

$$\nabla \mathsf{L}_{\text{rw-min}}(w) = \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma_{\text{core}}^2)} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(-y, \sigma_{\text{spu}}^2)} \left[ \frac{1}{1 - p_{\text{maj}}} \left( \frac{-y \exp(-y w^\top x)}{1 + \exp(-y w^\top x)} \right) x \right]. \tag{107}$$

Now we compute $\nabla \mathsf{L}_{\text{rw}}(w^\star) = p_{\text{maj}} \nabla \mathsf{L}_{\text{rw-maj}}(w^\star) + (1 - p_{\text{maj}}) \nabla \mathsf{L}_{\text{rw-min}}(w^\star)$. First we compute wrt the spurious attribute $\nabla_{\text{spu}} \mathsf{L}_{\text{rw}}(w^\star)$. For convenience, let $c = \frac{2}{\sigma_{\text{core}}^2}$.

$$
\begin{aligned}
\nabla_{\text{spu}} \mathsf{L}_{\text{rw-maj}}(w^\star) &= \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma_{\text{core}}^2)} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(y, \sigma_{\text{spu}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{-y \exp(-y c x_{\text{core}})}{1 + \exp(-y c x_{\text{core}})} \right) x_{\text{spu}} \right] \\
&= \frac{1}{2} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(1, \sigma_{\text{spu}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{- \exp(-c x_{\text{core}})}{1 + \exp(-c x_{\text{core}})} \right) x_{\text{spu}} \right] \\
&\quad + \frac{1}{2} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(-1, \sigma_{\text{core}}^2)} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(-1, \sigma_{\text{spu}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{\exp(c x_{\text{core}})}{1 + \exp(c x_{\text{core}})} \right) x_{\text{spu}} \right] \\
&= \frac{1}{2} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{- \exp(-c x_{\text{core}})}{1 + \exp(-c x_{\text{core}})} \right) \right] - \frac{1}{2} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(-1, \sigma_{\text{core}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{\exp(c x_{\text{core}})}{1 + \exp(c x_{\text{core}})} \right) \right] \\
&= \frac{1}{2} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{- \exp(-c x_{\text{core}})}{1 + \exp(-c x_{\text{core}})} \right) \right] \underbrace{- \frac{1}{2} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{\exp(-c x_{\text{core}})}{1 + \exp(-c x_{\text{core}})} \right) \right]}_{\text{Replacing } x_{\text{core}} \sim \mathcal{N}(-1, \sigma_{\text{core}}^2) \text{ with } -x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \\
&= \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{- \exp(-c x_{\text{core}})}{1 + \exp(-c x_{\text{core}})} \right) \right] \\
\nabla_{\text{spu}} \mathsf{L}_{\text{rw-min}}(w^\star) &= \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma_{\text{core}}^2)} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(-y, \sigma_{\text{spu}}^2)} \left[ \frac{1}{1 - p_{\text{maj}}} \left( \frac{-y \exp(-y c x_{\text{core}})}{1 + \exp(-y c x_{\text{core}})} \right) x_{\text{spu}} \right] \\
&= \frac{1}{2} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{\exp(-c x_{\text{core}})}{1 + \exp(-c x_{\text{core}})} \right) \right] + \frac{1}{2} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(-1, \sigma_{\text{core}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{\exp(c x_{\text{core}})}{1 + \exp(c x_{\text{core}})} \right) \right] \\
&= \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \left[ \frac{1}{1 - p_{\text{maj}}} \left( \frac{\exp(-c x_{\text{core}})}{1 + \exp(-c x_{\text{core}})} \right) \right]
\end{aligned}
$$

Now we take the weighted combination of $\nabla_{\text{spu}} \mathsf{L}_{\text{rw-maj}}(w^\star)$ and $\nabla_{\text{spu}} \mathsf{L}_{\text{rw-min}}(w^\star)$, based on the fraction of the majority and minority samples in the population, which makes the two terms cancel out.

$$\nabla_{\text{spu}} \mathsf{L}_{\text{rw}} = p_{\text{maj}} \nabla_{\text{spu}} \mathsf{L}_{\text{rw-maj}}(w^\star) + (1 - p_{\text{maj}}) \nabla_{\text{spu}} \mathsf{L}_{\text{rw-min}}(w^\star) = 0. \tag{108}$$

Now we compute $\nabla_{\text{core}} L_{\text{rw}}(w^\star)$.

$$
\begin{aligned}
\nabla_{\text{core}} L_{\text{rw-maj}}(w^\star) &= \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma_{\text{core}}^2)} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(y, \sigma_{\text{spu}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{-y \exp(-y c x_{\text{core}})}{1 + \exp(-y c x_{\text{core}})} \right) x_{\text{core}} \right] \\
&= \frac{1}{2} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{-\exp(-c x_{\text{core}})}{1 + \exp(-c x_{\text{core}})} \right) x_{\text{core}} \right] + \frac{1}{2} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(-1, \sigma_{\text{core}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{\exp(c x_{\text{core}})}{1 + \exp(c x_{\text{core}})} \right) x_{\text{core}} \right] \\
&= \frac{1}{2} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{-\exp(-c x_{\text{core}})}{1 + \exp(-c x_{\text{core}})} \right) x_{\text{core}} \right] + \frac{1}{2} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(-1, \sigma_{\text{core}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{1}{1 + \exp(-c x_{\text{core}})} \right) x_{\text{core}} \right] \\
&= \frac{1}{2 p_{\text{maj}}} \frac{1}{\sigma_{\text{core}} \sqrt{2\pi}} \int_{-\infty}^{\infty} \frac{\exp(-c x_{\text{core}}) \exp\left( \frac{-(x-1)^2}{2\sigma_{\text{core}}^2} \right) - \exp\left( \frac{-(x+1)^2}{2\sigma_{\text{core}}^2} \right)}{1 + \exp(-c x_{\text{core}})} x_{\text{core}} \, dx_{\text{core}} \\
&= \frac{1}{2 p_{\text{maj}}} \frac{1}{\sigma_{\text{core}} \sqrt{2\pi}} \int_{-\infty}^{\infty} 0 \, dx_{\text{core}}, \text{ Substituting } c = \frac{2}{\sigma_{\text{core}}^2} \\
&= 0.
\end{aligned}
$$

Similarly, we get $\nabla_{\text{core}} L_{\text{rw-min}}(w^\star) = 0$ and hence proved that $\nabla_{\text{core}} L_{\text{rw}}(w^\star) = 0$. $\qquad \square$

**Lemma 12.** *The following is true.*

$$
\text{Cov}[\nabla \ell_{\text{rw}}(x, y, w^\star)] \preceq \text{diag} \left( \frac{\sigma_{\text{core}}^2 + 1}{p_{\text{maj}}(1 - p_{\text{maj}})}, \frac{\sigma_{\text{spu}}^2 + 1}{p_{\text{maj}}(1 - p_{\text{maj}})} \right). \tag{109}
$$

We now compute the asymptotic variance which involves computing $\nabla^2 L(w^\star)$ and $\text{Cov}[\nabla \ell_{\text{rw}}(w^\star)]$.

*Proof.* First, we show that the off-diagonal entries of $\text{Cov}[\ell_{\text{rw}}(x, y, w^\star)]$ are zero.

$$
\begin{aligned}
&\mathbb{E}[\nabla_{\text{core}} \ell_{\text{rw}}(x, y, w^\star) \nabla_{\text{spu}} \ell_{\text{rw}}(x, y, w^\star)] - \mathbb{E}[\nabla_{\text{core}} \ell_{\text{rw}}(x, y, w^\star)] \mathbb{E}[\nabla_{\text{spu}} \ell_{\text{rw}}(x, y, w^\star)] \\
&= \mathbb{E}[\nabla_{\text{core}} \ell_{\text{rw}}(x, y, w^\star) \nabla_{\text{spu}} \ell_{\text{rw}}(x, y, w^\star)] \\
&= p_{\text{maj}} \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma_{\text{core}}^2)} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(y, \sigma_{\text{spu}}^2)} \left[ \frac{1}{p_{\text{maj}}^2} \left( \frac{-y \exp(-y c x_{\text{core}})}{1 + \exp(-y c x_{\text{core}})} \right)^2 x_{\text{core}} x_{\text{spu}} \right] \\
&\quad + (1 - p_{\text{maj}}) \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma_{\text{core}}^2)} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(-y, \sigma_{\text{spu}}^2)} \left[ \frac{1}{(1 - p_{\text{maj}})^2} \left( \frac{-y \exp(-y c x_{\text{core}})}{1 + \exp(-y c x_{\text{core}})} \right)^2 x_{\text{core}} x_{\text{spu}} \right] \\
&= \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma_{\text{core}}^2)} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{-y \exp(-y c x_{\text{core}})}{1 + \exp(-y c x_{\text{core}})} \right)^2 y \right] \\
&\quad - \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma_{\text{core}}^2)} \left[ \frac{1}{1 - p_{\text{maj}}} \left( \frac{-y \exp(-y c x_{\text{core}})}{1 + \exp(-y c x_{\text{core}})} \right)^2 y \right] \\
&= \frac{1 - 2 p_{\text{maj}}}{2 p_{\text{maj}}(1 - p_{\text{maj}})} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \left[ \left( \frac{\exp(-c x_{\text{core}})}{1 + \exp(-c x_{\text{core}})} \right)^2 \right] - \frac{1 - 2 p_{\text{maj}}}{2 p_{\text{maj}}(1 - p_{\text{maj}})} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(-1, \sigma_{\text{core}}^2)} \left[ \left( \frac{\exp(c x_{\text{core}})}{1 + \exp(c x_{\text{core}})} \right)^2 \right] \\
&= \frac{1 - 2 p_{\text{maj}}}{2 p_{\text{maj}}(1 - p_{\text{maj}})} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \left[ \left( \frac{\exp(-c x_{\text{core}})}{1 + \exp(-c x_{\text{core}})} \right)^2 \right] - \frac{1 - 2 p_{\text{maj}}}{2 p_{\text{maj}}(1 - p_{\text{maj}})} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \left[ \left( \frac{\exp(-c x_{\text{core}})}{1 + \exp(-c x_{\text{core}})} \right)^2 \right] = 0.
\end{aligned}
$$

Now, we bound the diagonal elements.

$$\mathbb{E}[\nabla_{\mathsf{core}}(\ell_{\mathsf{rw}}(x, y, w^\star))^2] - (\mathbb{E}[\nabla_{\mathsf{core}}\ell_{\mathsf{rw}}(x, y, w^\star)])^2$$
$$= \mathbb{E}[\nabla_{\mathsf{core}}(\ell_{\mathsf{rw}}(x, y, w^\star))^2]$$
$$= p_{\mathsf{maj}}\mathbb{E}_y\mathbb{E}_{x_{\mathsf{core}}\sim\mathcal{N}(y, \sigma_{\mathsf{core}}^2)}\left[\frac{1}{p_{\mathsf{maj}}^2}\left(\frac{-y\exp(-ycx_{\mathsf{core}})}{1 + \exp(-ycx_{\mathsf{core}})}\right)^2 x_{\mathsf{core}}^2\right]$$
$$+ (1 - p_{\mathsf{maj}})\mathbb{E}_y\mathbb{E}_{x_{\mathsf{core}}\sim\mathcal{N}(y, \sigma_{\mathsf{core}}^2)}\left[\frac{1}{(1 - p_{\mathsf{maj}})^2}\left(\frac{-y\exp(-ycx_{\mathsf{core}})}{1 + \exp(-ycx_{\mathsf{core}})}\right)^2 x_{\mathsf{core}}^2\right]$$
$$= \frac{1}{p_{\mathsf{maj}}(1 - p_{\mathsf{maj}})}\mathbb{E}_y\mathbb{E}_{x_{\mathsf{core}}\sim\mathcal{N}(y, \sigma_{\mathsf{core}}^2)}\left[\left(\frac{-y\exp(-ycx_{\mathsf{core}})}{1 + \exp(-ycx_{\mathsf{core}})}\right)^2 x_{\mathsf{core}}^2\right]$$
$$= \frac{1}{2p_{\mathsf{maj}}(1 - p_{\mathsf{maj}})}\mathbb{E}_{x_{\mathsf{core}}\sim\mathcal{N}(1, \sigma_{\mathsf{core}}^2)}\left[\left(\frac{-\exp(-cx_{\mathsf{core}})}{1 + \exp(-cx_{\mathsf{core}})}\right)^2 x_{\mathsf{core}}^2\right] + \frac{1}{2p_{\mathsf{maj}}(1 - p_{\mathsf{maj}})}\mathbb{E}_{x_{\mathsf{core}}\sim\mathcal{N}(-1, \sigma_{\mathsf{core}}^2)}\left[\left(\frac{-\exp(cx_{\mathsf{core}})}{1 + \exp(cx_{\mathsf{core}})}\right)^2 x_{\mathsf{core}}^2\right]$$
$$= \frac{1}{p_{\mathsf{maj}}(1 - p_{\mathsf{maj}})}\mathbb{E}_{x_{\mathsf{core}}\sim\mathcal{N}(1, \sigma_{\mathsf{core}}^2)}\left[\left(\frac{-\exp(-cx_{\mathsf{core}})}{1 + \exp(-cx_{\mathsf{core}})}\right)^2 x_{\mathsf{core}}^2\right]$$
$$\leq \frac{1}{p_{\mathsf{maj}}(1 - p_{\mathsf{maj}})}\mathbb{E}_{x_{\mathsf{core}}\sim\mathcal{N}(1, \sigma_{\mathsf{core}}^2)}[x_{\mathsf{core}}^2] = \frac{\sigma_{\mathsf{core}}^2 + 1}{p_{\mathsf{maj}}(1 - p_{\mathsf{maj}})}.$$

Finally,

$$\mathbb{E}[\nabla_{\mathsf{spu}}(\ell_{\mathsf{rw}}(x, y, w^\star))^2] - (\mathbb{E}[\nabla_{\mathsf{spu}}\ell_{\mathsf{rw}}(x, y, w^\star)])^2$$
$$= \mathbb{E}[\nabla_{\mathsf{spu}}(\ell_{\mathsf{rw}}(x, y, w^\star))^2]$$
$$= p_{\mathsf{maj}}\mathbb{E}_y\mathbb{E}_{x_{\mathsf{core}}\sim\mathcal{N}(y, \sigma_{\mathsf{core}}^2)}\mathbb{E}_{x_{\mathsf{spu}}\sim\mathcal{N}(y, \sigma_{\mathsf{spu}}^2)}\left[\frac{1}{p_{\mathsf{maj}}^2}\left(\frac{-y\exp(-ycx_{\mathsf{core}})}{1 + \exp(-ycx_{\mathsf{core}})}\right)^2 x_{\mathsf{spu}}^2\right]$$
$$+ (1 - p_{\mathsf{maj}})\mathbb{E}_y\mathbb{E}_{x_{\mathsf{core}}\sim\mathcal{N}(y, \sigma_{\mathsf{core}}^2)}\mathbb{E}_{x_{\mathsf{spu}}\sim\mathcal{N}(-y, \sigma_{\mathsf{spu}}^2)}\left[\frac{1}{(1 - p_{\mathsf{maj}})^2}\left(\frac{-y\exp(-ycx_{\mathsf{core}})}{1 + \exp(-ycx_{\mathsf{core}})}\right)^2 x_{\mathsf{spu}}^2\right]$$
$$\leq \frac{1}{p_{\mathsf{maj}}}\mathbb{E}_y\mathbb{E}_{x_{\mathsf{spu}}\sim\mathcal{N}(y, \sigma_{\mathsf{spu}}^2)}[x_{\mathsf{spu}}^2] + \frac{1}{1 - p_{\mathsf{maj}}}\mathbb{E}_y\mathbb{E}_{x_{\mathsf{spu}}\sim\mathcal{N}(-y, \sigma_{\mathsf{spu}}^2)}[x_{\mathsf{spu}}^2] = \frac{\sigma_{\mathsf{spu}}^2 + 1}{p_{\mathsf{maj}}(1 - p_{\mathsf{maj}})}.$$

$\square$

**Lemma 13.** *The following is true.*

$$\nabla^2 L_{\mathsf{rw}}(x, y, w^\star)] \succeq \mathrm{diag}\left(\frac{\exp\left(\frac{-4}{(\sigma_{\mathsf{core}}^2 + 8)\sigma_{\mathsf{core}}^2}\right)(\sigma_{\mathsf{core}}^2 + 9)}{4(1 + 8/\sigma_{\mathsf{core}}^2)^{3/2}}, \frac{\exp\left(\frac{-4}{(\sigma_{\mathsf{core}}^2 + 8)\sigma_{\mathsf{core}}^2}\right)(\sigma_{\mathsf{spu}}^2 + 1)}{4\sqrt{1 + 8/\sigma_{\mathsf{core}}^2}}\right). \tag{110}$$

*Proof.* We use the following expression for computing the population gradient.

$$\nabla^2 \log\left(\frac{1}{1 + \exp(-yw^\top x)}\right) = \nabla\left(\frac{-y\exp(-yw^\top x)}{1 + \exp(-yw^\top x)}\right)x = \nabla\left(\frac{-y}{1 + \exp(yw^\top x)}\right)x = \left(\frac{\exp(yw^\top x)}{(1 + \exp(yw^\top x))^2}\right)xx^\top. \tag{111}$$

Recall the definition of the population majority and minority losses (Equation 102).

$$\nabla^2 \mathsf{L}_{\text{rw-maj}}(w) = \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma^2_{\text{core}})} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(y, \sigma^2_{\text{spu}})} \left[ \frac{1}{p_{\text{maj}}} \left( \frac{\exp(y w^\top x)}{(1 + \exp(y w^\top x))^2} \right) x x^\top \right]. \tag{112}$$

$$\nabla^2 \mathsf{L}_{\text{rw-min}}(w) = \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma^2_{\text{core}})} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(-y, \sigma^2_{\text{spu}})} \left[ \frac{1}{1 - p_{\text{maj}}} \left( \frac{\exp(y w^\top x)}{(1 + \exp(y w^\top x))^2} \right) x x^\top \right]. \tag{113}$$

Like previously, we first compute the off-diagonal entries.

$$\begin{aligned}
[\nabla^2 \mathsf{L}_{\text{rw-maj}}(w^\star)]_{\text{spu, core}} &= \frac{1}{p_{\text{maj}}} \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma^2_{\text{core}})} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(y, \sigma^2_{\text{spu}})} \left[ \left( \frac{\exp(y w^{\star\top} x)}{(1 + \exp(y w^{\star\top} x))^2} \right) x_{\text{core}} x_{\text{spu}} \right] \\
&\quad + \frac{1}{p_{\text{maj}}} \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma^2_{\text{core}})} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(-y, \sigma^2_{\text{spu}})} \left[ \left( \frac{\exp(y w^{\star\top} x)}{(1 + \exp(y w^{\star\top} x))^2} \right) x_{\text{core}} x_{\text{spu}} \right] \\
&= \frac{1}{p_{\text{maj}}} \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma^2_{\text{core}})} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(y, \sigma^2_{\text{spu}})} \left[ \left( \frac{\exp(y w^{\star\top} x)}{(1 + \exp(y w^{\star\top} x))^2} \right) x_{\text{core}} x_{\text{spu}} \right] \\
&\quad - \frac{1}{p_{\text{maj}}} \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma^2_{\text{core}})} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(y, \sigma^2_{\text{spu}})} \left[ \left( \frac{\exp(y w^{\star\top} x)}{(1 + \exp(y w^{\star\top} x))^2} \right) x_{\text{core}} x_{\text{spu}} \right] \\
&= 0
\end{aligned}$$

$$[\nabla^2 \mathsf{L}_{\text{rw-min}}(w^\star)]_{\text{spu, core}} = 0, \text{ Similar calculation as above}$$
$$[\nabla^2 \mathsf{L}_{\text{rw}}(w^\star)]_{\text{spu, core}} = 0.$$

Now, we bound the diagonal entries. Recall that $w^\star_{\text{spu}} = 0$ and $w^\star_{\text{core}} = c$ where $c = \frac{2}{\sigma^2_{\text{core}}}$.

$$\begin{aligned}
[\nabla^2 \mathsf{L}_{\text{rw-maj}}(w^\star)]_{\text{core, core}} &= \frac{1}{p_{\text{maj}}} \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma^2_{\text{core}})} \left[ \left( \frac{\exp(y c x_{\text{core}})}{(1 + \exp(y c x_{\text{core}}))^2} \right) x^2_{\text{core}} \right] \\
&= \frac{1}{2 p_{\text{maj}}} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma^2_{\text{core}})} \left[ \left( \frac{\exp(c x_{\text{core}})}{(1 + \exp(c x_{\text{core}}))^2} \right) x^2_{\text{core}} \right] + \frac{1}{2 p_{\text{maj}}} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(-1, \sigma^2_{\text{core}})} \left[ \left( \frac{\exp(-c x_{\text{core}})}{(1 + \exp(-c x_{\text{core}}))^2} \right) x^2_{\text{core}} \right] \\
&= \frac{1}{p_{\text{maj}}} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma^2_{\text{core}})} \left[ \left( \frac{\exp(c x_{\text{core}})}{(1 + \exp(c x_{\text{core}}))^2} \right) x^2_{\text{core}} \right] \\
&\geq \frac{1}{p_{\text{maj}}} \frac{1}{4} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma^2_{\text{core}})} \left[ \exp(-c^2 x^2_{\text{core}}) x^2_{\text{core}} \right] \\
&= \frac{1}{p_{\text{maj}}} \frac{1}{4 \sigma_{\text{core}} \sqrt{2\pi}} \int_{-\infty}^{\infty} \exp(-c^2 x^2_{\text{core}}) \exp\left( \frac{-(x_{\text{core}} - 1)^2}{2 \sigma^2_{\text{core}}} \right) x^2_{\text{core}} \, dx_{\text{core}} \\
&= \frac{1}{p_{\text{maj}}} \frac{1}{4 \sigma_{\text{core}} \sqrt{2\pi}} \int_{-\infty}^{\infty} \exp\left( - \frac{8 x^2_{\text{core}} / \sigma^2_{\text{core}}}{2 \sigma^2_{\text{core}}} \right) \exp\left( \frac{-(x_{\text{core}} - 1)^2}{2 \sigma^2_{\text{core}}} \right) x^2_{\text{core}} \, dx_{\text{core}} \\
&= \frac{1}{p_{\text{maj}}} \frac{\exp\left( \frac{-8}{(\sigma^2_{\text{core}} + 8) \sigma^2_{\text{core}}} \right)}{4 \sigma_{\text{core}} \sqrt{2\pi}} \int_{-\infty}^{\infty} \exp\left( \frac{-(\sqrt{1 + 8/\sigma^2_{\text{core}}} x_{\text{core}} - \frac{1}{\sqrt{1 + 8/\sigma^2_{\text{core}}}})^2}{2 \sigma^2_{\text{core}}} \right) x^2_{\text{core}} \, dx_{\text{core}} \\
&= \frac{1}{p_{\text{maj}}} \frac{\exp\left( \frac{-8}{(\sigma^2_{\text{core}} + 8) \sigma^2_{\text{core}}} \right) (\sigma^2_{\text{core}} + 9)}{4(1 + 8/\sigma^2_{\text{core}})^{5/2}}.
\end{aligned}$$

$$[\nabla^2 \mathsf{L}_{\text{rw-min}}(w^\star)]_{\text{core, core}} = \frac{1}{1 - p_{\text{maj}}} \frac{\exp\left( \frac{-8}{(\sigma^2_{\text{core}} + 8) \sigma^2_{\text{core}}} \right) (\sigma^2_{\text{core}} + 9)}{4(1 + 8/\sigma^2_{\text{core}})^{5/2}}, \text{ By symmetry.}$$

$$\begin{aligned}
[\nabla^2 \mathsf{L}_{\text{rw}}(w^\star)]_{\text{core, core}} &= p_{\text{maj}} [\nabla^2 \mathsf{L}_{\text{rw-maj}}(w^\star)]_{\text{core, core}} + (1 - p_{\text{maj}}) [\nabla^2 \mathsf{L}_{\text{rw-min}}(w^\star)]_{\text{core, core}} \\
&= \frac{\exp\left( \frac{-8}{(\sigma^2_{\text{core}} + 8) \sigma^2_{\text{core}}} \right) (\sigma^2_{\text{core}} + 9)}{4(1 + 8/\sigma^2_{\text{core}})^{5/2}}.
\end{aligned}$$

Finally, we calculate $[\nabla^2 L_{\text{rw-maj}}(w^\star)]_{\text{spu, spu}}$ as follows.

$$[\nabla^2 L_{\text{rw-maj}}(w^\star)]_{\text{spu, spu}} = \frac{1}{p_{\text{maj}}} \mathbb{E}_y \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(y, \sigma_{\text{core}}^2)} \mathbb{E}_{x_{\text{spu}} \sim \mathcal{N}(y, \sigma_{\text{spu}}^2)} \left[ \left( \frac{\exp(ycx_{\text{core}})}{(1 + \exp(ycx_{\text{core}}))^2} \right) x_{\text{spu}}^2 \right]$$

$$= \frac{1}{2p_{\text{maj}}} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} \left[ \left( \frac{\exp(cx_{\text{core}})}{(1 + \exp(cx_{\text{core}}))^2} \right) \right] (\sigma_{\text{spu}}^2 + 1)$$

$$+ \frac{1}{2p_{\text{maj}}} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(-1, \sigma_{\text{core}}^2)} \left[ \left( \frac{\exp(-cx_{\text{core}})}{(1 + \exp(-cx_{\text{core}}))^2} \right) \right] (\sigma_{\text{spu}}^2 + 1)$$

$$\geq \frac{1}{4p_{\text{maj}}} \mathbb{E}_{x_{\text{core}} \sim \mathcal{N}(1, \sigma_{\text{core}}^2)} [\exp(-c^2 x_{\text{core}}^2)](\sigma_{\text{spu}}^2 + 1)$$

$$= \frac{1}{4p_{\text{maj}}} \frac{\exp\left( \frac{-4}{(\sigma_{\text{core}}^2 + 8)\sigma_{\text{core}}^2} \right)}{\sqrt{1 + 8/\sigma_{\text{core}}^2}} (\sigma_{\text{spu}}^2 + 1)$$

$$[\nabla^2 L_{\text{rw-min}}(w^\star)]_{\text{spu, spu}} = \frac{1}{4(1 - p_{\text{maj}})} \frac{\exp\left( \frac{-4}{(\sigma_{\text{core}}^2 + 8)\sigma_{\text{core}}^2} \right)}{\sqrt{1 + 8/\sigma_{\text{core}}^2}} (\sigma_{\text{spu}}^2 + 1), \text{ By symmetry.}$$

$$[\nabla^2 L_{\text{rw}}(w^\star)]_{\text{spu, spu}} = \frac{\exp\left( \frac{-4}{(\sigma_{\text{core}}^2 + 8)\sigma_{\text{core}}^2} \right)(\sigma_{\text{spu}}^2 + 1)}{4\sqrt{1 + 8/\sigma_{\text{core}}^2}}.$$

$\square$

**Proposition 5.** *The asymptotic distribution of the reweighted logistic regression estimator is as follows.*

$$\sqrt{n}(\hat{w} - w^\star) \to^d \mathcal{N}(0, V), \tag{95}$$

$$V \preceq \text{diag} \left( \frac{16 \exp\left( \frac{8}{(\sigma_{\text{core}}^2 + 8)\sigma_{\text{core}}^2} \right)(\sigma_{\text{core}}^2 + 1)(1 + 8/\sigma_{\text{core}}^2)^3}{p_{\text{maj}}(1 - p_{\text{maj}})(\sigma_{\text{core}}^2 + 9)^2}, \frac{16 \exp\left( \frac{8}{(\sigma_{\text{core}}^2 + 8)\sigma_{\text{core}}^2} \right)(1 + 8/\sigma_{\text{core}}^2)}{p_{\text{maj}}(1 - p_{\text{maj}})(\sigma_{\text{spu}}^2 + 1)} \right). \tag{96}$$

*For $\sigma_{\text{core}} \geq 1$, we have*

$$V \preceq \text{diag} \left( \frac{C_1}{p_{\text{maj}}(1 - p_{\text{maj}})}, \frac{C_2}{p_{\text{maj}}(1 - p_{\text{maj}})} \right), \tag{97}$$

*for some constants $C_1, C_2$.*

*Proof.* By asymptotic normality, we have $\sqrt{n}(\hat{w} - w^\star) \to \mathcal{N}(0, \nabla^2 L(w^\star)^{-1} \text{Cov}[\nabla \ell(x, y, w^\star)] \nabla^2 L(w^\star)^{-1})$. Combining Lemma 12 and Lemma 13, we get the expression in Equation 96. Each term is decreasing in $\sigma_{\text{core}}$, and hence we get the final result by substituting $\sigma_{\text{core}}^2 = 1$ to obtain the constants $C_1, C_2$ (and noting that $\sigma_{\text{spu}}^2 \geq 0$). $\square$