
Undersampling is a Minimax Optimal Robustness Intervention in Nonparametric Classification

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 While a broad range of techniques have been proposed to tackle distribution shift,
2 the simple baseline of training on an *undersampled* balanced dataset often achieves
3 close to state-of-the-art-accuracy across several popular benchmarks. This is rather
4 surprising, since undersampling algorithms discard excess majority group data.
5 To understand this phenomenon, we ask if learning is fundamentally constrained
6 by a lack of minority group samples. We prove that this is indeed the case in
7 the setting of nonparametric binary classification. Our results show that in the
8 worst case, an algorithm cannot outperform undersampling unless there is a high
9 degree of overlap between the train and test distributions (which is unlikely to be
10 the case in real-world datasets), or if the algorithm leverages additional structure
11 about the distribution shift. In particular, in the case of label shift we show that
12 there is always an undersampling algorithm that is minimax optimal. In the case
13 of group-covariate shift we show that there is an undersampling algorithm that is
14 minimax optimal when the overlap between the group distributions is small. We
15 also perform an experimental case study on a label shift dataset and find that in line
16 with our theory, the test accuracy of robust neural network classifiers is constrained
17 by the number of minority samples.

18 1 Introduction

19 A key challenge facing the machine learning community is to design models that are robust to
20 distribution shift. When there is a mismatch between the train and test distributions, current models
21 are often brittle and perform poorly on rare examples [Hovy and Søgaard, 2015, Blodgett et al.,
22 2016, Tatman, 2017, Hashimoto et al., 2018, Alcorn et al., 2019]. In this paper, our focus is on
23 group-structured distribution shifts. In the training set, we have many samples from a *majority* group
24 and relatively few samples from the *minority* group, while during test time we are equally likely to
25 get a sample from either group.

26 To tackle such distribution shifts, a naïve algorithm is one that first *undersamples* the training data
27 by discarding excess majority group samples [Kubat and Matwin, 1997, Wallace et al., 2011] and
28 then trains a model on this resulting dataset. The samples that remain in this undersampled dataset
29 constitute i.i.d. draws from the test distribution. Therefore, while a classifier trained on this pruned
30 dataset cannot suffer biases due to distribution shift, this algorithm is clearly wasteful, as it discards
31 training samples. This perceived inefficiency of undersampling has led to the design of several
32 algorithms to combat such distribution shift [Chawla et al., 2002, Lipton et al., 2018, Sagawa et al.,
33 2020, Cao et al., 2019, Menon et al., 2020, Ye et al., 2020, Kini et al., 2021, Wang et al., 2022].
34 In spite of this algorithmic progress, the simple baseline of training models on an undersampled
35 dataset remains competitive. In the case of label shift, where one class label is overrepresented in the

36 training data, this has been observed by Cui et al. [2019], Cao et al. [2019], and Yang and Xu [2020].
 37 While in the case of group-covariate shift, a study by Idrissi et al. [2022] showed that the empirical
 38 effectiveness of these more complicated algorithms is limited.

39 For example, Idrissi et al. [2022] showed that on the group-covariate shift CelebA dataset the worst-
 40 group accuracy of a ResNet-50 model on the undersampled CelebA dataset which *discards 97%* of
 41 the available training data is as good as methods that use all of available data such as importance-
 42 weighted ERM [Shimodaira, 2000], Group-DRO [Sagawa et al., 2020] and Just-Train-Twice [Liu
 43 et al., 2021]. In Table 1, we report the performance of the undersampled classifier compared to the
 44 state-of-the-art-methods in the literature across several label shift and group-covariate shift datasets.
 45 We find that, although undersampling isn’t always the optimal robustness algorithm, it is typically a
 very competitive baseline and within 1–4% the performance of the best method.

Table 1: Performance of undersampled classifier compared to the best classifier across several popular label shift and group-covariate shift datasets. When reporting worst-group accuracy we denote it by a *. When available, we report the 95% confidence interval. We find that the undersampled classifier is always within 1–4% of the best performing robustness algorithm, except on the MultiNLI dataset.

Shift Type	Dataset/Paper	Test/Worst-Group* Accuracy	
		Best	Undersampled
Label	Imb. CIFAR10 (step 10) [Cao et al., 2019]	87.81	84.59
	Imb. CIFAR100 (step 10) [Cao et al., 2019]	58.71	55.06
Group-Covariate	CelebA [Idrissi et al., 2022]	$86.9 \pm 1.1^*$	$85.6 \pm 2.3^*$
	Waterbirds [Idrissi et al., 2022]	$87.6 \pm 1.6^*$	$89.1 \pm 1.1^*$
	MultiNLI [Idrissi et al., 2022]	$78.0 \pm 0.7^*$	$68.9 \pm 0.8^*$
	CivilComments [Idrissi et al., 2022]	$72.0 \pm 1.9^*$	$71.8 \pm 1.4^*$

46

47 Inspired by the strong performance of undersampling in these experiments, we ask:

48 *Is the performance of a model under distribution shift fundamentally*
 49 *constrained by the lack of minority group samples?*

50 To answer this question we analyze the *minimax excess risk*. We lower bound the minimax excess risk
 51 to prove that the performance of *any* algorithm is lower bounded only as a function of the minority
 52 samples (n_{\min}). This shows that even if a robust algorithm optimally trades off between the bias and
 53 the variance, it is fundamentally constrained by the variance on the minority group which decreases
 54 only with n_{\min} .

55 **Our Contributions.** In our paper, we consider the well-studied setting of nonparametric binary
 56 classification [Tsybakov, 2010]. By operating in this nonparametric regime we are able to study the
 57 properties of undersampling in rich data distributions, but are able to circumvent the complications
 58 that arise due to the optimization and implicit bias of parametric models.

59 We provide insights into this question in the label shift scenario, where one of the labels is overrep-
 60 resented in the training data, $P_{\text{train}}(y = 1) \geq P_{\text{train}}(y = -1)$, whereas the test samples are equally
 61 likely to come from either class. Here the class-conditional distribution $P(x | y)$ is Lipschitz in x .
 62 We show that in the label shift setting there is a fundamental constraint, and that the minimax excess
 63 risk of *any robust learning method* is lower bounded by $1/n_{\min}^{1/3}$. That is, minority group samples
 64 fundamentally constrain performance under distribution shift. Furthermore, by leveraging previous
 65 results about nonparametric density estimation [Freedman and Diaconis, 1981] we show a matching
 66 upper bound on the excess risk of a standard binning estimator trained on an undersampled dataset to
 67 demonstrate that undersampling is optimal (see Theorem D.1).

68 Further, we experimentally show in a label shift dataset (Imbalanced Binary CIFAR10) that the
 69 accuracy of popular classifiers generally follow the trends predicted by our theory (see Appendix C).
 70 When the minority samples are increased, the accuracy of these classifiers increases drastically,
 71 whereas when the number of majority samples are increased the gains in the accuracy are marginal.

72 We also study the covariate shift case. In this setting, there has been extensive work studying the
73 effectiveness of transfer [Kpotufe and Martinet, 2018, Hanneke and Kpotufe, 2019] from train to test
74 distributions, often focusing on deriving specific conditions under which this transfer is possible. In
75 this work, we demonstrate that when the overlap (defined in terms of total variation distance) between
76 the group distributions P_a and P_b is small, transfer is difficult, and that the minimax excess risk of any
77 robust learning algorithm is lower bounded by $1/n_{\min}^{1/3}$ (see Theorem B.1). While this prior work
78 also shows the impossibility of using majority group samples in the extreme case with no overlap, our
79 results provide a simple lower bound that shows that the amount of overlap needed to make transfer
80 feasible is unrealistic. We also show that this lower bound is tight, by proving an upper bound on the
81 excess risk of the binning estimator acting on the undersampled dataset (see Theorem D.2).

82 Taken together, our results underline the need to move beyond designing “general-purpose” robustness
83 algorithms (like importance-weighting [Cao et al., 2019, Menon et al., 2020, Kini et al., 2021, Wang
84 et al., 2022], g-DRO [Sagawa et al., 2020], JTT [Liu et al., 2021], SMOTE [Chawla et al., 2002], etc.)
85 that are agnostic to the structure in the distribution shift. Our worst case analysis highlights that to
86 successfully beat undersampling, an algorithm must leverage additional structure in the distribution
87 shift.

88 **Organization.** We present our minimax lower bounds on the label shift in the main paper. The
89 matching upper bounds we proved in the appendix. The upper and lower bounds in the group-
90 covariate shift are presented in the appendix. Discussion of related work and simulations studying
91 the minority group sample dependence in robust neural networks classifiers are also in the appendix.

92 2 Setting

93 The setting for our study is nonparametric binary classification with Lipschitz data distributions.
94 We are given n training datapoints $\mathcal{S} := \{(x_1, y_1), \dots, (x_n, y_n)\} \in ([0, 1] \times \{-1, 1\})^n$ that are all
95 drawn from a *train* distribution P_{train} . During test time, the data shall be drawn from a *different*
96 distribution P_{test} . To present a clean analysis, we study the case where the features x are bounded
97 scalars, however, it is easy to extend our results to the high-dimensional setting.

98 Given a classifier $f : \mathbb{R} \rightarrow \{-1, 1\}$, we shall be interested in the test error (risk) of this classifier
99 under the test distribution P_{test} :

$$R(f; P_{\text{test}}) := \mathbb{E}_{(x,y) \sim P_{\text{test}}} [\mathbf{1}(f(x) \neq y)].$$

100 We assume that P_{train} consists of a mixture of two groups of unequal size, and P_{test} contains equal
101 numbers of samples from both groups. Given a majority group distribution P_{maj} and a minority
102 group distribution P_{min} , the learner has access to n_{maj} majority group samples and n_{min} minority
103 group samples: $\mathcal{S}_{\text{maj}} \sim P_{\text{maj}}^{n_{\text{maj}}}$ and $\mathcal{S}_{\text{min}} \sim P_{\text{min}}^{n_{\text{min}}}$. Here $n_{\text{maj}} > n/2$ and $n_{\text{min}} < n/2$ with
104 $n_{\text{maj}} + n_{\text{min}} = n$. The full training dataset is $\mathcal{S} = \mathcal{S}_{\text{maj}} \cup \mathcal{S}_{\text{min}} = \{(x_1, y_1), \dots, (x_n, y_n)\}$. We
105 assume that the learner has access to the knowledge whether a particular sample (x_i, y_i) comes from
106 the majority or minority group.

107 The test samples will be drawn from $P_{\text{test}} = \frac{1}{2}P_{\text{maj}} + \frac{1}{2}P_{\text{min}}$, a uniform mixture over P_{maj} and P_{min} .
108 Thus, the training dataset is an imbalanced draw from the distributions P_{maj} and P_{min} , whereas the
109 test samples are balanced draws. We let $\rho := n_{\text{maj}}/n_{\text{min}} > 1$ denote the imbalance ratio in the
110 training data.

111 We focus on two-types of distribution shifts: label shift that we describe below, and group-covariate
112 shift that we describe in Appendix G.1.

Label Shift. In this setting, the imbalance in the training data comes from there being more samples
from one class over another. Without loss of generality, we shall assume that the class $y = 1$ is the
majority class. Then, we define the majority and the minority class distributions as

$$P_{\text{maj}}(x, y) = P_1(x)\mathbf{1}(y = 1) \quad \text{and} \quad P_{\text{min}} = P_{-1}(x)\mathbf{1}(y = -1),$$

113 where P_1, P_{-1} are class-conditional distributions over the interval $[0, 1]$. We assume that class-
114 conditional distributions P_i have densities on $[0, 1]$ and that they are 1-Lipschitz: for any $x, x' \in [0, 1]$,

$$|P_i(x) - P_i(x')| \leq |x - x'|.$$

115 We denote the class of pairs of distributions $(P_{\text{maj}}, P_{\text{min}})$ that satisfy these conditions by \mathcal{P}_{LS} . We
 116 note that such Lipschitzness assumptions are common in the literature [see Tsybakov, 2010].

117 3 Lower Bounds on the Minimax Excess Risk

118 In this section, we shall prove our lower bounds that show that the performance of any algorithm is
 119 constrained by the number of minority samples n_{min} . Before we state our lower bounds, we need to
 120 introduce the notion of excess risk and minimax excess risk.

121 **Excess Risk and Minimax Excess Risk.** We measure the performance of an algorithm \mathcal{A} through
 122 its excess risk defined in the following way. Given an algorithm \mathcal{A} that takes as input a dataset \mathcal{S}
 123 and returns a classifier $\mathcal{A}^{\mathcal{S}}$, and a pair of distributions $(P_{\text{maj}}, P_{\text{min}})$ with $P_{\text{test}} = \frac{1}{2}P_{\text{maj}} + \frac{1}{2}P_{\text{min}}$, the
 124 *expected excess risk* is given by

$$\text{Excess Risk}[\mathcal{A}; (P_{\text{maj}}, P_{\text{min}})] := \mathbb{E}_{\mathcal{S} \sim P_{\text{maj}}^{n_{\text{maj}}} \times P_{\text{min}}^{n_{\text{min}}}} [R(\mathcal{A}^{\mathcal{S}}; P_{\text{test}}) - R(f^*(P_{\text{test}}); P_{\text{test}})], \quad (1)$$

125 where $f^*(P_{\text{test}})$ is the Bayes classifier that minimizes the risk $R(\cdot; P_{\text{test}})$. The first term corresponds
 126 to the expected risk for the algorithm when given n_{maj} samples from P_{maj} and n_{min} samples from
 127 P_{min} , whereas the second term corresponds to the Bayes error for the problem.

128 Excess risk does not let us characterize the inherent difficulty of a problem, since for any particular
 129 data distribution $(P_{\text{maj}}, P_{\text{min}})$ the best possible algorithm \mathcal{A} to minimize the excess risk would be the
 130 trivial mapping $\mathcal{A}^{\mathcal{S}} = f^*(P_{\text{test}})$. Therefore, to prove meaningful lower bounds on the performance of
 131 algorithms we need to define the notion of minimax excess risk [see Wainwright, 2019, Chapter 15].
 132 Given a class of pairs of distributions \mathcal{P} define

$$\text{Minimax Excess Risk}(\mathcal{P}) := \inf_{\mathcal{A}} \sup_{(P_{\text{maj}}, P_{\text{min}}) \in \mathcal{P}} \text{Excess Risk}[\mathcal{A}; (P_{\text{maj}}, P_{\text{min}})], \quad (2)$$

133 where the infimum is over all measurable estimators \mathcal{A} . The minimax excess risk is the excess risk of
 134 the “best” algorithm in the worst case over the class of problems defined by \mathcal{P} .

135 We demonstrate the hardness of the label shift problem in general by establishing a lower bound on
 136 the minimax excess risk.

137 **Theorem 3.1.** *Let \mathcal{P}_{LS} be the class of pairs of distributions $(P_{\text{maj}}, P_{\text{min}})$ that satisfy the label-shift*
 138 *assumptions. The minimax excess risk over this class is lower bounded as follows:*

$$\text{Minimax Excess Risk}(\mathcal{P}_{\text{LS}}) = \inf_{\mathcal{A}} \sup_{(P_{\text{maj}}, P_{\text{min}}) \in \mathcal{P}_{\text{LS}}} \text{Excess Risk}[\mathcal{A}; (P_{\text{maj}}, P_{\text{min}})] \geq \frac{1}{600} \frac{1}{n_{\text{min}}^{1/3}}. \quad (3)$$

139 We establish this result in Appendix F. We show that rather surprisingly, the lower bound on the
 140 minimax excess risk scales only with the number of minority class samples $n_{\text{min}}^{1/3}$, and does
 141 not depend on n_{maj} . Intuitively, this is because any learner must predict which class-conditional
 142 distribution $(P(x | 1) \text{ or } P(x | -1))$ assigns higher likelihood at that x . To interpret this result,
 143 consider the extreme scenario where $n_{\text{maj}} \rightarrow \infty$ but n_{min} is finite. In this case, the learner has
 144 full information about the majority class distribution. However, the learning task continues to be
 145 challenging since any learner would be uncertain about whether the minority class distribution assigns
 146 higher or lower likelihood at any given x . This uncertainty underlies the reason why the minimax
 147 rate of classification is constrained by the number of minority samples n_{min} .

148 We also note that the theorem can be trivially extended to higher dimensions. In this case the
 149 exponents degrade to $1/3d$ rather than $1/3$ as is to be expected in nonparametric classification.

150 **Discussion.** We showed that undersampling is an optimal robustness intervention in nonparametric
 151 classification in the absence of significant overlap between group distributions or without additional
 152 structure beyond Lipschitz continuity. At a high level our results highlight the need to reason about
 153 the specific structure in the distribution shift and design algorithms that are tailored to take advantage
 154 of this structure. This would require us to step away from the common practice in robust machine
 155 learning where the focus is to design “universal” robustness interventions that are agnostic to the
 156 structure in the shift. Alongside this, our results also dictate the need for datasets and benchmarks
 157 with the propensity for transfer from train to test time.

References

- 158
- 159 M. Alcorn, Q. Li, Z. Gong, C. Wang, L. Mai, W.-S. Ku, and A. Nguyen. Strike (with) a pose: Neural
160 networks are easily fooled by strange poses of familiar objects. In *Computer Vision and Pattern
161 Recognition (CVPR)*, 2019.
- 162 M. Arjovsky, K. Chaudhuri, and D. Lopez-Paz. Throwing away data improves worst-class error in
163 imbalanced classification. *arXiv preprint arXiv:2205.11672*, 2022.
- 164 S. Ben-David and R. Urner. On the hardness of domain adaptation and the utility of unlabeled target
165 samples. In *Algorithmic Learning Theory (ALT)*, 2012.
- 166 S. Ben-David and R. Urner. Domain adaptation—can quantity compensate for quality? *Annals of
167 Mathematics and Artificial Intelligence*, 2014.
- 168 S. Ben-David, J. Blitzer, K. Crammer, and F. Pereira. Analysis of representations for domain
169 adaptation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2006.
- 170 S. Ben-David, J. Blitzer, K. Crammer, A. Kulesza, F. Pereira, and J. W. Vaughan. A theory of learning
171 from different domains. *Machine learning*, 2010.
- 172 C. Berlind and R. Urner. Active nearest neighbors in changing environments. In *International
173 Conference on Machine Learning (ICML)*, 2015.
- 174 S. L. Blodgett, L. Green, and B. O’Connor. Demographic dialectal variation in social media: A
175 case study of african-american english. In *Empirical Methods in Natural Language Processing
176 (EMNLP)*, 2016.
- 177 J. Byrd and Z. Lipton. What is the effect of importance weighting in deep learning? In *International
178 Conference on Machine Learning (ICML)*, 2019.
- 179 C. Canonne. A short note on an inequality between KL and TV. *arXiv preprint arXiv:2202.07198*,
180 2022.
- 181 K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma. Learning imbalanced datasets with label-
182 distribution-aware margin loss. In *Advances in Neural Information Processing Systems (NeurIPS)*,
183 2019.
- 184 N. Chawla, K. Bowyer, L. Hall, and P. Kegelmeyer. Smote: Synthetic minority over-sampling
185 technique. *Journal of Artificial Intelligence Research*, 2002.
- 186 Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie. Class-balanced loss based on effective number of
187 samples. In *Computer Vision and Pattern Recognition (CVPR)*, 2019.
- 188 S. B. David, T. Lu, T. Luu, and D. Pál. Impossibility theorems for domain adaptation. In *International
189 Conference on Artificial Intelligence and Statistics (AISTATS)*, 2010.
- 190 L. Devroye and L. Györfi. *Nonparametric density estimation: the L_1 view*. Wiley Series in Probability
191 and Mathematical Statistics, 1985.
- 192 D. Freedman and P. Diaconis. On the histogram as a density estimator: L_2 theory. *Zeitschrift für
193 Wahrscheinlichkeitstheorie und verwandte Gebiete*, 1981.
- 194 S. Hanneke and S. Kpotufe. On the value of target data in transfer learning. In *Advances in Neural
195 Information Processing Systems (NeurIPS)*, 2019.
- 196 T. Hashimoto, M. Srivastava, H. Namkoong, and P. Liang. Fairness without demographics in repeated
197 loss minimization. In *International Conference on Machine Learning (ICML)*, 2018.
- 198 D. Hovy and A. Søgaard. Tagging performance correlates with author age. In *Association for
199 Computational Linguistics (ACL)*, 2015.
- 200 B. Y. Idrissi, M. Arjovsky, M. Pezeshki, and D. Lopez-Paz. Simple data balancing achieves competi-
201 tive worst-group-accuracy. In *Causal Learning and Reasoning*, 2022.

- 202 G. R. Kini, O. Paraskevas, S. Oymak, and C. Thrampoulidis. Label-imbalanced and group-sensitive
 203 classification under overparameterization. In *Advances in Neural Information Processing Systems*
 204 (*NeurIPS*), 2021.
- 205 S. Kpotufe and G. Martinet. Marginal singularity, and the benefits of labels in covariate-shift. In
 206 *Conference On Learning Theory (COLT)*, 2018.
- 207 M. Kubat and S. Matwin. Addressing the curse of imbalanced training sets: one-sided selection. In
 208 *International Conference on Machine Learning (ICML)*, 1997.
- 209 T. Li, A. Beirami, M. Sanjabi, and V. Smith. Tilted empirical risk minimization. In *International*
 210 *Conference on Learning Representations (ICLR)*, 2020.
- 211 Z. Lipton, Y.-X. Wang, and A. Smola. Detecting and correcting for label shift with black box
 212 predictors. In *International Conference on Machine Learning (ICML)*, 2018.
- 213 E. Liu, B. Haghighi, A. Chen, A. Raghunathan, P. W. Koh, S. Sagawa, P. Liang, and C. Finn. Just
 214 train twice: Improving group robustness without training group information. In *International*
 215 *Conference on Machine Learning (ICML)*, 2021.
- 216 S. Maity, Y. Sun, and M. Banerjee. Minimax optimal approaches to the label shift problem. *arXiv*
 217 *preprint arXiv:2003.10443*, 2020.
- 218 A. K. Menon, S. Jayasumana, A. S. Rawat, H. Jain, A. Veit, and S. Kumar. Long-tail learning via
 219 logit adjustment. In *International Conference on Learning Representations (ICLR)*, 2020.
- 220 S. Sagawa, P. W. Koh, T. Hashimoto, and P. Liang. Distributionally robust neural networks. In
 221 *International Conference on Learning Representations (ICLR)*, 2020.
- 222 H. Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood
 223 function. *Journal of Statistical Planning and Inference*, 2000.
- 224 R. Tatman. Gender and dialect bias in youtube’s automatic captions. In *ACL Workshop on Ethics in*
 225 *Natural Language Processing*, 2017.
- 226 A. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, 2010.
- 227 M. Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge University
 228 Press, 2019.
- 229 B. Wallace, K. Small, C. Brodley, and T. Trikalinos. Class imbalance, redux. In *International*
 230 *Conference on Data Mining (ICDM)*, 2011.
- 231 K. A. Wang, N. Chatterji, S. Haque, and T. Hashimoto. Is importance weighting incompatible with
 232 interpolating classifiers? In *International Conference on Learning Representations (ICLR)*, 2022.
- 233 L. Wasserman. Lecture notes in nonparametric classification, 2019. URL <https://www.stat.cmu.edu/~larry/=sml/nonparclass.pdf>. [Online; accessed 12-May-2022].
- 234
- 235 Wikipedia contributors. Poisson binomial distribution — Wikipedia, the free encyclope-
 236 dia, 2022. URL [https://en.wikipedia.org/w/index.php?title=Poisson_binomial_](https://en.wikipedia.org/w/index.php?title=Poisson_binomial_distribution&oldid=1071847908)
 237 [distribution&oldid=1071847908](https://en.wikipedia.org/w/index.php?title=Poisson_binomial_distribution&oldid=1071847908). [Online; accessed 5-May-2022].
- 238 D. Xu, Y. Ye, and C. Ruan. Understanding the role of importance weighting for deep learning. In
 239 *International Conference on Learning Representations (ICLR)*, 2020.
- 240 Y. Yang and Z. Xu. Rethinking the value of labels for improving class-imbalanced learning. *Advances*
 241 *in Neural Information Processing Systems (NeurIPS)*, 2020.
- 242 H.-J. Ye, H.-Y. Chen, D.-C. Zhan, and W.-L. Chao. Identifying and compensating for feature deviation
 243 in imbalanced deep learning. *arXiv preprint arXiv:2001.01385*, 2020.

244 **A Related Work**

245 On several group-covariate shift benchmarks (CelebA, CivilComments, Waterbirds), Idrissi et al.
 246 [2022] showed that training ResNet classifiers on an undersampled dataset either outperforms or
 247 performs as well as other popular reweighting methods like Group-DRO [Sagawa et al., 2020],
 248 reweighted ERM, and Just-Train-Twice [Liu et al., 2021]. They find Group-DRO performs compara-
 249 bly to undersampling, while both tend to outperform methods that don't utilize group information.

250 One classic method to tackle distribution shift is importance weighting [Shimodaira, 2000], which
 251 reweights the loss of the minority group samples to yield an unbiased estimate of the loss. However,
 252 recent work [Byrd and Lipton, 2019, Xu et al., 2020] has demonstrated the ineffectiveness of such
 253 methods when applied to overparameterized neural networks. Many followup papers [Cao et al.,
 254 2019, Ye et al., 2020, Menon et al., 2020, Kini et al., 2021, Wang et al., 2022] have introduced
 255 methods that modify the loss function in various ways to address this. However, despite this progress
 256 undersampling remains a competitive alternative to these importance weighted classifiers.

257 Our theory draws from the rich literature on non-parametric classification [Tsybakov, 2010]. Apart
 258 from borrowing this setting of nonparametric classification, we also utilize upper bounds on the
 259 estimation error of the simple histogram estimator [Freedman and Diaconis, 1981, Devroye and
 260 Györfi, 1985] to prove our upper bounds in the label shift case. Finally, we note that to prove
 261 our minimax lower bounds we proceed by using the general recipe of reducing from estimation to
 262 testing [Wainwright, 2019, Chapter 15]. One difference from this standard framework is that our
 263 training samples shall be drawn from a different distribution than the test samples used to define the
 264 risk.

265 There is rich literature that studies domain adaptation and transfer learning under label shift [Maity
 266 et al., 2020] and covariate shift [Ben-David et al., 2006, David et al., 2010, Ben-David et al., 2010,
 267 Ben-David and Uner, 2012, 2014, Berlind and Uner, 2015, Kpotufe and Martinet, 2018, Hanneke
 268 and Kpotufe, 2019]. The principal focus of this line of work was to understand the value of unlabeled
 269 data from the target domain, rather than to characterize the relative value of the number of labeled
 270 samples from the majority and minority groups. Among these papers, most closely related to our
 271 work are those in the covariate shift setting [Kpotufe and Martinet, 2018, Hanneke and Kpotufe,
 272 2019]. Their lower bound results can be reinterpreted to show that under covariate shift in the absence
 273 of overlap, the minimax excess risk is lower bounded by $1/n_{\min}^{1/3}$. We provide a more detailed
 274 comparison with their results after presenting our lower bounds in Section B.

275 Finally, we note that Arjovsky et al. [2022] recently showed that undersampling can improve the
 276 worst-class accuracy of linear SVMs in the presence of label shift. In comparison, our results hold
 277 for arbitrary classifiers with the rich nonparametric data distributions.

278 **B Group-Covariate Shift Lower Bounds**

279 First we define group-covariate shifts.

Group-Covariate Shift. In this setting, we have two groups $\{a, b\}$, and corresponding to each of
 these groups is a distribution (with densities) over the features $P_a(x)$ and $P_b(x)$. We let a correspond
 to the majority group and b correspond to the minority group. Then, we define

$$P_{\text{maj}}(x, y) = P_a(x)P(y | x) \quad \text{and} \quad P_{\text{min}}(x, y) = P_b(x)P(y | x).$$

280 We assume that for $y \in \{-1, 1\}$, for all $x, x' \in [0, 1]$:

$$|P(y | x) - P(y | x')| \leq |x - x'|,$$

281 that is, the distribution of the label given the feature is 1-Lipschitz, and it varies slowly over the
 282 domain.

283 To quantify the shift between the train and test distribution, we define a notion of overlap between the
 284 group distributions P_a and P_b as follows:

$$\text{Overlap}(P_a, P_b) := 1 - \text{TV}(P_a, P_b)$$

285 where $\text{TV}(P_a, P_b) := \sup_{E \subseteq [0,1]} |P_a(E) - P_b(E)|$, denotes the total variation distance between
 286 P_a and P_b . Notice that when P_a and P_b have disjoint supports, $\text{TV}(P_a, P_b) = 1$ and therefore

287 $\text{Overlap}(P_a, P_b) = 0$. On the other hand when $P_a = P_b$, $\text{TV}(P_a, P_b) = 0$ and $\text{Overlap}(P_a, P_b) = 1$.
 288 When the overlap is 1, the majority and minority distributions are identical and hence we have no
 289 shift between train and test. Observe that $\text{Overlap}(P_a, P_b) = \text{Overlap}(P_{\text{maj}}, P_{\text{min}})$ since $P(y | x)$ is
 290 shared across P_{maj} and P_{min} .

291 Given a level of overlap $\tau \in [0, 1]$ we denote the class of pairs of distributions $(P_{\text{maj}}, P_{\text{min}})$ with
 292 overlap at least τ by $\mathcal{P}_{\text{GS}}(\tau)$. It is easy to check that, $\mathcal{P}_{\text{GS}}(\tau) \subseteq \mathcal{P}_{\text{GS}}(0)$ at any overlap level $\tau \in [0, 1]$.

293 Next, we shall state our lower bound on the minimax excess risk that demonstrates the hardness of the
 294 group-covariate shift problem. In the theorem below $c > 0$ shall be an absolute constant independent
 295 of n_{maj} , n_{min} and τ .

296 **Theorem B.1.** *Consider the group shift setting described in Section B. Given any overlap $\tau \in [0, 1]$
 297 recall that $\mathcal{P}_{\text{GS}}(\tau)$ is the class of distributions such that $\text{Overlap}(P_{\text{maj}}, P_{\text{min}}) \geq \tau$. The minimax
 298 excess risk in this setting is lower bounded as follows:*

$$\begin{aligned} \text{Minimax Excess Risk}(\mathcal{P}_{\text{GS}}(\tau)) &= \inf_{\mathcal{A}} \sup_{(P_{\text{maj}}, P_{\text{min}}) \in \mathcal{P}_{\text{GS}}(\tau)} \text{Excess Risk}[\mathcal{A}; (P_{\text{maj}}, P_{\text{min}})] \\ &\geq \frac{1}{200(n_{\text{min}} \cdot (2 - \tau) + n_{\text{maj}} \cdot \tau)^{1/3}} \geq \frac{1}{200n_{\text{min}}^{1/3}(\rho \cdot \tau + 2)^{1/3}}, \end{aligned} \quad (4)$$

299 where $\rho = n_{\text{maj}}/n_{\text{min}} > 1$.

300 We prove this theorem in Appendix G.

301 We see that in the *low overlap* setting ($\tau \ll 1/\rho$), the minimax excess risk is lower bounded by
 302 $1/n_{\text{min}}^{1/3}$, and we are fundamentally constrained by the number of samples in minority group. To
 303 see why this is the case, consider the extreme example with $\tau = 0$ where P_a has support $[0, 0.5]$
 304 and P_b has support $[0.5, 1]$. The n_{maj} majority group samples from P_a provide information about
 305 the correct label predict in the interval $[0, 0.5]$ (the support of P_a). However, since the distribution
 306 $P(y | x)$ is 1-Lipschitz in the worst case these samples provide very limited information about the
 307 correct predictions in $[0.5, 1]$ (the support of P_b). Thus, predicting on the support of P_b requires
 308 samples from the minority group and this results in the n_{min} dependent rate. In fact, in this extreme
 309 case ($\tau = 0$) even if $n_{\text{maj}} \rightarrow \infty$, the minimax excess risk is still bounded away from zero. This
 310 intuition also carries over to the case when the overlap is small but non-zero and our lower bound
 311 shows that minority samples are much more valuable than majority samples at reducing the risk.

312 On the other hand, when the overlap is high ($\tau \gg 1/\rho$) the minimax excess risk is lower bounded
 313 by $1/(n_{\text{min}}(2 - \tau) + n_{\text{maj}}\tau)^{1/3}$ and the extra majority samples are quite beneficial. This is roughly
 314 because the supports of P_a and P_b have large overlap and hence samples from the majority group
 315 are useful in helping make predictions even in regions where P_b is large. In the extreme case when
 316 $\tau = 1$, we have that $P_a = P_b$ and therefore recover the classic i.i.d. setting with no distribution shift.
 317 Here, the lower bound scales with $1/n^{1/3}$, as one might expect.

318 Identical to the label shift case, the theorem can be extended to hold in higher dimensions with the
 319 exponents being $1/3d$ rather than $1/3$.

320 Previous work on transfer learning with covariate shift has considered other more elaborate notions
 321 of *transferability* [Kpotufe and Martinet, 2018, Hanneke and Kpotufe, 2019] than overlap between
 322 group distributions considered here. In the case of no overlap ($\tau = 0$), previous results [Kpotufe and
 323 Martinet, 2018, Theorem 1 with $\alpha = 1, \beta = 0$ and $\gamma = \infty$] yield the same lower bound of $1/n_{\text{min}}^{1/3}$.
 324 Beyond the case of no overlap ($\tau = 0$), our lower bound is key to drawing the simple conclusion that
 325 even when overlap is small between group distributions, minority samples alone dictate the rate of
 326 convergence. On the other hand, when the overlap is large our bound tells us that all samples can
 327 help reduce the risk.

328 C Minority Sample Dependence in Practice

329 Inspired by our worst-case theoretical predictions in nonparametric classification, we ask: how does
 330 the accuracy of neural network classifiers trained using robust algorithms evolve as a function of the
 331 majority and minority samples?

332 To explore this question, we conduct a small case study using the imbalanced binary CIFAR10
 333 dataset [Byrd and Lipton, 2019, Wang et al., 2022] that is constructed using the ‘‘cat’’ and ‘‘dog’’

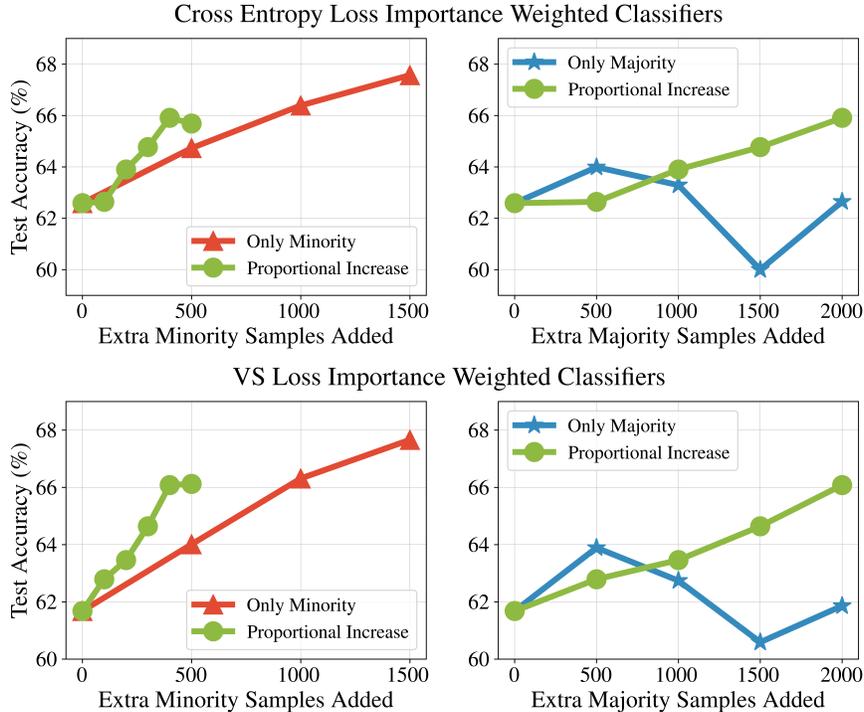


Figure 1: Convolutional neural network classifiers trained on the Imbalanced Binary CIFAR10 dataset with a 5:1 label imbalance. (Top) Models trained using the importance weighted cross entropy loss with early stopping. (Bottom) Models trained using the importance weighted VS loss [Kini et al., 2021] with early stopping. We report the average test accuracy calculated on a balanced test set over 5 random seeds. We start off with 2500 cat examples and 500 dog examples in the training dataset. We find that in accordance with our theory, for both of the classifiers adding only minority class samples (red) leads to large gain in accuracy ($\sim 6\%$), while adding majority class samples (blue) leads to little or no gain. In fact, adding majority samples sometimes hurts test accuracy due to the added bias. When we add majority and minority samples in a 5:1 ratio (green), the gain is largely due to the addition of minority samples and is only marginally higher ($< 2\%$) than adding only minority samples. The green curves correspond to the same classifiers in both the left and right panels.

334 classes. The test set consists of all of the 1000 cat and 1000 dog test examples. To form our initial
 335 train and validation sets, we take 2500 cat examples but only 500 dog examples from the official train
 336 set, corresponding to a 5:1 label imbalance. We then use 80% of those examples for training and the
 337 rest for validation. In our experiment, we either (a) add only minority samples; (b) add only majority
 338 samples; (c) add both majority and minority samples in a 5:1 ratio. We consider competitive robust
 339 classifiers proposed in the literature that are convolutional neural networks trained either by using
 340 (i) the importance weighted cross entropy loss, or (ii) the importance weighted VS loss [Kini et al.,
 341 2021]. We early stop using the importance weighted validation loss in both cases. The additional
 342 experimental details are presented in Appendix I.

343 Our results in Figure 1 are generally consistent with our theoretical predictions. By adding only
 344 minority class samples the test accuracy of both classifiers increases by a great extent (6%), while by
 345 adding only majority class samples the test accuracy remains constant or in some cases even decreases
 346 owing to the added bias of the classifiers. When we add samples to both groups proportionately, the
 347 increase in the test accuracy appears to largely to be due to the increase in the number of minority
 348 class samples and on the left panels, we see that the difference between adding only extra minority
 349 group samples (red) and both minority and majority group samples (green) is small. Thus, we find
 350 that the accuracy for these neural network classifiers is also constrained by the number of minority
 351 class samples. Similar conclusions hold for classifiers trained using the tilted loss [Li et al., 2020]
 352 and group-DRO objective [Sagawa et al., 2020] (see Appendix H).

353 **D Upper Bounds on the Excess Risk for the Undersampled Binning**
 354 **Estimator**

355 We will show that an undersampled estimator matches the rates in the previous section showing
 356 that undersampling is an optimal robustness intervention. We start by defining the undersampling
 357 procedure and the undersampling binning estimator.

358 **Undersampling Procedure.** Given training data $\mathcal{S} := \{(x_1, y_1), \dots, (x_n, y_n)\}$, generate a new
 359 undersampled dataset \mathcal{S}_{US} by

- 360 • including all n_{\min} samples from \mathcal{S}_{\min} and,
- 361 • including n_{\min} samples from \mathcal{S}_{maj} by sampling uniformly at random without replacement.

362 This procedure ensures that in the undersampled dataset \mathcal{S}_{US} , the groups are balanced, and that
 363 $|\mathcal{S}_{\text{US}}| = 2n_{\min}$.

364 The undersampling binning estimator defined next will first run this undersampling procedure to
 365 obtain \mathcal{S}_{US} and just uses these samples to output a classifier.

366 **Undersampled Binning Estimator** The undersampled binning estimator \mathcal{A}_{USB} takes as input a
 367 dataset \mathcal{S} and a positive integer K corresponding to the number of bins, and returns a classifier
 368 $\mathcal{A}_{\text{USB}}^{S,K} : [0, 1] \rightarrow \{-1, 1\}$. This estimator is defined as follows:

- 369 1. First, we compute the undersampled dataset \mathcal{S}_{US} .
- 370 2. Given this dataset \mathcal{S}_{US} , let $n_{1,j}$ be the number of points with label $+1$ that lie in the interval
 371 $I_j = [\frac{j-1}{K}, \frac{j}{K}]$. Also, define $n_{-1,j}$ analogously. Then set

$$\mathcal{A}_j = \begin{cases} 1 & \text{if } n_{1,j} > n_{-1,j}, \\ -1 & \text{otherwise.} \end{cases}$$

- 372 3. Define the classifier $\mathcal{A}_{\text{USB}}^{S,K}$ such that if $x \in I_j$ then

$$\mathcal{A}_{\text{USB}}^{S,K}(x) = \mathcal{A}_j. \tag{5}$$

373 Essentially in each bin I_j , we set the prediction to be the majority label among the samples
 374 that fall in this bin.

375 Whenever the number of bins K is clear from the context we shall denote $\mathcal{A}_{\text{USB}}^{S,K}$ by $\mathcal{A}_{\text{USB}}^S$. Below we
 376 establish upper bounds on the excess risk of this simple estimator.

377 **D.1 Label Shift Upper Bounds**

378 We now establish an upper bound on the excess risk of \mathcal{A}_{USB} in the label shift setting (see Section 2).
 379 Below we let $c, C > 0$ be absolute constants independent of problem parameters like n_{maj} and n_{\min} .

380 **Theorem D.1.** *Consider the label shift setting described in Section 2. For any $(P_{\text{maj}}, P_{\min}) \in \mathcal{P}_{\text{LS}}$
 381 the expected excess risk of the Undersampling Binning Estimator (Eq. (5)) with number of bins with
 382 $K = c\lceil n_{\min}^{1/3} \rceil$ is upper bounded by*

$$\text{Excess Risk}[\mathcal{A}_{\text{USB}}; (P_{\text{maj}}, P_{\min})] = \mathbb{E}_{\mathcal{S} \sim P_{\text{maj}}^{n_{\text{maj}}} \times P_{\min}^{n_{\min}}} [R(\mathcal{A}_{\text{USB}}^S; P_{\text{test}}) - R(f^*; P_{\text{test}})] \leq \frac{C}{n_{\min}^{1/3}}.$$

383 We prove this result in Appendix F. This upper bound combined with the lower bound in Theorem 3.1
 384 shows that an undersampling approach is minimax optimal up to constants in the presence of label
 385 shift.

386 Our analysis leaves open the possibility of better algorithms when the learner has additional infor-
 387 mation about the structure of the label shift beyond Lipschitz continuity. We also note that it is
 388 straightforward to generalize the upper bound to higher dimensions with the exponent being $1/3d$
 389 instead of $1/3$.

390 D.2 Group-Covariate Shift Upper Bounds

391 Next, we present our upper bounds on the excess risk of the undersampled binning estimator in the
 392 group-covariate shift setting (see Section B). In the theorem below, $C > 0$ is an absolute constant
 393 independent of the problem parameters n_{maj} , n_{min} and τ .

394 **Theorem D.2.** *Consider the group shift setting described in Section B. For any overlap $\tau \in [0, 1]$
 395 and for any $(P_{\text{maj}}, P_{\text{min}}) \in \mathcal{P}_{\text{GS}}(\tau)$ the expected excess risk of the Undersampling Binning Estimator
 396 (Eq. (5)) with number of bins with $K = \lceil n_{\text{min}}^{1/3} \rceil$ is*

$$\text{Excess Risk}[\mathcal{A}_{\text{USB}}; (P_{\text{maj}}, P_{\text{min}})] = \mathbb{E}_{\mathcal{S} \sim P_{\text{maj}}^{n_{\text{maj}}} \times P_{\text{min}}^{n_{\text{min}}}} [R(\mathcal{A}_{\text{USB}}^{\mathcal{S}}; P_{\text{test}}) - R(f^*; P_{\text{test}})] \leq \frac{C}{n_{\text{min}}^{1/3}}.$$

397 We provide a proof for this theorem in Appendix G. Compared to the lower bound established in
 398 Theorem B.1 which scales as $1 / ((2 - \tau)n_{\text{min}} + n_{\text{maj}}\tau)^{1/3}$, the upper bound for the undersampled
 399 binning estimator always scales with $1/n_{\text{min}}^{1/3}$ since it operates on the undersampled dataset $(\mathcal{S}_{\text{US}})$.

400 Thus, we have shown that in the absence of overlap ($\tau \ll 1/\rho = n_{\text{min}}/n_{\text{maj}}$) there is an under-
 401 sampling algorithm that is minimax optimal up to constants. However when there is high overlap
 402 ($\tau \gg 1/\rho$) there is a non-trivial gap between the upper and lower bounds:

$$\frac{\text{Upper Bound}}{\text{Lower Bound}} = c(\rho \cdot \tau + 2)^{1/3}.$$

403 Again this upper bound can be generalized to higher dimensions.

404 E Technical Tools

405 In this section we avail ourselves of some technical tools that shall be used in all of the proofs below.

406 E.1 Reduction to lower bounds over a finite class

407 The lower bound on the minimax excess risk will be established via the usual route of first identifying
 408 a “hard” finite set of problem instances and then establishing the lower bound over this finite class.
 409 One difference from the usual setup in proving such lower bounds [see Wainwright, 2019, Chapter 15]
 410 is that the training samples are drawn from an imbalanced distribution, whereas the test samples are
 411 drawn from a balanced one.

412 Let \mathcal{P} be a class of pairs of distributions, where each element $(P_{\text{maj}}, P_{\text{min}}) \in \mathcal{P}$ is a pair of dis-
 413 tributions over $[0, 1] \times \{-1, 1\}$. As before, we let P_{test} denote the uniform mixture over P_{maj}
 414 and P_{min} . We let \mathcal{V} denote a finite index set. Corresponding to each element $v \in \mathcal{V}$ there is a
 415 $P_v = (P_{v,\text{maj}}, P_{v,\text{min}}) \in \mathcal{P}$ with $P_{v,\text{test}} = (P_{v,\text{maj}} + P_{v,\text{min}})/2$. Finally, also define a pair of random
 416 variables (V, S) as follows:

- 417 1. V is a uniform random variable over the set \mathcal{V} .
- 418 2. $(S \mid V = v) \sim P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}$, is an independent draw of n_{maj} samples from $P_{v,\text{maj}}$ and
 419 n_{min} samples from $P_{v,\text{min}}$.

420 We shall let Q denote the joint distribution of the random variables (V, S) , and let Q_S denote the
 421 marginal distribution of S .

422 With this notation in place, we now present a lemma that lower bounds the minimax excess risk in
 423 terms of quantities defined over the finite class of “hard” instances P_v .

424 **Lemma E.1.** *Let the random variables (V, S) be as defined above. The minimax excess risk is lower
 425 bounded as follows:*

$$\begin{aligned} \text{Minimax Excess Risk}(\mathcal{P}) &= \inf_{\mathcal{A}} \sup_{(P_{\text{maj}}, P_{\text{min}}) \in \mathcal{P}} \mathbb{E}_{\mathcal{S} \sim P_{\text{maj}}^{n_{\text{maj}}} \times P_{\text{min}}^{n_{\text{min}}}} [R(\mathcal{A}^{\mathcal{S}}; P_{\text{test}}) - R(f^*(P_{\text{test}}); P_{\text{test}})] \\ &\geq \mathfrak{R}_{\mathcal{V}} - \mathfrak{B}_{\mathcal{V}}, \end{aligned}$$

426 where $\mathfrak{R}_{\mathcal{V}}$ and Bayes-error $\mathfrak{B}_{\mathcal{V}}$ are defined as

$$\begin{aligned} \mathfrak{R}_{\mathcal{V}} &:= \mathbb{E}_{S \sim Q_S} \left[\inf_h \mathbb{P}_{(x,y) \sim \sum_{v \in \mathcal{V}} Q(v|S) P_{v,\text{test}}} (h(x) \neq y) \right], \\ \mathfrak{B}_{\mathcal{V}} &:= \mathbb{E}_V [R(f^*(P_{V,\text{test}}); P_{V,\text{test}})]. \end{aligned}$$

427 *Proof.* By the definition of Minimax Excess Risk,

$$\begin{aligned}
\text{Minimax Excess Risk} &= \inf_{\mathcal{A}} \sup_{(P_{\text{maj}}, P_{\text{min}}) \in \mathcal{P}} \mathbb{E}_{S \sim P_{\text{maj}}^{n_{\text{maj}}} \times P_{\text{min}}^{n_{\text{min}}}} [R(\mathcal{A}^S; P_{\text{test}})] - R(f^*(P_{\text{test}}); P_{\text{test}}) \\
&\geq \inf_{\mathcal{A}} \sup_{v \in \mathcal{V}} \mathbb{E}_{S|v \sim P_{v, \text{maj}}^{n_{\text{maj}}} \times P_{v, \text{min}}^{n_{\text{min}}}} [R(\mathcal{A}^S; P_{v, \text{test}})] - R(f^*(P_{v, \text{test}}); P_{v, \text{test}}) \\
&\geq \inf_{\mathcal{A}} \mathbb{E}_V \left[\mathbb{E}_{S|V \sim P_{V, \text{maj}}^{n_{\text{maj}}} \times P_{V, \text{min}}^{n_{\text{min}}}} [R(\mathcal{A}^S; P_{V, \text{test}})] - R(f^*(P_{V, \text{test}}); P_{V, \text{test}}) \right] \\
&= \inf_{\mathcal{A}} \mathbb{E}_V \left[\mathbb{E}_{S|V \sim P_{V, \text{maj}}^{n_{\text{maj}}} \times P_{V, \text{min}}^{n_{\text{min}}}} [R(\mathcal{A}^S; P_{V, \text{test}})] \right] - \underbrace{\mathbb{E}_V [R(f^*(P_{V, \text{test}}); P_{V, \text{test}})]}_{=\mathfrak{B}_V}.
\end{aligned}$$

428 We continue lower bounding the first term as follows

$$\begin{aligned}
\inf_{\mathcal{A}} \mathbb{E}_V \left[\mathbb{E}_{S|V \sim P_{V, \text{maj}}^{n_{\text{maj}}} \times P_{V, \text{min}}^{n_{\text{min}}}} [R(\mathcal{A}^S; P_{V, \text{test}})] \right] &= \inf_{\mathcal{A}} \mathbb{E}_{(V, S) \sim \mathcal{Q}} [\mathbb{P}_{(x, y) \sim P_{V, \text{test}}} (\mathcal{A}^S(x) \neq y)] \\
&= \inf_{\mathcal{A}} \mathbb{E}_{S \sim \mathcal{Q}_S} \mathbb{E}_{V \sim \mathcal{Q}(\cdot|S)} [\mathbb{P}_{(x, y) \sim P_{V, \text{test}}} (\mathcal{A}^S(x) \neq y)] \\
&\stackrel{(i)}{\geq} \mathbb{E}_{S \sim \mathcal{Q}_S} \left[\inf_h \mathbb{E}_{V \sim \mathcal{Q}(\cdot|S)} [\mathbb{P}_{(x, y) \sim P_{V, \text{test}}} (h(x) \neq y)] \right] \\
&= \mathbb{E}_{S \sim \mathcal{Q}_S} \left[\inf_h \mathbb{P}_{(x, y) \sim \sum_{v \in \mathcal{V}} \mathcal{Q}(v|S) P_{v, \text{test}}} (h(x) \neq y) \right] \\
&= \mathfrak{R}_V,
\end{aligned}$$

429 where (i) follows since \mathcal{A}^S is a fixed classifier given the sample set S . This, combined with the
430 previous equation block completes the proof. \square

431 E.2 The Hat Function and its Properties

432 In this section, we define the *hat function* and establish some of its properties. This function will be
433 useful in defining “hard” problem instances to prove our lower bounds. Given a positive integer K
434 the hat function is defined as

$$\phi_K(x) = \begin{cases} |x + \frac{1}{4K}| - \frac{1}{4K} & \text{for } x \in \left[-\frac{1}{2K}, 0\right], \\ \frac{1}{4K} - |x - \frac{1}{4K}| & \text{for } x \in \left[0, \frac{1}{2K}\right], \\ 0 & \text{otherwise.} \end{cases} \quad (6)$$

435 When K is clear from context, we omit the subscript.

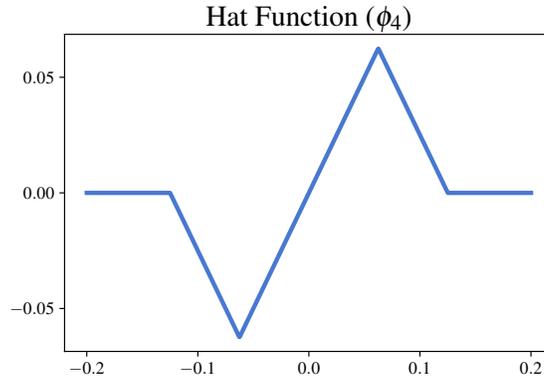


Figure 2: The hat function with $K = 4$.

436 We first notice that this function is 1-Lipschitz and odd, so

$$\int_{-\frac{1}{2K}}^{\frac{1}{2K}} \phi_K(x) dx = 0.$$

437 We also compute some other key quantities for ϕ .

438 **Lemma E.2.** For any positive integer K ,

$$\int_{-\frac{1}{2K}}^{\frac{1}{2K}} |\phi_K(x)| dx = \frac{1}{8K^2}.$$

439 *Proof.* We suppress K in the notation. We have that,

$$\int_{-\frac{1}{2K}}^{\frac{1}{2K}} |\phi(x)| dx = \int_{-\frac{1}{2K}}^0 \left| \frac{1}{4K} - \left| x + \frac{1}{4K} \right| \right| dx + \int_0^{\frac{1}{2K}} \left| \left| x - \frac{1}{4K} \right| - \frac{1}{4K} \right| dx.$$

440 The integrand $\left| \frac{1}{4K} - \left| x + \frac{1}{4K} \right| \right|$ over $x \in [-\frac{1}{2K}, 0]$ defines a triangle with base $\frac{1}{2K}$ and height $\frac{1}{4K}$,
 441 thus it has area $\frac{1}{16K^2}$. Therefore,

$$\int_{-\frac{1}{2K}}^0 \left| \frac{1}{4K} - \left| x + \frac{1}{4K} \right| \right| dx = \frac{1}{16K^2}.$$

442 The same holds for the second term. Thus, by adding them up we get that $\int_{-\frac{1}{2K}}^{\frac{1}{2K}} |\phi(x)| dx =$
 443 $\frac{1}{8K^2}$. \square

444 **Lemma E.3.** For any positive integer K ,

$$\int_0^{\frac{1}{K}} \log \left(\frac{1 + \phi_K(x - \frac{1}{2K})}{1 - \phi_K(x - \frac{1}{2K})} \right) \left(1 + \phi_K \left(x - \frac{1}{2K} \right) \right) dx \leq \frac{1}{3K^3}$$

445 and

$$\int_0^{\frac{1}{K}} \log \left(\frac{1 - \phi_K(x - \frac{1}{2K})}{1 + \phi_K(x - \frac{1}{2K})} \right) \left(1 - \phi_K \left(x - \frac{1}{2K} \right) \right) dx \leq \frac{1}{3K^3}.$$

446 *Proof.* Let us suppress K in the notation. We prove the first bound below and the second bound
 447 follows by an identical argument. We have that

$$\begin{aligned} & \int_0^{\frac{1}{K}} \log \left(\frac{1 + \phi(x - \frac{1}{2K})}{1 - \phi(x - \frac{1}{2K})} \right) \left(1 + \phi \left(x - \frac{1}{2K} \right) \right) dx \\ &= \int_{-\frac{1}{2K}}^{\frac{1}{2K}} \log \left(\frac{1 + \phi(x)}{1 - \phi(x)} \right) (1 + \phi(x)) dx \\ &= \int_0^{\frac{1}{2K}} \log \left(\frac{1 + \phi(x)}{1 - \phi(x)} \right) (1 + \phi(x)) dx + \int_{-\frac{1}{2K}}^0 \log \left(\frac{1 + \phi(x)}{1 - \phi(x)} \right) (1 + \phi(x)) dx \\ &= \int_0^{\frac{1}{2K}} \log \left(\frac{1 + \phi(x)}{1 - \phi(x)} \right) (1 + \phi(x)) dx - \int_{\frac{1}{2K}}^0 \log \left(\frac{1 + \phi(-x)}{1 - \phi(-x)} \right) (1 + \phi(-x)) dx \\ &= \int_0^{\frac{1}{2K}} \log \left(\frac{1 + \phi(x)}{1 - \phi(x)} \right) (1 + \phi(x)) dx + \int_0^{\frac{1}{2K}} \log \left(\frac{1 - \phi(x)}{1 + \phi(x)} \right) (1 - \phi(x)) dx, \end{aligned}$$

448 where the last equality follows since ϕ is an odd function. Now, we may collect the integrands to get
 449 that,

$$\begin{aligned} & \int_0^{\frac{1}{K}} \log \left(\frac{1 + \phi(x - \frac{1}{2K})}{1 - \phi(x - \frac{1}{2K})} \right) \left(1 + \phi \left(x - \frac{1}{2K} \right) \right) dx \\ &= 2 \int_0^{\frac{1}{2K}} \log \left(\frac{1 + \phi(x)}{1 - \phi(x)} \right) \phi(x) dx \\ &= 2 \int_0^{\frac{1}{2K}} \log \left(1 + \frac{2\phi(x)}{1 - \phi(x)} \right) \phi(x) dx \\ &\leq 2 \int_0^{\frac{1}{2K}} \frac{2\phi(x)^2}{1 - \phi(x)} dx, \end{aligned}$$

450 where the last inequality follows since $\log(1+x) \leq x$ for all x . Now we observe that $\phi(x) \leq x \leq \frac{1}{2}$
 451 for $x \in [0, \frac{1}{2K}]$, and in particular, $\frac{1}{1-\phi(x)} \leq 2$. Thus,

$$\begin{aligned} & \int_0^{\frac{1}{2K}} \log \left(\frac{1 + \phi(x - \frac{1}{2K})}{1 - \phi(x - \frac{1}{2K})} \right) \left(1 + \phi \left(x - \frac{1}{2K} \right) \right) dx \\ & \leq 8 \int_0^{\frac{1}{2K}} \phi(x)^2 dx \\ & \leq 8 \int_0^{\frac{1}{2K}} x^2 dx \\ & = \frac{1}{3K^3}. \end{aligned}$$

452 This proves the first bound. The second bound follows analogously. \square

453 F Proofs in the Label Shift Setting

454 Throughout this section we operate in the label shift setting (see Section 2).

455 First, in Appendix F.1 through a sequence of lemmas we prove the minimax lower bound Theorem 3.1.
 456 Next, in Appendix F.2 we prove Theorem D.1 which is an upper bound on the excess risk of the
 457 undersampled binning estimator (see Eq. (5)) with $\lceil n_{\min} \rceil^{1/3}$ bins by invoking previous results on
 458 nonparametric density estimation [Freedman and Diaconis, 1981, Devroye and Györfi, 1985].

459 F.1 Proof of Theorem 3.1

460 In this section, we provide a proof of the minimax lower bound in the label shift setting.

461 We will proceed by constructing a class of distributions where the separation between any two
 462 distributions in the class is small enough such that it is hard to distinguish between them with finite
 463 minority class samples. In particular, we split the interval $[0, 1]$ into sub-intervals and each class
 464 distribution on each sub-interval either has slightly more probability mass on the left side of the
 465 sub-interval, on the right, or completely uniform. Since the minority class sample size is limited, no
 466 classifier will be able to tell which distribution the minority class is generated from, and hence will
 467 suffer high excess risk.

468 We construct the “hard” set of distributions as follows. Fix K to be an integer that will be specified
 469 in the sequel as a function of n_{\min} . Let the index set be $\mathcal{V} = \{-1, 0, 1\}^K \times \{-1, 0, 1\}^K$. For
 470 $v \in \mathcal{V}$, we will let $v_1 \in \{-1, 0, 1\}^K$ be the first K coordinates and $v_{-1} \in \{-1, 0, 1\}^K$ be the last K
 471 coordinates. That is, $v = (v_1, v_{-1})$.

472 For every $v \in \mathcal{P}$ we shall define pair of class-conditional distributions $P_{v,1}$ and $P_{v,-1}$ as follows: for
 473 $x \in I_j = [\frac{j-1}{K}, \frac{j}{K}]$,

$$\begin{aligned} P_{v,1}(x) &= 1 + v_{1,j} \phi \left(x - \frac{j+1/2}{K} \right) \\ P_{v,-1}(x) &= 1 + v_{-1,j} \phi \left(x - \frac{j+1/2}{K} \right), \end{aligned}$$

474 where ϕ is defined in Eq. 6. Notice that $P_{v,1}$ only depends on v_1 while $P_{v,-1}$ only depends on v_{-1} .
 475 We continue to define

$$\begin{aligned} P_{v,\text{maj}}(x, y) &= P_{v,1}(x) \mathbf{1}(y = 1) \\ P_{v,\text{min}}(x, y) &= P_{v,-1}(x) \mathbf{1}(y = -1), \end{aligned}$$

476 and

$$P_{v,\text{test}}(x, y) = \frac{P_{v,\text{maj}}(x, y) + P_{v,\text{min}}(x, y)}{2} = \frac{P_{v,1}(x) \mathbf{1}(y = 1) + P_{v,-1}(x) \mathbf{1}(y = -1)}{2}.$$

477 Observe that in the test distribution it is equally likely for the label to be $+1$ or -1 .

478 Recall that as described in Section E.1, V shall be a uniform random variable over \mathcal{V} and $S \mid V \sim$
 479 $P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}$. We shall let Q denote the joint distribution of (V, S) and let Q_S denote the marginal
 480 over S .

481 With this construction in place, we first show that the minimax excess risk is lower bounded by

482 **Lemma F.1.** *For any positive integers $K, n_{\text{maj}}, n_{\text{min}}$, the minimax excess risk is lower bounded as*
 483 *follows:*

$$\begin{aligned} & \text{Minimax Excess Risk}(\mathcal{P}_{\text{LS}}) \\ &= \inf_{\mathcal{A}} \sup_{(P_{\text{maj}}, P_{\text{min}}) \in \mathcal{P}_{\text{LS}}} \mathbb{E}_{S \sim P_{\text{maj}}^{n_{\text{maj}}} \times P_{\text{min}}^{n_{\text{min}}}} [R(\mathcal{A}^S; P_{\text{test}}) - R(f^*; P_{\text{test}})] \\ &\geq \frac{1}{36K} - \frac{1}{2} \mathbb{E}_{S \sim Q_S} \left[\text{TV} \left(\sum_{v \in \mathcal{V}} Q(v \mid S) P_{v,1}, \sum_{v \in \mathcal{V}} Q(v \mid S) P_{v,-1} \right) \right]. \end{aligned} \quad (7)$$

484 *Proof.* By invoking Lemma E.1 we get that

$$\begin{aligned} & \text{Minimax Excess Risk}(\mathcal{P}_{\text{LS}}) \\ &\geq \underbrace{\mathbb{E}_{S \sim Q_S} [\inf_h \mathbb{P}_{(x,y) \sim \sum_{v \in \mathcal{V}} Q(v \mid S) P_{v,\text{test}}} (h(x) \neq y)]}_{=: \mathfrak{R}_{\mathcal{V}}} - \underbrace{\mathbb{E}_V [R(f^*(P_{V,\text{test}}); P_{V,\text{test}})]}_{=: \mathfrak{B}_{\mathcal{V}}}. \end{aligned}$$

485 We proceed by calculating alternate expressions for $\mathfrak{R}_{\mathcal{V}}$ and $\mathfrak{B}_{\mathcal{V}}$ to get our desired lower bound on
 486 the minimax excess risk.

487 **Calculation of $\mathfrak{R}_{\mathcal{V}}$:** Immediately by Le Cam's lemma [Wainwright, 2019, Eq. 15.13], we get that

$$\begin{aligned} \mathfrak{R}_{\mathcal{V}} &= \mathbb{E}_{S \sim Q_S} \left[\inf_h \mathbb{P}_{(x,y) \sim \sum_{v \in \mathcal{V}} Q(v \mid S) P_{v,\text{test}}} (h(x) \neq y) \right] \\ &= \frac{1}{2} \mathbb{E}_{S \sim Q_S} \left[1 - \text{TV} \left(\sum_{v \in \mathcal{V}} Q(v \mid S) P_{v,1}, \sum_{v \in \mathcal{V}} Q(v \mid S) P_{v,-1} \right) \right]. \end{aligned} \quad (8)$$

488 **Calculation of $\mathfrak{B}_{\mathcal{V}}$:** Again by invoking Le Cam's lemma [Wainwright, 2019, Eq. 15.13], we get that
 489 for any class conditional distributions P_1, P_{-1} ,

$$R(f^*; P_{\text{test}}) = \frac{1}{2} - \frac{1}{2} \text{TV}(P_1, P_{-1}).$$

490 So by taking expectations, we get that

$$\mathfrak{B}_{\mathcal{V}} = \mathbb{E}_V [R(f^*(P_{V,\text{test}}); P_{V,\text{test}})] = \mathbb{E}_V \left[\frac{1}{2} - \frac{1}{2} \text{TV}(P_{V,1}, P_{V,-1}) \right]. \quad (9)$$

491 We now compute $\mathbb{E}_V [\text{TV}(P_{V,1}, P_{V,-1})]$ as follows:

$$\begin{aligned} \mathbb{E}_V [\text{TV}(P_{V,1}, P_{V,-1})] &= \frac{1}{2} \mathbb{E}_V \left[\int_{x=0}^1 |P_{V,1}(x) - P_{V,-1}(x)| dx \right] \\ &= \frac{1}{2} \mathbb{E}_V \left[\sum_{j=1}^K \int_{\frac{j-1}{K}}^{\frac{j}{K}} |V_{1,j} - V_{-1,j}| \left| \phi \left(x - \frac{j+1/2}{K} \right) \right| dx \right] \\ &= \frac{1}{2} \sum_{j=1}^K \mathbb{E}_V \left[\int_{\frac{j-1}{K}}^{\frac{j}{K}} |V_{1,j} - V_{-1,j}| \left| \phi \left(x - \frac{j+1/2}{K} \right) \right| dx \right] \\ &\stackrel{(i)}{=} \frac{1}{16K^2} \sum_{j=1}^K \mathbb{E}_V [|V_{1,j} - V_{-1,j}|], \end{aligned}$$

492 where (i) follows by Lemma E.2. Observe that $V_{1,j}, V_{-1,j}$ are independent uniform random variables
 493 on $\{-1, 0, 1\}$, it is therefore straightforward to compute that

$$\mathbb{E}_V [|V_{1,j} - V_{-1,j}|] = \frac{8}{9}.$$

494 This yields that

$$\mathbb{E}_V [\text{TV}(\mathbb{P}_{V,1}, \mathbb{P}_{V,-1})] = \frac{1}{18K}.$$

495 Plugging this into Eq. (9) allows us to conclude that

$$\mathfrak{B}_V = \mathbb{E}_V [R(f^*(\mathbb{P}_{V,\text{test}}); \mathbb{P}_{V,\text{test}})] = \frac{1}{2} \left(1 - \frac{1}{18K}\right). \quad (10)$$

496 Combining Eqs. (8) and (10) establishes the claimed result.

497

□

498 In light of this previous lemma we now aim to upper bound the expected total variation distance in
499 Eq. (7).

500 **Lemma F.2.** *Suppose that v is drawn uniformly from the set $\{-1, 1\}^K$, and that $S \mid v$ is drawn from
501 $\mathbb{P}_{v,\text{maj}}^{n_{\text{maj}}} \times \mathbb{P}_{v,\text{min}}^{n_{\text{min}}}$ then,*

$$\mathbb{E}_S \left[\text{TV} \left(\sum_{v \in \mathcal{V}} \mathbb{Q}(v \mid S) \mathbb{P}_{v,1}, \sum_{v \in \mathcal{V}} \mathbb{Q}(v \mid S) \mathbb{P}_{v,-1} \right) \right] \leq \frac{1}{18K} - \frac{1}{144K} \exp\left(-\frac{n_{\text{min}}}{3K^3}\right).$$

502 *Proof.* Let $\psi := \mathbb{E}_S [\text{TV}(\sum_{v \in \mathcal{V}} \mathbb{Q}(v \mid S) \mathbb{P}_{v,1}, \sum_{v \in \mathcal{V}} \mathbb{Q}(v \mid S) \mathbb{P}_{v,-1})]$. Then,

$$\begin{aligned} \psi &= \mathbb{E}_S \left[\text{TV} \left(\sum_{v \in \mathcal{V}} \mathbb{Q}(v \mid S) \mathbb{P}_{v,1}, \sum_{v \in \mathcal{V}} \mathbb{Q}(v \mid S) \mathbb{P}_{v,-1} \right) \right] \\ &= \frac{1}{2} \mathbb{E}_S \left[\int_{x=0}^1 \left| \sum_{v \in \mathcal{V}} \mathbb{Q}(v \mid S) (\mathbb{P}_{v,1}(x) - \mathbb{P}_{v,-1}(x)) \right| dx \right] \\ &= \frac{1}{2} \mathbb{E}_S \left[\sum_{j=1}^K \int_{x=\frac{j-1}{K}}^{\frac{j}{K}} \left| \sum_{v \in \mathcal{V}} \mathbb{Q}(v \mid S) (\mathbb{P}_{v,1}(x) - \mathbb{P}_{v,-1}(x)) \right| dx \right] \\ &= \frac{1}{2} \mathbb{E}_S \left[\sum_{j=1}^K \int_{x=\frac{j-1}{K}}^{\frac{j}{K}} \left| \sum_{v \in \mathcal{V}} \mathbb{Q}(v \mid S) (v_{1,j} - v_{-1,j}) \phi\left(x - \frac{j+1/2}{K}\right) \right| dx \right], \end{aligned}$$

503 where the last equality is by the definition of $\mathbb{P}_{v,1}$ and $\mathbb{P}_{v,-1}$. Continuing we get that,

$$\begin{aligned} \psi &= \frac{1}{2} \sum_{j=1}^K \left[\int_{x=\frac{j-1}{K}}^{\frac{j}{K}} \left| \phi\left(x - \frac{j+1/2}{K}\right) \right| dx \right] \mathbb{E}_S \left[\left| \sum_{v \in \mathcal{V}} \mathbb{Q}(v \mid S) (v_{1,j} - v_{-1,j}) \right| \right] \\ &\stackrel{(i)}{=} \frac{1}{16K^2} \mathbb{E}_S \left[\sum_{j=1}^K \left| \sum_{v \in \mathcal{V}} \mathbb{Q}(v \mid S) (v_{1,j} - v_{-1,j}) \right| \right] \\ &= \frac{1}{16K^2} \sum_{j=1}^K \int \left| \sum_{v \in \mathcal{V}} \mathbb{Q}(v \mid S) (v_{1,j} - v_{-1,j}) \right| d\mathbb{Q}_S(S) \\ &= \frac{1}{16K^2} \sum_{j=1}^K \int \left| \sum_{v \in \mathcal{V}} \mathbb{Q}(v, S) (v_{1,j} - v_{-1,j}) \right| dS \\ &\stackrel{(ii)}{=} \frac{1}{16K^2 |\mathcal{V}|} \sum_{j=1}^K \int \left| \sum_{v \in \mathcal{V}} \mathbb{Q}(S \mid v) (v_{1,j} - v_{-1,j}) \right| dS, \end{aligned}$$

504 where (i) follows by the calculation in Lemma E.2 and (ii) follows since v is a uniform random
505 variable over the set \mathcal{V} .

506 The distributions $P_{v,1}$ and $P_{v,-1}$ are symmetrically defined over all intervals $I_j = [\frac{j-1}{K}, \frac{j}{K}]$, and
 507 hence all of the summands in the RHS above are equal. Thus,

$$\psi = \frac{1}{16K|\mathcal{V}|} \int \left| \sum_{v \in \mathcal{V}} Q(S | v)(v_{1,1} - v_{-1,1}) \right| dS. \quad (11)$$

508 Before we continue further, let us define

$$\mathcal{V}^+ = \{v \in \mathcal{V} \mid v_{1,1} > v_{-1,1}\}.$$

509 For every $v \in \mathcal{V}^+$, let $\tilde{v} \in \mathcal{V}$ be such that is the same as v on all coordinates, except $\tilde{v}_{1,1} = -v_{1,1}$
 510 and $\tilde{v}_{-1,1} = -v_{-1,1}$. Then continuing from Eq. (11) we find that,

$$\begin{aligned} \psi &\stackrel{(i)}{=} \frac{1}{16K|\mathcal{V}|} \int \left| \sum_{v \in \mathcal{V}^+} (v_{1,1} - v_{-1,1})(Q(S | v) - Q(S | \tilde{v})) \right| dS \\ &\stackrel{(ii)}{\leq} \frac{1}{16K|\mathcal{V}|} \int \sum_{v \in \mathcal{V}^+} (v_{1,1} - v_{-1,1}) |Q(S | v) - Q(S | \tilde{v})| dS \\ &= \frac{1}{16K|\mathcal{V}|} \sum_{v \in \mathcal{V}^+} (v_{1,1} - v_{-1,1}) \int |Q(S | v) - Q(S | \tilde{v})| dS \\ &= \frac{1}{8K|\mathcal{V}|} \underbrace{\sum_{v \in \mathcal{V}^+} (v_{1,1} - v_{-1,1}) \text{TV}(Q(S | v), Q(S | \tilde{v}))}_{=: \Xi}, \end{aligned} \quad (12)$$

511 where (i) we use the definition of \mathcal{V}^+ and \tilde{v} , (ii) follows since $v_{1,1} > v_{-1,1}$ for $v \in \mathcal{V}^+$.

512 Now we further partition \mathcal{V}^+ into 3 sets $\mathcal{V}^{(1,0)}$, $\mathcal{V}^{(0,-1)}$, $\mathcal{V}^{(1,-1)}$ as follows

$$\begin{aligned} \mathcal{V}^{(1,0)} &= \{v \in \mathcal{V} \mid v_{1,1} = 1, v_{-1,1} = 0\}, \\ \mathcal{V}^{(0,-1)} &= \{v \in \mathcal{V} \mid v_{1,1} = 0, v_{-1,1} = -1\}, \\ \mathcal{V}^{(1,-1)} &= \{v \in \mathcal{V} \mid v_{1,1} = 1, v_{-1,1} = -1\}. \end{aligned}$$

513 Note that $Q(S | v) = P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}$, and therefore

$$\begin{aligned} \Xi &= \sum_{v \in \mathcal{V}^+} (v_{1,1} - v_{-1,1}) \text{TV} \left(P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}, P_{\tilde{v},\text{maj}}^{n_{\text{maj}}} \times P_{\tilde{v},\text{min}}^{n_{\text{min}}} \right) \\ &\stackrel{(i)}{=} \sum_{v \in \mathcal{V}^{(1,0)}} \text{TV} \left(P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}, P_{\tilde{v},\text{maj}}^{n_{\text{maj}}} \times P_{\tilde{v},\text{min}}^{n_{\text{min}}} \right) \\ &\quad + \sum_{v \in \mathcal{V}^{(0,-1)}} \text{TV} \left(P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}, P_{\tilde{v},\text{maj}}^{n_{\text{maj}}} \times P_{\tilde{v},\text{min}}^{n_{\text{min}}} \right) \\ &\quad + 2 \sum_{v \in \mathcal{V}^{(1,-1)}} \text{TV} \left(P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}, P_{\tilde{v},\text{maj}}^{n_{\text{maj}}} \times P_{\tilde{v},\text{min}}^{n_{\text{min}}} \right), \end{aligned} \quad (13)$$

514 where (i) follows since $v_1, v_{-1} \in \{-1, 0, 1\}^K$ and by the definition of the sets $\mathcal{V}^{(1,0)}$, $\mathcal{V}^{(0,-1)}$ and
 515 $\mathcal{V}^{(1,-1)}$.

516 Now by the Bretagnolle–Huber inequality [see Canonne, 2022, Corollary 4],

$$\begin{aligned} \text{TV} \left(P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}, P_{\tilde{v},\text{maj}}^{n_{\text{maj}}} \times P_{\tilde{v},\text{min}}^{n_{\text{min}}} \right) &= \text{TV} \left(P_{\tilde{v},\text{maj}}^{n_{\text{maj}}} \times P_{\tilde{v},\text{min}}^{n_{\text{min}}}, P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}} \right) \\ &\leq 1 - \frac{1}{2} \exp \left(-\text{KL} \left(P_{\tilde{v},\text{maj}}^{n_{\text{maj}}} \times P_{\tilde{v},\text{min}}^{n_{\text{min}}} \parallel P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}} \right) \right), \end{aligned}$$

517 where we flip the arguments in the first step for simplicity later.

518 Next, by the chain rule for KL-divergence, we have that

$$\text{KL}(P_{\tilde{v},\text{maj}}^{n_{\text{maj}}} \times P_{\tilde{v},\text{min}}^{n_{\text{min}}} \parallel P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}) = n_{\text{maj}} \text{KL}(P_{\tilde{v},\text{maj}} \parallel P_{v,\text{maj}}) + n_{\text{min}} \text{KL}(P_{\tilde{v},\text{min}} \parallel P_{v,\text{min}}).$$

519 Using these, let us upper bound the first term in Eq. (13) corresponding to $v \in \mathcal{V}^{(0,-1)}$. For
520 $v \in \mathcal{V}^{(0,-1)}$, notice that $\text{KL}(\mathbb{P}_{\tilde{v},\text{maj}}\|\mathbb{P}_{v,\text{maj}}) = 0$ since $v_{1,j} = \tilde{v}_{1,j}$ for all $j \in \{1, \dots, K\}$. For the
521 second term, $\text{KL}(\mathbb{P}_{\tilde{v},\text{min}}\|\mathbb{P}_{v,\text{min}})$, only $v_{1,1}$ and $\tilde{v}_{1,1}$ differ, so

$$\begin{aligned} \text{KL}(\mathbb{P}_{\tilde{v},\text{min}}\|\mathbb{P}_{v,\text{min}}) &= \int_0^1 \mathbb{P}_{v,-1}(x) \log \left(\frac{\mathbb{P}_{v,-1}(x)}{\mathbb{P}_{\tilde{v},-1}(x)} \right) dx \\ &= \int_0^{\frac{1}{K}} \log \left(\frac{1 + \phi_K(x - \frac{1}{2K})}{1 - \phi_K(x - \frac{1}{2K})} \right) \left(1 + \phi_K \left(x - \frac{1}{2K} \right) \right) dx \\ &\leq \frac{1}{3K^3}, \end{aligned}$$

522 where the last inequality is a result of the calculation in Lemma E.3.

523 Therefore, we get

$$\sum_{v \in \mathcal{V}^{(0,-1)}} \text{TV} \left(\mathbb{P}_{v,\text{maj}}^{n_{\text{maj}}} \times \mathbb{P}_{v,\text{min}}^{n_{\text{min}}}, \mathbb{P}_{\tilde{v},\text{maj}}^{n_{\text{maj}}} \times \mathbb{P}_{\tilde{v},\text{min}}^{n_{\text{min}}} \right) \leq 9^{K-1} \left(1 - \frac{1}{2} \exp \left(-\frac{n_{\text{min}}}{3K^3} \right) \right).$$

524 For the terms in Eq. (13) corresponding to $\mathcal{V}^{(0,-1)}$, $\mathcal{V}^{(1,-1)}$, we simply take the trivial bound to get

$$\begin{aligned} \sum_{v \in \mathcal{V}^{(0,-1)}} \text{TV} \left(\mathbb{P}_{v,\text{maj}}^{n_{\text{maj}}} \times \mathbb{P}_{v,\text{min}}^{n_{\text{min}}}, \mathbb{P}_{\tilde{v},\text{maj}}^{n_{\text{maj}}} \times \mathbb{P}_{\tilde{v},\text{min}}^{n_{\text{min}}} \right) &\leq 9^{K-1}, \\ \sum_{v \in \mathcal{V}^{(1,-1)}} \text{TV} \left(\mathbb{P}_{v,\text{maj}}^{n_{\text{maj}}} \times \mathbb{P}_{v,\text{min}}^{n_{\text{min}}}, \mathbb{P}_{\tilde{v},\text{maj}}^{n_{\text{maj}}} \times \mathbb{P}_{\tilde{v},\text{min}}^{n_{\text{min}}} \right) &\leq 9^{K-1}. \end{aligned}$$

525 Plugging these bounds into Eq. (13) we get that,

$$\Xi \leq 4 \cdot 9^{K-1} - \frac{9^{K-1}}{2} \exp \left(-\frac{n_{\text{min}}}{3K^3} \right).$$

526 Now using this bound on Ξ in Eq. (12) and observing that $|\mathcal{V}| = 9^K$, we get that,

$$\begin{aligned} \psi &= \mathbb{E}_S \left[\text{TV} \left(\sum_{v \in \mathcal{V}} Q(v | S) P_{v,1}, \sum_{v \in \mathcal{V}} Q(v | S) P_{v,-1} \right) \right] \\ &\leq \frac{1}{8 \cdot 9^K K} \left(4 \cdot 9^{K-1} - \frac{9^{K-1}}{2} \exp \left(-\frac{n_{\text{min}}}{3K^3} \right) \right) \\ &= \frac{1}{18K} - \frac{1}{144K} \exp \left(-\frac{n_{\text{min}}}{3K^3} \right), \end{aligned}$$

527 completing the proof. \square

528 Finally, we combine Lemma F.1 and Lemma F.2 to establish the minimax lower bound in this label
529 shift setting. We recall the statement of the theorem here.

530 **Theorem 3.1.** *Let \mathcal{P}_{LS} be the class of pairs of distributions $(\mathbb{P}_{\text{maj}}, \mathbb{P}_{\text{min}})$ that satisfy the label-shift
531 assumptions. The minimax excess risk over this class is lower bounded as follows:*

$$\text{Minimax Excess Risk}(\mathcal{P}_{\text{LS}}) = \inf_{\mathcal{A}} \sup_{(\mathbb{P}_{\text{maj}}, \mathbb{P}_{\text{min}}) \in \mathcal{P}_{\text{LS}}} \text{Excess Risk}[\mathcal{A}; (\mathbb{P}_{\text{maj}}, \mathbb{P}_{\text{min}})] \geq \frac{1}{600} \frac{1}{n_{\text{min}}^{1/3}}. \quad (3)$$

532 *Proof.* By Lemma F.1 we know that,

$$\text{Minimax Excess Risk}(\mathcal{P}_{\text{LS}}) \geq \frac{1}{36K} - \frac{1}{2} \mathbb{E}_{S \sim Q_S} \left[\text{TV} \left(\sum_{v \in \mathcal{V}} Q(v | S) P_{v,1}, \sum_{v \in \mathcal{V}} Q(v | S) P_{v,-1} \right) \right].$$

533 Next by the calculation in Lemma F.2 we have that

$$\begin{aligned} \text{Minimax Excess Risk}(\mathcal{P}_{\text{LS}}) &\geq \frac{1}{36K} - \frac{1}{2} \left(\frac{1}{18K} - \frac{1}{144K} \exp \left(-\frac{n_{\text{min}}}{3K^3} \right) \right) \\ &= \frac{1}{288K} \exp \left(-\frac{n_{\text{min}}}{3K^3} \right). \end{aligned}$$

534 Setting $K = \lceil n_{\min}^{1/3} \rceil$ yields the following

$$\begin{aligned}
\text{Minimax Excess Risk}(\mathcal{P}_{\text{LS}}) &\geq \frac{1}{288 \lceil n_{\min}^{1/3} \rceil} \exp\left(-\frac{n_{\min}}{3 \lceil n_{\min}^{1/3} \rceil^3}\right) \\
&\geq \frac{\exp\left(-\frac{n_{\min}}{3 \lceil n_{\min}^{1/3} \rceil^3}\right)}{288} \frac{n_{\min}^{1/3}}{\lceil n_{\min}^{1/3} \rceil} \frac{1}{n_{\min}^{1/3}} \\
&\stackrel{(i)}{\geq} \frac{0.7 \exp\left(-\frac{1}{3}\right)}{288} \frac{1}{n_{\min}^{1/3}} \\
&\geq \frac{1}{600} \frac{1}{n_{\min}^{1/3}},
\end{aligned}$$

535 where (i) follows since $n_{\min}^{1/3} / \lceil n_{\min}^{1/3} \rceil \geq 0.7$ for $n_{\min} \geq 1$. \square

536 F.2 Proof of Theorem D.1

537 In this section, we derive an upper bound on the excess risk of the undersampled binning estimator
538 \mathcal{A}_{USB} (Eq. (5)) in the label shift setting. Recall that given a dataset \mathcal{S} this estimator first calculates
539 the undersampled dataset \mathcal{S}_{US} , where the number of points from the minority group (n_{\min}) is equal to
540 the number of points from the majority group (n_{\min}), and the size of the dataset is $2n_{\min}$. Throughout
541 this section, $(P_{\text{maj}}, P_{\text{min}})$ shall be an arbitrary element of \mathcal{P}_{LS} .

542 To bound the excess risk of the undersampling algorithm, we will relate it to density estimation.

543 Recall that $n_{1,j}$ denotes the number of points in \mathcal{S}_{US} with label +1 that lie in I_j , and $n_{-1,j}$ is defined
544 analogously.

545 Given a positive integer K , for $x \in I_j = [\frac{j-1}{K}, \frac{j}{K}]$, by the definition of the undersampled binning
546 estimator (Eq. (5))

$$\mathcal{A}_{\text{USB}}^{\mathcal{S}}(x) = \begin{cases} 1 & \text{if } n_{1,j} > n_{-1,j}, \\ -1 & \text{otherwise.} \end{cases}$$

547 Recall that since we have undersampled, $\sum_j n_{1,j} = \sum_j n_{-1,j} = n_{\min}$. Therefore, define the simple
548 histogram estimators for $P_1(x) = P(x | y = 1)$ and $P_{-1}(x) = P(x | y = -1)$ as follows: for
549 $x \in I_j$,

$$\hat{P}_1^{\mathcal{S}}(x) := \frac{n_{1,j}}{K n_{\min}} \quad \text{and} \quad \hat{P}_{-1}^{\mathcal{S}}(x) := \frac{n_{-1,j}}{K n_{\min}}.$$

550 With this histogram estimator in place, we may define an estimator for $\eta(x) := P_{\text{test}}(y = 1 | x)$ as
551 follows,

$$\hat{\eta}^{\mathcal{S}}(x) := \frac{\hat{P}_1^{\mathcal{S}}(x)}{\hat{P}_1^{\mathcal{S}}(x) + \hat{P}_{-1}^{\mathcal{S}}(x)}.$$

552 Observe that, for $x \in I_j$

$$\hat{\eta}^{\mathcal{S}}(x) > 1/2 \iff n_{1,j} > n_{-1,j} \iff \mathcal{A}_{\text{USB}}^{\mathcal{S}}(x) = 1.$$

553 Defining an estimator $\hat{\eta}^{\mathcal{S}}$ for the $P_{\text{test}}(y = 1 | x)$ in this way will allow us to relate the excess risk of
554 \mathcal{A}_{USB} to the estimation error in $\hat{P}_1^{\mathcal{S}}$ and $\hat{P}_{-1}^{\mathcal{S}}$.

555 Before proving the theorem we restate it here.

556 **Theorem D.1.** *Consider the label shift setting described in Section 2. For any $(P_{\text{maj}}, P_{\text{min}}) \in \mathcal{P}_{\text{LS}}$
557 the expected excess risk of the Undersampling Binning Estimator (Eq. (5)) with number of bins with
558 $K = c \lceil n_{\min}^{1/3} \rceil$ is upper bounded by*

$$\text{Excess Risk}[\mathcal{A}_{\text{USB}}; (P_{\text{maj}}, P_{\text{min}})] = \mathbb{E}_{\mathcal{S} \sim P_{\text{maj}}^{n_{\text{maj}}} \times P_{\text{min}}^{n_{\text{min}}}} [R(\mathcal{A}_{\text{USB}}^{\mathcal{S}}; P_{\text{test}}) - R(f^*; P_{\text{test}})] \leq \frac{C}{n_{\min}^{1/3}}.$$

559 *Proof.* By the definition of the excess risk

$$\text{Excess Risk}[\mathcal{A}_{\text{USB}}; (\mathbf{P}_{\text{maj}}, \mathbf{P}_{\text{min}})] := \mathbb{E}_{\mathcal{S} \sim \mathcal{P}_{\text{maj}}^{n_{\text{maj}}} \times \mathcal{P}_{\text{min}}^{n_{\text{min}}}} [R(\mathcal{A}_{\text{USB}}^{\mathcal{S}}; \mathbf{P}_{\text{test}}) - R(f^*; \mathbf{P}_{\text{test}})].$$

560 By invoking [Wasserman, 2019, Theorem 1] we may upper bound the excess risk given a draw of \mathcal{S}
561 by

$$R(\mathcal{A}_{\text{USB}}^{\mathcal{S}}; \mathbf{P}_{\text{test}}) - R(f^*; \mathbf{P}_{\text{test}}) \leq 2 \int |\hat{\eta}^{\mathcal{S}}(x) - \eta(x)| \mathbf{P}_{\text{test}}(x) dx.$$

562 Continuing using the definition of $\hat{\eta}^{\mathcal{S}}$ above and because $\eta = \mathbf{P}_1 / (\mathbf{P}_1 + \mathbf{P}_{-1})$ we have that,

$$\begin{aligned} & R(\mathcal{A}_{\text{USB}}^{\mathcal{S}}; \mathbf{P}_{\text{test}}) - R(f^*; \mathbf{P}_{\text{test}}) \\ &= 2 \int_0^1 \left| \frac{\hat{\mathbf{P}}_1^{\mathcal{S}}(x)}{\hat{\mathbf{P}}_1^{\mathcal{S}}(x) + \hat{\mathbf{P}}_{-1}^{\mathcal{S}}(x)} - \frac{\mathbf{P}_1(x)}{\mathbf{P}_1(x) + \mathbf{P}_{-1}(x)} \right| \left(\frac{\mathbf{P}_1(x) + \mathbf{P}_{-1}(x)}{2} \right) dx \\ &= \int_0^1 \left| \left(\frac{\mathbf{P}_1(x) + \mathbf{P}_{-1}(x)}{\hat{\mathbf{P}}_1^{\mathcal{S}}(x) + \hat{\mathbf{P}}_{-1}^{\mathcal{S}}(x)} \right) \hat{\mathbf{P}}_1^{\mathcal{S}}(x) - \mathbf{P}_1(x) \right| dx \\ &\stackrel{(i)}{\leq} \int_0^1 \left| \hat{\mathbf{P}}_1^{\mathcal{S}}(x) - \mathbf{P}_1(x) \right| dx + \int_0^1 \left| \frac{\mathbf{P}_1(x) + \mathbf{P}_{-1}(x)}{\hat{\mathbf{P}}_1^{\mathcal{S}}(x) + \hat{\mathbf{P}}_{-1}^{\mathcal{S}}(x)} - 1 \right| \hat{\mathbf{P}}_1^{\mathcal{S}}(x) dx \\ &= \int_0^1 \left| \hat{\mathbf{P}}_1^{\mathcal{S}}(x) - \mathbf{P}_1(x) \right| dx + \int_0^1 \left| \hat{\mathbf{P}}_1^{\mathcal{S}}(x) + \hat{\mathbf{P}}_{-1}^{\mathcal{S}}(x) - \mathbf{P}_1(x) - \mathbf{P}_{-1}(x) \right| \frac{\hat{\mathbf{P}}_1^{\mathcal{S}}(x)}{\hat{\mathbf{P}}_1^{\mathcal{S}}(x) + \hat{\mathbf{P}}_{-1}^{\mathcal{S}}(x)} dx \\ &\leq 2 \int_0^1 \left| \hat{\mathbf{P}}_1^{\mathcal{S}}(x) - \mathbf{P}_1(x) \right| dx + \int_0^1 \left| \hat{\mathbf{P}}_{-1}^{\mathcal{S}}(x) - \mathbf{P}_{-1}(x) \right| dx \\ &\stackrel{(ii)}{\leq} 2 \sqrt{\int_0^1 \left(\hat{\mathbf{P}}_1^{\mathcal{S}}(x) - \mathbf{P}_1(x) \right)^2 dx} + \sqrt{\int_0^1 \left(\hat{\mathbf{P}}_{-1}^{\mathcal{S}}(x) - \mathbf{P}_{-1}(x) \right)^2 dx}, \end{aligned}$$

563 where (i) follows by the triangle inequality, (ii) is by the Cauchy–Schwarz inequality.

564 Taking expectation over the samples \mathcal{S} and by invoking Jensen’s inequality we find that,

$$\begin{aligned} & \text{Excess Risk}(\mathcal{A}^{\mathcal{S}}; (\mathbf{P}_{\text{maj}}, \mathbf{P}_{\text{min}})) \\ &= \mathbb{E}_{\mathcal{S}} [R(\mathcal{A}_{\text{USB}}^{\mathcal{S}}; \mathbf{P}_{\text{test}}) - R(f^*; \mathbf{P}_{\text{test}})] \\ &\leq 2 \sqrt{\mathbb{E}_{\mathcal{S}} \left[\int \left(\hat{\mathbf{P}}_1^{\mathcal{S}}(x) - \mathbf{P}_1(x) \right)^2 dx \right]} + \sqrt{\mathbb{E}_{\mathcal{S}} \left[\int \left(\hat{\mathbf{P}}_{-1}^{\mathcal{S}}(x) - \mathbf{P}_{-1}(x) \right)^2 dx \right]}. \end{aligned}$$

565 We note that $\hat{\mathbf{P}}_j^{\mathcal{S}}$ only depends on n_{min} i.i.d. draws from class j . Thus by [Freedman and Diaconis,
566 1981, Theorem 1.7], if $K = c \lceil n_{\text{min}} \rceil^{1/3}$ then

$$\mathbb{E}_{\mathcal{S}} \left[\int \left(\hat{\mathbf{P}}_j^{\mathcal{S}}(x) - \mathbf{P}_j(x) \right)^2 dx \right] \leq \frac{C}{n_{\text{min}}^{2/3}}.$$

567 Plugging this into the previous inequality yields the desired result. \square

568 G Proof in the Group-Covariate Shift Setting

569 Throughout this section we operate in the group-covariate shift setting (see Section B).

570 We will proceed similarly to Section F. We shall construct a family of class-conditional distributions
571 such that it will be necessary for adequate samples in each sub-interval of $[0, 1]$ to be able to learn the
572 maximally likely label in that sub-interval. On the other hand, we will construct the group-covariate
573 distributions to be separated from one another. As a consequence, sub-intervals with high probability
574 mass under the minority group distribution will have low probability mass under the majority group
575 distribution. Hence, these sub-intervals will not have enough training sample points for any classifier
576 to be able to learn the maximally likely label and as a result shall suffer high excess risk.

577 First in Appendix G.1, we prove Theorem B.1, the minimax lower bound through a sequence of
578 lemmas. Second in Appendix G.2, we prove Theorem D.2 that upper bound on the excess risk of the
579 undersampled binning estimator with $\lceil n_{\text{min}} \rceil^{1/3}$ bins.

580 **G.1 Proof of Theorem B.1**

581 In this section, we provide a proof of the minimax lower bound in the group shift setting.

582 We construct the “hard” set of distributions as follows. Let the index set be $\mathcal{V} = \{-1, 1\}^K$. For every
583 $v \in \mathcal{V}$ define a distribution as follows: for $x \in I_j = [\frac{j-1}{K}, \frac{j}{K}]$,

$$P_v(y = 1 | x) := \frac{1}{2} \left[1 + v_j \phi \left(x - \frac{j + 1/2}{K} \right) \right],$$

584 where ϕ is defined in Eq. 6. Given a $\tau \in [0, 1]$ we also construct the group distributions as follows:

$$P_a(x) = \begin{cases} 2 - \tau & \text{if } x \in [0, 0.5) \\ \tau & \text{if } x \in [0.5, 1], \end{cases}$$

585 and let

$$P_b(x) = 2 - P_a(x).$$

586 We can verify that

$$\text{Overlap}(P_a, P_b) = 1 - \text{TV}(P_a, P_b) = 1 - \frac{1}{2} \int_{x=0}^1 |P_a(x) - P_b(x)| dx = \tau.$$

587 We continue to define

$$\begin{aligned} P_{v,\text{maj}}(x, y) &= P_v(y | x) P_a(x) \\ P_{v,\text{min}}(x, y) &= P_v(y | x) P_b(x), \end{aligned}$$

588 and

$$P_{v,\text{test}}(x, y) = P_v(y | x) \left(\frac{P_a(x) + P_b(x)}{2} \right).$$

589 Observe that $(P_a(x) + P_b(x))/2 = 1$, the uniform distribution over $[0, 1]$.

590 Recall that as described in Section E.1, V shall be a uniform random variable over \mathcal{V} and $S | V \sim$
591 $P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}$. We shall let Q denote the joint distribution of (V, S) and let Q_S denote the marginal
592 over S .

593 With this construction in place, we present the following lemma that lower bounds the minimax
594 excess risk by a sum of $\exp(-\text{KL}(Q(S | v_j = 1) || Q(S | v_j = -1)))$ over the intervals. Intuitively,
595 $\text{KL}(Q(S | v_j = 1) || Q(S | v_j = -1))$ is a measure of how difficult it is to identify whether $v_j = 1$ or
596 $v_j = -1$ from the samples.

597 **Lemma G.1.** *For any positive integers $K, n_{\text{maj}}, n_{\text{min}}$ and $\tau \in [0, 1]$, the minimax excess risk is lower
598 bounded as follows:*

$$\begin{aligned} \text{Minimax Excess Risk}(\mathcal{P}_{\text{GS}}(\tau)) &= \inf_{\mathcal{A}} \sup_{(P_{\text{maj}}, P_{\text{min}}) \in \mathcal{P}_{\text{GS}}(\tau)} \mathbb{E}_{S \sim P_{\text{maj}}^{n_{\text{maj}}} \times P_{\text{min}}^{n_{\text{min}}}} [R(\mathcal{A}^S; P_{\text{test}}) - R(f^*; P_{\text{test}})] \\ &\geq \frac{1}{32K^2} \sum_{j=1}^K \exp(-\text{KL}(Q(S | v_j = 1) || Q(S | v_j = -1))). \end{aligned}$$

599 *Proof.* By invoking Lemma E.1, we know that the minimax excess risk is lower bounded by

$$\begin{aligned} &\text{Minimax Excess Risk}(\mathcal{P}_{\text{GS}}(\tau)) \\ &\geq \underbrace{\mathbb{E}_{S \sim Q_S} [\inf_h \mathbb{P}_{(x,y) \sim \sum_{v \in \mathcal{V}} Q(v|S) P_{v,\text{test}}}(h(x) \neq y)]}_{=\mathfrak{R}_{\mathcal{V}}} - \underbrace{\mathbb{E}_V [R(f^*(P_{V,\text{test}}); P_{V,\text{test}})]}_{=\mathfrak{B}_{\mathcal{V}}}, \end{aligned}$$

600 where V is a uniform random variable over the set \mathcal{V} , $S | V = v$ is a draw from $P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}$, and
601 Q denotes the joint distribution over (V, S) .

602 We shall lower bound this minimax risk in parts. First, we shall establish a lower bound on $\mathfrak{R}_{\mathcal{V}}$, and
603 then an upper bound on the Bayes risk $\mathfrak{B}_{\mathcal{V}}$.

604 **Lower bound on $\mathfrak{R}_\mathcal{V}$.** Unpacking $\mathfrak{R}_\mathcal{V}$ using its definition we get that,

$$\begin{aligned}
\mathfrak{R}_\mathcal{V} &= \mathbb{E}_{S \sim \mathcal{Q}_S} [\inf_h \mathbb{P}_{(x,y) \sim \sum_{v \in \mathcal{V}} \mathcal{Q}(v|S) \mathcal{P}_{v,\text{test}}} (h(x) \neq y)] \\
&= \mathbb{E}_{S \sim \mathcal{Q}_S} \left[\inf_h \int_0^1 \mathcal{P}_{\text{test}}(x) \mathbb{P}_{y \sim \sum_{v \in \mathcal{V}} \mathcal{Q}(v|S) \mathcal{P}_v(\cdot|x)} [h(x) \neq y] dx \right] \\
&\stackrel{(i)}{=} \mathbb{E}_{S \sim \mathcal{Q}_S} \left[\int_0^1 \mathcal{P}_{\text{test}}(x) \min \left\{ \sum_{v \in \mathcal{V}} \mathcal{Q}(v|S) \mathcal{P}_v(1|x), \sum_{v \in \mathcal{V}} \mathcal{Q}(v|S) \mathcal{P}_v(-1|x) \right\} dx \right] \\
&\stackrel{(ii)}{=} \frac{1}{2} - \mathbb{E}_{S \sim \mathcal{Q}_S} \left[\int_0^1 \mathcal{P}_{\text{test}}(x) \left| \frac{1}{2} - \sum_{v \in \mathcal{V}} \mathcal{Q}(v|S) \mathcal{P}_v(1|x) \right| dx \right] \\
&\stackrel{(iii)}{=} \frac{1}{2} - \int_0^1 \mathcal{P}_{\text{test}}(x) \mathbb{E}_{S \sim \mathcal{Q}_S} \left[\left| \frac{1}{2} - \sum_{v \in \mathcal{V}} \mathcal{Q}(v|S) \mathcal{P}_v(1|x) \right| \right] dx, \tag{14}
\end{aligned}$$

605 where (i) follows by taking h to be the pointwise minimizer over x , (ii) follows since $\mathcal{P}_v(-1|x) =$
606 $1 - \mathcal{P}_v(1|x)$ and $\min\{s, 1-s\} = (1 - |1-2s|)/2$ for all $s \in [0, 1]$, and (iii) follows by Fubini's
607 theorem which allows us to switch the order of the integrals.

608 If $x \in I_j = [\frac{j-1}{K}, \frac{j}{K}]$ for some $j \in \{1, \dots, K\}$ we let j_x denote the value of this index j . With this
609 notation in place let us continue to upper bound integrand in the second term in the RHS above as
610 follows:

$$\begin{aligned}
&\mathbb{E}_{S \sim \mathcal{Q}_S} \left[\left| \frac{1}{2} - \sum_{v \in \mathcal{V}} \mathcal{Q}(v|S) \mathcal{P}_v(1|x) \right| \right] \\
&\stackrel{(i)}{=} \mathbb{E}_{S \sim \mathcal{Q}_S} \left[\left| \phi \left(x - \frac{j_x + 1/2}{K} \right) \right| \left| \mathcal{Q}(v_{j_x} = 1|S) - \mathcal{Q}(v_{j_x} = -1|S) \right| \right] \\
&= \left| \phi \left(x - \frac{j_x + 1/2}{K} \right) \right| \mathbb{E}_{S \sim \mathcal{Q}_S} [\left| \mathcal{Q}(v_{j_x} = 1|S) - \mathcal{Q}(v_{j_x} = -1|S) \right|] \\
&\stackrel{(ii)}{=} \left| \phi \left(x - \frac{j_x + 1/2}{K} \right) \right| \mathbb{E}_{S \sim \mathcal{Q}_S} \left[\left| \frac{\mathcal{Q}(S|v_{j_x} = 1) \mathcal{Q}_V(v_{j_x} = 1)}{\mathcal{Q}_S(S)} - \frac{\mathcal{Q}(S|v_{j_x} = -1) \mathcal{Q}_V(v_{j_x} = -1)}{\mathcal{Q}_S(S)} \right| \right] \\
&\stackrel{(iii)}{=} \frac{1}{2} \left| \phi \left(x - \frac{j_x + 1/2}{K} \right) \right| \text{TV}(\mathcal{Q}(S|v_{j_x} = 1), \mathcal{Q}(S|v_{j_x} = -1)), \tag{15}
\end{aligned}$$

611 where (i) follows since $\mathcal{P}_v(1|x) = (1 + v_{j_x} \phi(x - (j_x + 1/2)/K))/2$ and by marginalizing $\mathcal{Q}(v|S)$
612 over the indices $j \neq j_x$, (ii) follows by using Bayes' rule and (iii) follows since the total-variation
613 distance is half the ℓ_1 distance. Now by the Bretagnolle–Huber inequality [see Canonne, 2022,
614 Corollary 4] we get that,

$$\begin{aligned}
&\text{TV}(\mathcal{Q}(S|v_{j_x} = 1), \mathcal{Q}(S|v_{j_x} = -1)) \\
&\leq 1 - \frac{\exp(-\text{KL}(\mathcal{Q}(S|v_{j_x} = 1) \parallel \mathcal{Q}(S|v_{j_x} = -1)))}{2}. \tag{16}
\end{aligned}$$

615 Combining Eqs. (14)-(16) we get that

$$\begin{aligned}
&\mathfrak{R}_\mathcal{V} \\
&\geq \frac{1}{2} - \frac{1}{2} \int_0^1 \mathcal{P}_{\text{test}}(x) \left| \phi \left(x - \frac{j_x + 1/2}{K} \right) \right| dx \\
&\quad + \frac{1}{4} \int_0^1 \mathcal{P}_{\text{test}}(x) \left| \phi \left(x - \frac{j_x + 1/2}{K} \right) \right| \exp(-\text{KL}(\mathcal{Q}(S|v_{j_x} = 1) \parallel \mathcal{Q}(S|v_{j_x} = -1))) dx. \tag{17}
\end{aligned}$$

616 **Upper bound on \mathfrak{B}_V :** The Bayes error is

$$\begin{aligned}
\mathfrak{B}_V &= \mathbb{E}_V[R(f^*(P_V); P_V)] \\
&= \mathbb{E}_V \left[\inf_f \mathbb{E}_{(x,y) \sim P_{v,\text{test}}} \mathbf{1}(f(x) \neq y) \right] \\
&= \mathbb{E}_V \left[\inf_f \int_{x=0}^1 \sum_{y \in \{-1,1\}} P_{\text{test}}(x) P_{V,\text{test}}(y | x) \mathbf{1}(f(x) = -y) \right] \\
&= \mathbb{E}_V \left[\int_{x=0}^1 P_{\text{test}}(x) \min_{y \in \{-1,1\}} P_{V,\text{test}}(y | x) \right] \\
&\stackrel{(i)}{=} \mathbb{E}_V \left[\frac{1}{2} \left(1 - \int_{x=0}^1 P_{\text{test}}(x) |P_{V,\text{test}}(1 | x) - P_{V,\text{test}}(-1 | x)| dx \right) \right] \\
&\stackrel{(ii)}{=} \mathbb{E}_V \left[\frac{1}{2} \left(1 - \int_{x=0}^1 P_{\text{test}}(x) \left| \phi \left(x - \frac{j_x + 1/2}{K} \right) \right| dx \right) \right] \\
&= \frac{1}{2} - \frac{1}{2} \int_{x=0}^1 P_{\text{test}}(x) \left| \phi \left(x - \frac{j_x + 1/2}{K} \right) \right| dx, \tag{18}
\end{aligned}$$

617 where (i) follows since $P_v(1 | x) = 1 - P_v(-1 | x)$ and $\min\{s, 1 - s\} = (1 - |1 - 2s|)/2$ for all
618 $s \in [0, 1]$, and (ii) follows by our construction of P_v above along with the fact that $P_v(1 | x) =$
619 $1 - P_v(-1 | x)$.

620 **Putting things together:** Combining Eqs. (17) and (18) allows us to conclude that

$$\begin{aligned}
&\text{Minimax Excess Risk}(\mathcal{P}_{\text{GS}}(\tau)) \\
&\geq \frac{1}{4} \int_0^1 P_{\text{test}}(x) \left| \phi \left(x - \frac{j_x + 1/2}{K} \right) \right| \exp(-\text{KL}(\mathbb{Q}(S | v_{j_x} = 1) \| \mathbb{Q}(S | v_{j_x} = -1))) dx \\
&= \frac{1}{4} \sum_{j=1}^K \int_{\frac{j-1}{K}}^{\frac{j}{K}} P_{\text{test}}(x) \left| \phi \left(x - \frac{j + 1/2}{K} \right) \right| \exp(-\text{KL}(\mathbb{Q}(S | v_j = 1) \| \mathbb{Q}(S | v_j = -1))) dx \\
&= \frac{1}{4} \sum_{j=1}^K \exp(-\text{KL}(\mathbb{Q}(S | v_j = 1) \| \mathbb{Q}(S | v_j = -1))) \left[\int_{\frac{j-1}{K}}^{\frac{j}{K}} P_{\text{test}}(x) \left| \phi \left(x - \frac{j + 1/2}{K} \right) \right| dx \right] \\
&\stackrel{(i)}{=} \frac{1}{32K^2} \sum_{j=1}^K \exp(-\text{KL}(\mathbb{Q}(S | v_j = 1) \| \mathbb{Q}(S | v_j = -1))),
\end{aligned}$$

621 where (i) follows by using Lemma E.2 along with the fact that $P_{\text{test}}(x) = 1$ in our construction to
622 show that the integral in the square brackets is equal to $1/8K^2$. This proves the result. \square

623 The next lemma upper bounds the KL divergence between $\mathbb{Q}(S | v_j = 1)$ and $\mathbb{Q}(S | v_j = -1)$ for
624 each $j \in \{1, \dots, K\}$. It shows that the KL divergence between these two posteriors is larger when
625 the expected number of samples in that bin is larger.

626 **Lemma G.2.** *Suppose that v is drawn uniformly from the set $\{-1, 1\}^K$, and that $S | v$ is drawn*
627 *from $P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}$. Then for any $j \in \{1, \dots, K/2\}$ and any $\tau \in [0, 1]$,*

$$\text{KL}(\mathbb{Q}(S | v_j = 1) \| \mathbb{Q}(S | v_j = -1)) \leq \frac{n_{\text{maj}}(2 - \tau) + n_{\text{min}}\tau}{3K^3},$$

628 and for any $j \in \{K/2 + 1, \dots, K\}$

$$\text{KL}(\mathbb{Q}(S | v_j = 1) \| \mathbb{Q}(S | v_j = -1)) \leq \frac{n_{\text{maj}}\tau + n_{\text{min}}(2 - \tau)}{3K^3}.$$

629 *Proof.* Let us consider the case when $j = 1$. The bound for all other $j \in \{2, \dots, K\}$ shall follow
630 analogously.

631 Given samples S , let $S = (S_1, \bar{S}_1)$ be a partition where S_1 are the samples that fall in the interval I_1 ,
 632 and \bar{S}_1 be the other samples. Similarly, given a vector $v \in \{-1, 1\}$, let $v = (v_1, \bar{v}_1)$, where v_1 is the
 633 first component and \bar{v}_1 denotes the other components $(2, \dots, K)$ of v .

634 First, we will show that

$$\mathbb{Q}(S | v_1) = \mathbb{Q}(S_1 | v_1)\mathbb{Q}(\bar{S}_1).$$

635 To see this, observe that

$$\mathbb{Q}(S | v_1) = \mathbb{Q}((S_1, \bar{S}_1) | v_1) = \mathbb{Q}(S_1 | v_1)\mathbb{Q}(\bar{S}_1 | v_1, S_1).$$

636 Further, if v is chosen uniformly over the hypercube $\{-1, 1\}^K$, then

$$\begin{aligned} \mathbb{Q}(\bar{S}_1 | v_1, S_1) &= \sum_{\bar{v}_1} \mathbb{Q}(\bar{S}_1, \bar{v}_1 | v_1, S_1) \\ &= \sum_{\bar{v}_1} \mathbb{Q}(\bar{S}_1 | v_1, \bar{v}_1, S_1)\mathbb{Q}(\bar{v}_1 | v_1, S_1) \\ &\stackrel{(i)}{=} \sum_{\bar{v}_1} \mathbb{Q}(\bar{S}_1 | v_1, \bar{v}_1, S_1)\mathbb{Q}(\bar{v}_1) \\ &\stackrel{(ii)}{=} \sum_{\bar{v}_1} \mathbb{Q}(\bar{S}_1 | v_1, \bar{v}_1)\mathbb{Q}(\bar{v}_1) \\ &\stackrel{(iii)}{=} \sum_{\bar{v}_1} \mathbb{Q}(\bar{S}_1 | \bar{v}_1)\mathbb{Q}(\bar{v}_1) \\ &= \mathbb{Q}(\bar{S}_1), \end{aligned}$$

637 where (i) follows since by Bayes' rule

$$\begin{aligned} \mathbb{Q}(\bar{v}_1 | v_1, S_1) &= \frac{\mathbb{Q}(\bar{v}_1 | v_1)\mathbb{Q}(S_1 | v_1, \bar{v}_1)}{\mathbb{Q}(S_1 | v_1)} \\ &= \frac{\mathbb{Q}(\bar{v}_1)\mathbb{Q}(S_1 | v_1, \bar{v}_1)}{\mathbb{Q}(S_1 | v_1)} \quad (\text{since } \bar{v}_1 \text{ is independent of } v_1) \\ &= \frac{\mathbb{Q}(\bar{v}_1)\mathbb{Q}(S_1 | v_1)}{\mathbb{Q}(S_1 | v_1)} = \mathbb{Q}(\bar{v}_1) \quad (\text{the samples in } S_1 \text{ depend only on } v_1). \end{aligned}$$

638 Inequality (ii) follows since the samples are drawn independently given $v = (v_1, \bar{v}_1)$. Finally, (iii)
 639 follows since \bar{S}_1 (the samples that lie outside the interval I_1) only depend on \bar{v}_1 since the marginal
 640 distribution of x is independent of v and the distribution of $y | x$ depends only on the value of v
 641 corresponding to the interval in which x lies.

642 Thus since, $\mathbb{Q}(S | v_1) = \mathbb{Q}(S_1 | v_1)\mathbb{Q}(\bar{S}_1)$ we have that

$$\text{KL}(\mathbb{Q}(S | v_1 = 1) || \mathbb{Q}(S | v_1 = -1)) = \text{KL}(\mathbb{Q}(S_1 | v_1 = 1) || \mathbb{Q}(S_1 | v_1 = -1)). \quad (19)$$

643 To bound this KL divergence, let us condition on the number of samples in S_1 from group a , (the
 644 majority group) $n_{1,a}$ and the number of samples from group b (the minority group), $n_{1,b}$. Now since
 645 $n_{1,a}$ and $n_{1,b}$ are independent of v_1 (which only affects the labels) we have that,

$$\begin{aligned} \mathbb{Q}(S_1 | v_1) &= \sum_{n_{1,a}, n_{1,b}} \mathbb{Q}(n_{1,a}, n_{1,b} | v_1)\mathbb{Q}(S_1 | v_1, n_{1,a}, n_{1,b}) \\ &= \sum_{n_{1,a}, n_{1,b}} \mathbb{Q}(n_{1,a}, n_{1,b})\mathbb{Q}(S_1 | v_1, n_{1,a}, n_{1,b}) \\ &= \mathbb{E}_{n_{1,a}, n_{1,b}} [\mathbb{Q}(S_1 | v_1, n_{1,a}, n_{1,b})]. \end{aligned}$$

646 Therefore, by the joint convexity of the KL-divergence and by Jensen's inequality we have that,

$$\begin{aligned} &\text{KL}(\mathbb{Q}(S_1 | v_1 = 1) || \mathbb{Q}(S_1 | v_1 = -1)) \\ &\leq \mathbb{E}_{n_{1,a}, n_{1,b}} [\text{KL}(\mathbb{Q}(S_1 | v_1 = 1, n_{1,a}, n_{1,b}) || \mathbb{Q}(S_1 | v_1 = -1, n_{1,a}, n_{1,b}))]. \quad (20) \end{aligned}$$

647 Now conditioned on $v_1, n_{1,a}$ and $n_{1,b}$, samples in S_1 are composed of 2 groups of samples $(S_{1,a}, S_{1,b})$.
648 The samples in each group $(S_{1,a}, S_{1,b})$ are drawn independently from the distributions $P_a(x | x \in$
649 $I_1)P_v(y | x)$ and $P_b(x | x \in I_1)P_v(y | x)$ respectively. Therefore,

$$\begin{aligned}
& \text{KL}(\mathbb{Q}(S_1 | v_1 = 1, n_{1,a}, n_{1,b}) \| \mathbb{Q}(S_1 | v_1 = -1, n_{1,a}, n_{1,b})) \\
& \stackrel{(i)}{=} n_{1,a} \text{KL}(P_a(x | x \in I_1)P_{v_1=1}(y | x) \| P_a(x | x \in I_1)P_{v_1=-1}(y | x)) \\
& \quad + n_{1,b} \text{KL}(P_b(x | x \in I_1)P_{v_1=1}(y | x) \| P_b(x | x \in I_1)P_{v_1=-1}(y | x)) \\
& \stackrel{(ii)}{=} (n_{1,a} + n_{1,b}) \mathbb{E}_{x \sim \text{Unif}(I_1)} [\text{KL}(P_{v_1=1}(y | x) \| P_{v_1=-1}(y | x))] \\
& \stackrel{(iii)}{=} \frac{n_{1,a} + n_{1,b}}{2} \mathbb{E}_{x \sim \text{Unif}(I_1)} \left[\sum_{y \in \{-1, 1\}} \left(1 + y\phi \left(x - \frac{1}{2K} \right) \right) \log \left(\frac{(1 + y\phi(x - \frac{1}{2K}))}{(1 + y\phi(x - \frac{1}{2K}))} \right) \right] \\
& = \frac{n_{1,a} + n_{1,b}}{2} \sum_{y \in \{-1, 1\}} \mathbb{E}_{x \sim \text{Unif}(I_1)} \left[\left(1 + y\phi \left(x - \frac{1}{2K} \right) \right) \log \left(\frac{(1 + y\phi(x - \frac{1}{2K}))}{(1 + y\phi(x - \frac{1}{2K}))} \right) \right] \\
& = \frac{n_{1,a} + n_{1,b}}{2K} \sum_{y \in \{-1, 1\}} \int_{x=0}^{\frac{1}{K}} \left[\left(1 + y\phi \left(x - \frac{1}{2K} \right) \right) \log \left(\frac{(1 + y\phi(x - \frac{1}{2K}))}{(1 + y\phi(x - \frac{1}{2K}))} \right) \right] dx \\
& \stackrel{(iv)}{\leq} \frac{n_{1,a} + n_{1,b}}{3K^2}, \tag{21}
\end{aligned}$$

650 where in (i) we let P_{v_1} denote the conditional distribution of y for $x \in I_1$ given v_1 , (ii) follows since
651 both P_a and P_b are constant in the interval, (iii) follows by our construction of P_v above, and finally
652 (iv) follows by invoking Lemma E.3 that ensures that the integral is bounded by $1/3K^2$.

653 Using this bound in Eq. (20), along with Eq. (19) we get that

$$\text{KL}(\mathbb{Q}(S | v_1 = 1) \| \mathbb{Q}(S | v_1 = -1)) \leq \frac{\mathbb{E}[n_{1,a} + n_{2,b}]}{3K^2}.$$

654 Now there are n_{maj} samples from group a in S and n_{min} samples from group b . Therefore,

$$\begin{aligned}
\mathbb{E}[n_{1,a}] &= n_{\text{maj}} P_a(x \in I_1) = \frac{n_{\text{maj}}(2 - \tau)}{K}, \\
\mathbb{E}[n_{1,b}] &= n_{\text{min}} P_b(x \in I_1) = \frac{n_{\text{min}}\tau}{K}.
\end{aligned}$$

655 Plugging this bound into Eq. (21) completes the proof by the first interval. An identical argument
656 holds for $j \in \{2, \dots, K/2\}$. For $j \in \{K/2 + 1, \dots, K\}$ the only change is that

$$\begin{aligned}
\mathbb{E}[n_{j,a}] &= n_{\text{maj}} P_a(x \in I_j) = \frac{n_{\text{maj}}\tau}{K}, \\
\mathbb{E}[n_{j,b}] &= n_{\text{min}} P_b(x \in I_j) = \frac{n_{\text{min}}(2 - \tau)}{K}.
\end{aligned}$$

657 □

658 Next, we combine the previous two lemmas to establish our stated lower bound. We first restate it
659 here.

660 **Theorem B.1.** Consider the group shift setting described in Section B. Given any overlap $\tau \in [0, 1]$
661 recall that $\mathcal{P}_{\text{GS}}(\tau)$ is the class of distributions such that $\text{Overlap}(P_{\text{maj}}, P_{\text{min}}) \geq \tau$. The minimax
662 excess risk in this setting is lower bounded as follows:

$$\begin{aligned}
\text{Minimax Excess Risk}(\mathcal{P}_{\text{GS}}(\tau)) &= \inf_{\mathcal{A}} \sup_{(P_{\text{maj}}, P_{\text{min}}) \in \mathcal{P}_{\text{GS}}(\tau)} \text{Excess Risk}[\mathcal{A}; (P_{\text{maj}}, P_{\text{min}})] \\
&\geq \frac{1}{200(n_{\text{min}} \cdot (2 - \tau) + n_{\text{maj}} \cdot \tau)^{1/3}} \geq \frac{1}{200n_{\text{min}}^{1/3}(\rho \cdot \tau + 2)^{1/3}}, \tag{4}
\end{aligned}$$

663 where $\rho = n_{\text{maj}}/n_{\text{min}} > 1$.

664 *Proof.* First, by Lemma G.1 we know that

$$\text{Minimax Excess Risk}(\mathcal{P}_{\text{GS}}(\tau)) \geq \frac{1}{32K^2} \sum_{j=1}^K \exp(-\text{KL}(\mathbb{Q}(S | v_j = 1) \| \mathbb{Q}(S | v_j = -1))).$$

665 Next, by invoking the bound on the KL divergences in the equation above by Lemma G.2 we get that

$$\begin{aligned} & \text{Minimax Excess Risk}(\mathcal{P}_{\text{GS}}(\tau)) \\ & \geq \frac{1}{64K} \left[\exp\left(-\frac{n_{\text{maj}}(2-\tau) + n_{\text{min}}\tau}{3K^3}\right) + \exp\left(-\frac{n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau}{3K^3}\right) \right] \\ & \geq \frac{1}{64K} \left[\exp\left(-\frac{n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau}{3K^3}\right) \right] \end{aligned}$$

666 Setting $K = \lceil (n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau)^{1/3} \rceil$ and recalling that $\tau \leq 1$ we get that

$$\begin{aligned} & \text{Minimax Excess Risk}(\mathcal{P}_{\text{GS}}(\tau)) \\ & \geq \frac{1}{64 \lceil (n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau)^{1/3} \rceil} \left[\exp\left(-\frac{n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau}{3 \lceil (n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau)^{1/3} \rceil^3}\right) \right] \\ & \stackrel{(i)}{\geq} \frac{\exp(-1/3)}{64} \frac{(n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau)^{1/3}}{\lceil (n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau)^{1/3} \rceil} \frac{1}{(n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau)^{1/3}} \\ & \stackrel{(ii)}{\geq} \frac{0.7 \exp(-1/3)}{64} \frac{1}{(n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau)^{1/3}} \\ & \geq \frac{1}{200} \frac{1}{(n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau)^{1/3}}, \end{aligned}$$

667 where (i) follows since $n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau / \lceil (n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau)^{1/3} \rceil^3 \leq 1$, and (ii) follows
668 since $0 \leq \tau \leq 1$ and $n_{\text{min}} \geq 1$ and hence $\frac{(n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau)^{1/3}}{\lceil (n_{\text{min}}(2-\tau) + n_{\text{maj}}\tau)^{1/3} \rceil} \geq 0.7$. \square

669 G.2 Proof of Theorem D.2

670 In this section, we derive an upper bound on the excess risk of the undersampled binning estimator
671 \mathcal{A}_{USB} (Eq. (5)). Recall that given a dataset \mathcal{S} this estimator first calculates the undersampled dataset
672 \mathcal{S}_{US} , where the number of points from the minority group (n_{min}) is equal to the number of points from
673 the majority group (n_{min}), and the size of the dataset is $2n_{\text{min}}$. Throughout this section, $(\mathbb{P}_{\text{maj}}, \mathbb{P}_{\text{min}})$
674 shall be an arbitrary element of $\mathcal{P}_{\text{GS}}(\tau)$ for any $\tau \in [0, 1]$. In this section, whenever we shall often
675 denote $\text{Excess Risk}(\mathcal{A}; (\mathbb{P}_{\text{maj}}, \mathbb{P}_{\text{min}}))$ by simply $\text{Excess Risk}(\mathcal{A})$.

676 Before we proceed, we introduce some additional notation. For any $j \in \{1, \dots, K\}$ and $I_j =$
677 $[\frac{j-1}{K}, \frac{j}{K}]$ let

$$q_{j,1} := \mathbb{P}_{\text{test}}(y = 1 | x \in I_j) = \int_{x \in I_j} \mathbb{P}(y = 1 | x) \mathbb{P}_{\text{test}}(x | x \in I_j) dx, \quad (22a)$$

$$q_{j,-1} := \mathbb{P}_{\text{test}}(y = -1 | x \in I_j) = \int_{x \in I_j} \mathbb{P}(y = -1 | x) \mathbb{P}_{\text{test}}(x | x \in I_j) dx. \quad (22b)$$

678 For the undersampled binning estimator \mathcal{A}_{USB} (defined above in Eq. (5)), define the *excess risk in an*
679 *interval* I_j as follows:

$$\begin{aligned} R_j(\mathcal{A}_{\text{USB}}^{\mathcal{S}}) & := p(y = -\mathcal{A}_j^{\mathcal{S}} | x \in I_j) - \min\{\mathbb{P}_{\text{test}}(y = 1 | x \in I_j), \mathbb{P}_{\text{test}}(y = -1 | x \in I_j)\} \\ & = q_{j,-\mathcal{A}_j^{\mathcal{S}}} - \min\{q_{j,1}, q_{j,-1}\}. \end{aligned}$$

680 The proof of the upper bound shall proceed in steps. First, in Lemma G.3 we will show that the
681 excess risk is equal to sum the excess risk over the intervals up to a factor of $2/K$ on account of the
682 distribution being 1-Lipschitz. Next, in Lemma G.4 we upper bound the risk over each interval. We
683 put these two together and to upper bound the risk.

684 **Lemma G.3.** *The expected excess risk of undersampled binning estimator \mathcal{A}_{USB} can be decomposed*
 685 *as follows*

$$\text{Excess Risk}(\mathcal{A}_{\text{USB}}) \leq \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{S} \sim \mathcal{P}_{\text{maj}}^{n_{\text{maj}}} \times \mathcal{P}_{\text{min}}^{n_{\text{min}}}} [R_j(\mathcal{A}_{\text{USB}}^{\mathcal{S}})] \cdot \mathbb{P}_{\text{test}}(I_j) + \frac{2}{K},$$

686 where $\mathbb{P}_{\text{test}}(I_j) := \int_{x \in I_j} \mathbb{P}_{\text{test}}(x) \, dx$.

687 *Proof.* Recall that by definition, the expected excess risk is

$$\mathbb{E}_{\mathcal{S} \sim \mathcal{P}_{\text{maj}}^{n_{\text{maj}}} \times \mathcal{P}_{\text{min}}^{n_{\text{min}}}} [R(\mathcal{A}^{\mathcal{S}}; \mathbb{P}_{\text{test}}) - R(f^*; \mathbb{P}_{\text{test}})].$$

688 Let us first decompose the Bayes risk $R(f^*)$,

$$\begin{aligned} R(f^*) &= \inf_f \mathbb{E}_{(x,y) \sim \mathbb{P}_{\text{test}}} [\mathbf{1}(f(x) \neq y)] \\ &= \inf_f \int_{x=0}^1 \sum_{y \in \{-1,1\}} \mathbf{1}(f(x) \neq y) \mathbb{P}_{\text{test}}(y | x) \mathbb{P}_{\text{test}}(x) \, dx \\ &= \int_{x=0}^1 \inf_{f(x) \in \{-1,1\}} \sum_{y \in \{-1,1\}} \mathbf{1}(f(x) \neq y) \mathbb{P}_{\text{test}}(y | x) \mathbb{P}_{\text{test}}(x) \, dx \\ &= \int_{x=0}^1 \inf_{f(x) \in \{-1,1\}} \mathbb{P}_{\text{test}}(y = -f(x) | x) \mathbb{P}_{\text{test}}(x) \, dx \\ &= \int_{x=0}^1 \min \{ \mathbb{P}_{\text{test}}(y = 1 | x), \mathbb{P}_{\text{test}}(y = -1 | x) \} \mathbb{P}_{\text{test}}(x) \, dx. \end{aligned} \quad (23)$$

689 The risk of the undersampled binning algorithm \mathcal{A}_{USB} is given by

$$\begin{aligned} R(\mathcal{A}_{\text{USB}}^{\mathcal{S}}) &= \int_{x=0}^1 \sum_{y \in \{-1,1\}} \mathbf{1}(\mathcal{A}_{\text{USB}}^{\mathcal{S}}(x) \neq y) \mathbb{P}_{\text{test}}(y | x) \mathbb{P}_{\text{test}}(x) \, dx \\ &= \int_{x=0}^1 \mathbb{P}_{\text{test}}(y = -\mathcal{A}_{\text{USB}}^{\mathcal{S}}(x) | x) \mathbb{P}_{\text{test}}(x) \, dx. \end{aligned}$$

690 Next, recall that the undersampled binning estimator is constant over the intervals I_j for $j \in$
 691 $\{1, \dots, K\}$ where it takes the value $\mathcal{A}_j^{\mathcal{S}}$ (to ease notation let us simply denote it by \mathcal{A}_j below), and
 692 therefore

$$R(\mathcal{A}_{\text{USB}}^{\mathcal{S}}) = \sum_{j=0}^{K-1} \int_{x \in I_j} \mathbb{P}_{\text{test}}(y = -\mathcal{A}_j | x) \mathbb{P}_{\text{test}}(x) \, dx.$$

693 This combined with Eq. (23) tells us that

$$\begin{aligned} &R(\mathcal{A}_{\text{USB}}^{\mathcal{S}}) - R(f^*) \\ &= \sum_{j=0}^{K-1} \int_{x \in I_j} (\mathbb{P}_{\text{test}}(y = -\mathcal{A}_j | x) - \min \{ \mathbb{P}_{\text{test}}(y = 1 | x), \mathbb{P}_{\text{test}}(y = -1 | x) \}) \mathbb{P}_{\text{test}}(x) \, dx. \end{aligned} \quad (24)$$

694 Recall the definition of $q_{j,1}$ and $q_{j,-1}$ from Eqs. (22a)-(22b) above. For any $x \in I_j = [\frac{j-1}{K}, \frac{j}{K}]$,
 695 $|\mathbb{P}_{\text{test}}(y | x) - q_{j,y}| \leq 1/K$, since the distribution $\mathbb{P}_{\text{test}}(y | x)$ is 1-Lipschitz and $q_{j,y}$ is its conditional
 696 mean. Therefore,

$$\begin{aligned} &R(\mathcal{A}_{\text{USB}}^{\mathcal{S}}) - R(f^*) \\ &\leq \sum_{j=0}^{K-1} \int_{x \in I_j} (q_{j,-\mathcal{A}_j} - \min \{ q_{j,1}, q_{j,-1} \}) \mathbb{P}_{\text{test}}(x) \, dx + \frac{2}{K} \sum_{j=0}^{K-1} \int_{x \in I_j} \mathbb{P}_{\text{test}}(x) \, dx \\ &= \sum_{j=0}^{K-1} \int_{x \in I_j} R_j(\mathcal{A}_{\text{USB}}^{\mathcal{S}}) \mathbb{P}_{\text{test}}(x) \, dx + \frac{2}{K}. \end{aligned}$$

697 Taking expectation over the training samples \mathcal{S} (where n_{min} samples are drawn independently from
 698 \mathbb{P}_{min} and n_{maj} samples are drawn independently from \mathbb{P}_{maj}) concludes the proof. \square

699 Next we provide an upper bound on the expected excess risk is an interval $R_j(\mathcal{A}_{\text{USB}}^S)$.

700 **Lemma G.4.** For any $j \in \{1, \dots, K\}$ with $I_j = [\frac{j-1}{K}, \frac{j}{K}]$,

$$\mathbb{E}_{\mathcal{S} \sim \mathbb{P}_{\text{maj}}^{n_{\text{maj}}} \times \mathbb{P}_{\text{min}}^{n_{\text{min}}}} [R_j(\mathcal{A}_{\text{USB}}^S)] \leq \frac{c}{\sqrt{n_{\text{min}} \mathbb{P}_{\text{test}}(I_j)}} + \frac{c}{K},$$

701 where c is an absolute constant, and $\mathbb{P}_{\text{test}}(I_j) := \int_{x \in I_j} \mathbb{P}_{\text{test}}(x) dx$.

702 *Proof.* Consider an arbitrary bucket $j \in \{1, \dots, K\}$.

703 Let us introduce some notation that shall be useful in the remainder of the proof. Analogous to $q_{j,1}$
704 and $q_{j,-1}$ defined above (see Eqs. (22a)-(22b)), define $q_{j,1}^a$ and $q_{j,1}^b$ as follows:

$$q_{j,1}^a := \mathbb{P}_a(y = 1 \mid x \in I_j) = \int_{x \in I_j} \mathbb{P}(y = 1 \mid x) \mathbb{P}_a(x \mid x \in I_j) dx, \quad (25a)$$

$$q_{j,1}^b := \mathbb{P}_b(y = 1 \mid x \in I_j) = \int_{x \in I_j} \mathbb{P}(y = 1 \mid x) \mathbb{P}_b(x \mid x \in I_j) dx. \quad (25b)$$

705 Essentially, $q_{j,1}^a$ is the probability that a sample is from group a and has label 1, conditioned on the
706 event that the sample falls in the interval I_j . Since

$$\mathbb{P}_{\text{test}}(x \mid x \in I_j) = \frac{1}{2} [\mathbb{P}_a(x \mid x \in I_j) + \mathbb{P}_b(x \mid x \in I_j)],$$

707 therefore

$$\begin{aligned} |q_{j,1} - q_{j,1}^a| &= \left| \int_{x \in I_j} \mathbb{P}(y = 1 \mid x) \mathbb{P}_{\text{test}}(x \mid x \in I_j) dx - \int_{x \in I_j} \mathbb{P}(y = 1 \mid x) \mathbb{P}_a(x \mid x \in I_j) dx \right| \\ &\leq \frac{1}{K}. \end{aligned} \quad (26)$$

708 This follows since $\mathbb{P}(y \mid x)$ is 1-Lipschitz and therefore can fluctuate by at most $1/K$ in the interval
709 I_j . Of course the same bound also holds for $|q_{j,1} - q_{j,1}^b|$.

710 With this notation in place let us present a bound on the expected value of $R_j(\mathcal{A}_{\text{USB}}^S)$. By definition

$$R_j(\mathcal{A}_{\text{USB}}^S) = q_{j,-\mathcal{A}_j^S} - \min\{q_{j,1}, q_{j,-1}\}.$$

711 First, note that $q_{j,1} := \mathbb{P}_{\text{test}}(y = 1 \mid x \in I_j) = 1 - q_{j,-1}$. Suppose that $q_{j,1} < 1/2$ and therefore
712 $q_{j,-1} > 1/2$ (the same bound shall hold in the other case). In this case, risk is incurred only when
713 $\mathcal{A}_j^S = 1$. That is,

$$\begin{aligned} \mathbb{E}_{\mathcal{S} \sim \mathbb{P}_{\text{maj}}^{n_{\text{maj}}} \times \mathbb{P}_{\text{min}}^{n_{\text{min}}}} [R_j(\mathcal{A}_{\text{USB}}^S)] &= |q_{j,-1} - q_{j,1}| \mathbb{P}_{\mathcal{S}}[\mathcal{A}_j^S = 1] \\ &= |1 - 2q_{j,1}| \mathbb{P}_{\mathcal{S}}[\mathcal{A}_j^S = 1]. \end{aligned} \quad (27)$$

714 Now by the definition of the undersampled binning estimator (see Eq. (5)), $\mathcal{A}_j^S = 1$ only when there
715 are more samples in the interval I_j with label 1 than -1 . However, we can bound the probability of
716 this happening since $q_{j,1}$ is smaller than $q_{j,-1}$.

717 Let n_j be the number of samples in the undersampled sample set \mathcal{S}_{US} in the interval I_j . Let $n_{1,j}$ be
718 the number of these samples with label 1, and $n_{-1,j} = n_j - n_{1,j}$ be the number of samples with
719 label -1 . Further, let $n_{a,j}$ be the number of samples in from group a such that they fall in the interval
720 I_j , and define $m_{b,j}$ analogously.

721 The probability of incurring risk is given by

$$\mathbb{P}[\mathcal{A}_j = 1] = \sum_{s=1}^{2n_{\text{min}}} \mathbb{P}[\mathcal{A}_j = 1 \mid n_j = s] \mathbb{P}[n_j = s], \quad (28)$$

722 where the sum is up to $2n_{\text{min}}$ since the size of the undersample dataset $|\mathcal{S}_{\text{US}}|$ is equal to $2n_{\text{min}}$.

723 Conditioned on the event that $n_j = s$ the probability of incurring risk is

$$\begin{aligned} \mathbb{P}[\mathcal{A}_j = 1 \mid n_j = s] &= \mathbb{P}[m_{1,j} > n_{-1,j} \mid n_j = s] = \mathbb{P}[n_{1,j} > n_j/2 \mid n_j = s] \\ &= \mathbb{P}[n_{1,j} > s/2 \mid n_j = s]. \end{aligned} \quad (29)$$

724 Now, note that $n_j = n_{a,j} + n_{b,j}$. Thus continuing, we have that

$$\begin{aligned} \mathbb{P}[n_{1,j} > s/2 \mid n_j = s] &= \sum_{s' \leq s} \mathbb{P}[n_{1,j} > s/2 \mid n_j = s, n_{b,j} = s'] \mathbb{P}[n_{b,j} = s'] \\ &= \sum_{s' \leq s} \mathbb{P}[n_{1,j} > s/2 \mid n_{a,j} = s - s', n_{b,j} = s'] \mathbb{P}[n_{b,j} = s']. \end{aligned}$$

725 In light of this previous equation, we want to control the probability that the number of samples with
726 label 1 in the interval I_j conditioned on the event that the number of samples from group a in this
727 interval is $s - s'$ and the number of samples from group b in this interval is s' . Recall that $q_{j,1}^a$ and
728 $q_{j,1}^b$ the probabilities of the label of the sample being 1 conditioned the event that sample is in the
729 interval I_j when it is group a and b respectively. So we define the random variables:

$$z_a[s - s'] \sim \text{Bin}(s - s', q_{j,1}^a), \quad z_b[s'] \sim \text{Bin}(s', q_{j,1}^b), \quad z[s] \sim \text{Bin}(s, \max\{q_{j,1}^a, q_{j,1}^b\}).$$

730 Then,

$$\begin{aligned} &\mathbb{P}[n_{1,j} > s/2 \mid n_j = s] \\ &= \sum_{s' \leq s} \mathbb{P}[n_{1,j} > s/2 \mid n_{j,a} = s - s', n_{j,b} = s'] \mathbb{P}[n_{j,b} = s'] \\ &= \sum_{s' \leq s} \mathbb{P}[z_a[s - s'] + z_b[s'] > s/2 \mid n_{a,j} = s - s', n_{b,j} = s'] \mathbb{P}[n_{b,j} = s'] \\ &\leq \sum_{s' \leq s} \mathbb{P}[z[s] > s/2 \mid n_{a,j} = s - s', n_{b,j} = s'] \mathbb{P}[n_{b,j} = s'] \\ &= \sum_{s' \leq s} \mathbb{P}[z[s] > s/2] \mathbb{P}[n_{b,j} = s'] \\ &= \mathbb{P}[z[s] > s/2] \\ &\stackrel{(i)}{\leq} \exp\left(-\frac{s}{2}(1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\})^2\right), \end{aligned} \quad (30)$$

731 where (i) follows by invoking Hoeffding's inequality [Wainwright, 2019, Proposition 2.5]. Combining
732 this with Eqs. (28) and (29) we get that

$$\mathbb{P}[\mathcal{A}_j = 1] \leq \sum_{s=1}^{2n_{\min}} \exp\left(-\frac{s}{2}(1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\})^2\right) \mathbb{P}[n_j = s].$$

733 Now n_j , which is the number of samples that lands in the interval I_j is equal to $n_{a,j} + n_{b,j}$. Now each
734 of $n_{a,j}$ and $n_{b,j}$ (the number of samples in this interval from each of the groups) are random variables
735 with distributions $\text{Bin}(n_{\min}, P_a(I_j))$ and $\text{Bin}(n_{\min}, P_b(I_j))$, where $P_a(I_j) = \int_{x \in I_j} P_a(x) dx$ and
736 $P_b(I_j) = \int_{x \in I_j} P_b(x) dx$. Therefore, n_j is distributed as a sum of two binomial distribution and is
737 therefore Poisson binomially distributed [Wikipedia contributors, 2022]. Using the formula for the
738 moment generating function (MGF) of a Poisson binomially distributed random variable we infer
739 that,

$$\begin{aligned} \mathbb{P}[\mathcal{A}_j = 1] &\leq \left(1 - P_a(I_j) + P_a(I_j) \exp\left(-\frac{(1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\})^2}{2}\right)\right)^{n_{\min}} \times \\ &\quad \left(1 - P_b(I_j) + P_b(I_j) \exp\left(-\frac{(1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\})^2}{2}\right)\right)^{n_{\min}}. \end{aligned}$$

740 Plugging this into Eq. (28) we get that,

$$\begin{aligned}
& \mathbb{E}_{\mathcal{S} \sim \mathcal{P}_{\text{maj}}^{n_{\text{maj}}} \times \mathcal{P}_{\text{min}}^{n_{\text{min}}}} [R_j(\mathcal{A}_{\text{USB}}^{\text{S}})] \\
& \leq |1 - 2q_{j,1}| \left[1 - \text{P}_a(I_j) + \text{P}_a(I_j) \exp\left(-\frac{(1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\})^2}{2}\right) \right]^{n_{\text{min}}} \times \\
& \quad \left[1 - \text{P}_b(I_j) + \text{P}_b(I_j) \exp\left(-\frac{(1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\})^2}{2}\right) \right]^{n_{\text{min}}} \\
& = |1 - 2q_{j,1}| \left[1 - \text{P}_a(I_j) \left(1 - \exp\left(-\frac{(1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\})^2}{2}\right) \right) \right]^{n_{\text{min}}} \times \\
& \quad \left[1 - \text{P}_b(I_j) \left(1 - \exp\left(-\frac{(1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\})^2}{2}\right) \right) \right]^{n_{\text{min}}}.
\end{aligned}$$

741 Since $|1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\}| \leq 1$,

$$1 - \exp\left(-\frac{(1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\})^2}{2}\right) \geq \frac{(1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\})^2}{4},$$

742 and therefore

$$\begin{aligned}
\mathbb{E}_{\mathcal{S} \sim \mathcal{P}_{\text{maj}}^{n_{\text{maj}}} \times \mathcal{P}_{\text{min}}^{n_{\text{min}}}} [R_j(\mathcal{A}_{\text{USB}}^{\text{S}})] & \leq |1 - 2q_{j,1}| \left[1 - \text{P}_a(I_j) \frac{(1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\})^2}{2} \right]^{n_{\text{min}}} \times \\
& \quad \left[1 - \text{P}_b(I_j) \frac{(1 - 2 \max\{q_{j,1}^a, q_{j,1}^b\})^2}{2} \right]^{n_{\text{min}}} \\
& \stackrel{(i)}{\leq} |1 - 2q_{j,1}| \left[1 - \text{P}_a(I_j) \frac{(1 - 2q_{j,1} - 2\gamma)^2}{2} \right]^{n_{\text{min}}} \times \\
& \quad \left[1 - \text{P}_b(I_j) \frac{(1 - 2q_{j,1} - 2\gamma)^2}{2} \right]^{n_{\text{min}}} \\
& \stackrel{(ii)}{\leq} |1 - 2q_{j,1}| \exp\left(-n_{\text{min}}(\text{P}_a(I_j) + \text{P}_b(I_j)) \frac{(1 - 2q_{j,1} - 2\gamma)^2}{2}\right),
\end{aligned}$$

743 where (i) follows since $|\max\{q_{j,1}^a, q_{j,1}^b\} - q_{j,1}| \leq 1/K$ by Eq. (26) and γ is such that $|\gamma| \leq 1/K$, and

744 (ii) follows since $(1 + z)^b \leq \exp(bz)$. Now the RHS above is maximized when $(1 - 2q_{j,1} - 2\gamma)^2 =$

745 $\frac{c}{n_{\text{min}}(\text{P}_a(I_j) + \text{P}_b(I_j))}$, for some constant c . Plugging this into the equation above we get that

$$\begin{aligned}
\mathbb{E}_{\mathcal{S} \sim \mathcal{P}_{\text{maj}}^{n_{\text{maj}}} \times \mathcal{P}_{\text{min}}^{n_{\text{min}}}} [R_j(\mathcal{A}_{\text{USB}}^{\text{S}})] & \leq \frac{c'}{\sqrt{n_{\text{min}}(\text{P}_a(I_j) + \text{P}_b(I_j))}} + c'|\gamma| \\
& \leq \frac{c'}{\sqrt{n_{\text{min}}(\text{P}_a(I_j) + \text{P}_b(I_j))}} + \frac{c'}{K}.
\end{aligned}$$

746 Finally, noting that $\text{P}_{\text{test}}(I_j) = (\text{P}_a(I_j) + \text{P}_b(I_j))/2$ completes the proof. \square

747 By combining the previous two lemmas we can now prove our upper bound on the risk of the
748 undersampled binning estimator. We begin by restating it.

749 **Theorem D.2.** Consider the group shift setting described in Section B. For any overlap $\tau \in [0, 1]$
750 and for any $(\text{P}_{\text{maj}}, \text{P}_{\text{min}}) \in \mathcal{P}_{\text{GS}}(\tau)$ the expected excess risk of the Undersampling Binning Estimator
751 (Eq. (5)) with number of bins with $K = \lceil n_{\text{min}}^{1/3} \rceil$ is

$$\text{Excess Risk}[\mathcal{A}_{\text{USB}}; (\text{P}_{\text{maj}}, \text{P}_{\text{min}})] = \mathbb{E}_{\mathcal{S} \sim \mathcal{P}_{\text{maj}}^{n_{\text{maj}}} \times \mathcal{P}_{\text{min}}^{n_{\text{min}}}} [R(\mathcal{A}_{\text{USB}}^{\text{S}}; \text{P}_{\text{test}}) - R(f^*; \text{P}_{\text{test}})] \leq \frac{C}{n_{\text{min}}^{1/3}}.$$

752 *Proof.* First by Lemma G.3 we know that

$$\text{Excess Risk}[\mathcal{A}_{\text{USB}}] \leq \sum_{j=0}^{K-1} \mathbb{E}_{\mathcal{S} \sim \mathcal{P}_{\text{maj}}^{n_{\text{maj}}} \times \mathcal{P}_{\text{min}}^{n_{\text{min}}}} [R_j(\mathcal{A}_{\text{USB}}^{\mathcal{S}})] \cdot \mathbb{P}_{\text{test}}(I_j) + \frac{2}{K}.$$

753 Next by using the bound on $\mathbb{E}_{\mathcal{S} \sim \mathcal{P}_{\text{maj}}^{n_{\text{maj}}} \times \mathcal{P}_{\text{min}}^{n_{\text{min}}}} [R_j(\mathcal{A}_{\text{USB}}^{\mathcal{S}})]$ established in Lemma G.4 we get that,

$$\begin{aligned} \text{Excess Risk}(\mathcal{A}_{\text{USB}}) &\leq c \sum_{j=0}^{K-1} \frac{1}{\sqrt{n_{\text{min}} \mathbb{P}_{\text{test}}(I_j)}} \mathbb{P}_{\text{test}}(I_j) + \frac{c}{K} \\ &= \frac{c}{\sqrt{n_{\text{min}}}} \sum_{j=0}^{K-1} \sqrt{\mathbb{P}_{\text{test}}(I_j)} + \frac{c}{K} \\ &\stackrel{(i)}{\leq} \frac{c}{\sqrt{n_{\text{min}}}} \sqrt{K} \sum_{j=0}^{K-1} \mathbb{P}_{\text{test}}(I_j) + \frac{c}{K} \\ &= c \sqrt{\frac{K}{n_{\text{min}}}} + \frac{c}{K}. \end{aligned}$$

754 where (i) follows since for any vector $z \in \mathbb{R}^K$, $\|z\|_1 \leq \sqrt{K} \|z\|_2$. Maximizing over K yields the
755 choice $K = \lceil n_{\text{min}}^{1/3} \rceil$, completing the proof.

756 □

757 H Additional Simulations

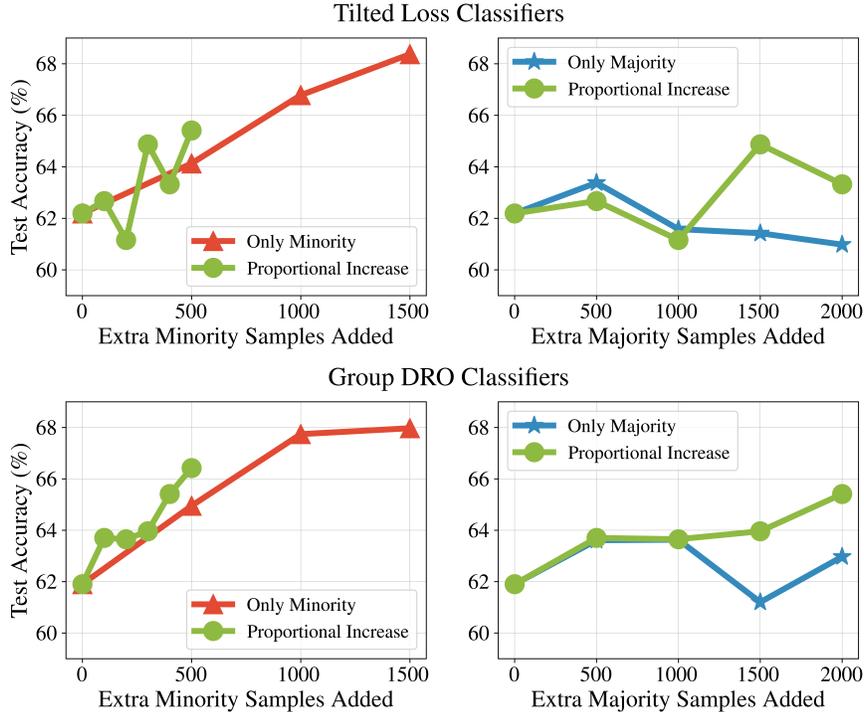


Figure 3: Convolutional neural network classifiers trained on the Imbalanced Binary CIFAR10 dataset with a 5:1 label imbalance. (Top) Models trained using the tilted loss [Li et al., 2020] with early stopping. (Bottom) Models trained using group-DRO [Sagawa et al., 2020] with early stopping. We report the average test accuracy calculated on a balanced test set over 5 random seeds. We start off with 2500 cat examples and 500 dog examples in the training dataset. We find similar trends to those obtained in Figure 1 even with these losses that are designed to optimize for the worst group accuracy.

758 **I Experimental Details for Figures 1 and 3**

759 We construct our label shift dataset from the original CIFAR10 dataset. We create a binary classi-
 760 fication task using the “cat” and “dog” classes. We use the official test examples as the balanced
 761 test set with 1000 cats and 1000 dogs. To form the initial train and validation sets, we use 2500 cat
 762 examples (half of the training set) and 500 dog examples, corresponding to a 5:1 label imbalance. We
 763 use 80% of those examples for training and the rest for validation. We are left with 2500 additional
 764 cat examples and 4500 dog examples from the original train set which we add into our training set to
 765 generate Figure 1.

766 We use the same convolutional neural network architecture as [Byrd and Lipton, 2019, Wang et al.,
 767 2022] with random initializations for this dataset. We train this model using SGD for 800 epochs
 768 with batchsize 64, a constant learning rate 0.001 and momentum 0.9. The importance weights
 769 used upweight the minority class samples in the training loss and validation loss is calculated to be
 770 $\frac{\# \text{Cat Train Examples}}{\# \text{Dog Train Examples}}$. We note that all of the experiments were performed on an internal cluster on 8
 771 GPUs.

772 **VS Loss:** Given a dataset $\{x_i, y_i\}_{i=1}^n$, the VS loss [Kini et al., 2021] is defined as follows

$$\mathcal{L}_{\text{VS}}(f) := \sum_{i=1}^n \log \left(1 + \exp \left(- \left(\frac{n_{g_i}}{n_{\max}} \right)^\gamma y_i f(x_i) - \frac{\tau n_{g_i}}{n} \right) \right),$$

773 where g_i denotes the group label, n_{g_i} corresponds to the number of samples from the group, n_{\max}
 774 is the number of samples in the largest group and n is the total number of samples. We set $\tau = 3$
 775 and $\gamma = 0.3$, the best hyperparameters identified by Wang et al. [2022] on this dataset for this neural
 776 network architecture.

777 **Tilted Loss:** The tilted loss [Li et al., 2020] is defined as

$$\mathcal{L}_{\text{Tilted}}(f) := \frac{1}{t} \log \left[\sum_{i=1}^n \exp(t\ell(y_i f(x_i))) \right],$$

778 where we take ℓ to be the logistic loss. In our experiments we set $t = 2$.

779 **Group-DRO:** We run group-DRO [Sagawa et al., 2020, Algorithm 1] with the logistic loss. We set
 780 adversarial step-size $\eta_q = 0.05$ which was the best hyperparameter identified by Wang et al. [2022].