

## 393 A Technical Tools

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

### 395 A.1 Reduction to lower bounds over a finite class

396 The lower bound on the minimax excess risk will be established via the usual route of first identifying  
 397 a “hard” finite set of problem instances and then establishing the lower bound over this finite class.  
 398 One difference from the usual setup in proving such lower bounds [see 22, Chapter 15] is that the  
 399 training samples are drawn from an imbalanced distribution, whereas the test samples are drawn from  
 400 a balanced one.

401 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-  
 402 tributions over  $[0, 1] \times \{-1, 1\}$ . As before, we let  $P_{\text{test}}$  denote the uniform mixture over  $P_{\text{maj}}$   
 403 and  $P_{\text{min}}$ . We let  $\mathcal{V}$  denote a finite index set. Corresponding to each element  $v \in \mathcal{V}$  there is a  
 404  $P_v = (P_{v,\text{maj}}, P_{v,\text{min}}) \in \mathcal{P}$  with  $P_{v,\text{test}} = (P_{v,\text{maj}} + P_{v,\text{min}})/2$ . Finally, also define a pair of random  
 405 variables  $(V, S)$  as follows:

- 406 1.  $V$  is a uniform random variable over the set  $\mathcal{V}$ .
- 407 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  $n_{\text{maj}}$  samples from  $P_{v,\text{maj}}$  and  
 408  $n_{\text{min}}$  samples from  $P_{v,\text{min}}$ .

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

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

413 **Lemma A.1.** *Let the random variables  $(V, S)$  be as defined above. The minimax excess risk is lower*  
 414 *bounded as follows:*

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

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

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

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

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

417 We continue lower bounding the first term as follows

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

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

420 **A.2 The Hat Function and its Properties**

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

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

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

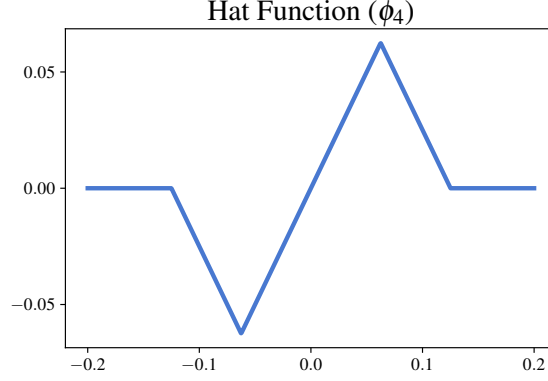


Figure 3: The hat function with  $K = 4$ .

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

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

426 We also compute some other key quantities for  $\phi$ .

427 **Lemma A.2.** For any positive integer  $K$ ,

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

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

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

429 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  $\frac{1}{4K}$ ,  
 430 thus it has area  $\frac{1}{16K^2}$ . Therefore,

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

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

433 **Lemma A.3.** For any positive integer  $K$ ,

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

434 and

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

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

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

437 where the last equality follows since  $\phi$  is an odd function. Now, we may collect the integrands to get  
438 that,

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

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

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

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

## 442 B Proofs in the Label Shift Setting

443 Throughout this section we operate in the label shift setting (Section 3.2.1).

444 First, in Appendix B.1 through a sequence of lemmas we prove the minimax lower bound Theorem 4.1.  
445 Next, in Appendix B.2 we prove Theorem 5.1 which is an upper bound on the excess risk of the  
446 undersampled binning estimator (see Eq. (5)) with  $\lceil n_{\min} \rceil^{1/3}$  bins by invoking previous results on  
447 nonparametric density estimation [9, 8].

448 **B.1 Proof of Theorem 4.1**

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

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

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

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

456 where  $\phi$  is defined in Eq. 6. Notice that  $P_{v,1}$  only depends on  $v_1$  while  $P_{v,-1}$  only depends on  $v_{-1}$ .  
 457 We continue to define We continue to define

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

458 and

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

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

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

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

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

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

466 *Proof.* By invoking Lemma A.1 we get that

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

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

469 **Calculation of  $\mathfrak{R}_{\mathcal{V}}$ :** Immediately by Le Cam’s lemma [22, Eq. 15.13], we get that

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

470 **Calculation of  $\mathfrak{B}_{\mathcal{V}}$ :** Again by invoking Le Cam's lemma [22, Eq. 15.13], we get that for any class  
 471 conditional distributions  $P_1, P_{-1}$ ,

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

472 So by taking expectations, we get that

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

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

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

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

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

476 This yields that

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

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

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

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

479 □

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

482 **Lemma B.2.** *Suppose that  $v$  is drawn uniformly from the set  $\{-1, 1\}^K$ , and that  $S \mid v$  is drawn from*  
 483  $P_{v,\text{maj}}^{n_{\text{maj}}} \times P_{v,\text{min}}^{n_{\text{min}}}$  *then,*

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

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

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

485 where the last equality is by the definition of  $P_{v,1}$  and  $P_{v,-1}$ . Continuing we get that,

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

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

488 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  
489 hence all of the summands in the RHS above are equal. Thus,

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

490 Before we continue further, let us define

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

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

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

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

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

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

495 Note that  $\mathbb{Q}(S | v) = \mathbb{P}_{v,\text{maj}}^{n_{\text{maj}}} \times \mathbb{P}_{v,\text{min}}^{n_{\text{min}}}$ , and therefore

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

496 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  
 497  $\mathcal{V}^{(1,-1)}$ .

498 Now by the Bretagnolle–Huber inequality [see 4, Corollary 4],

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

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

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

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

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

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

504 where the last inequality is a result of the calculation in Lemma A.3.

505 Therefore, we get

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

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

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

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

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

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

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

509 completing the proof.  $\square$

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

512 **Theorem 4.1.** *Consider the label shift setting described in Section 3.2.1. Recall that  $\mathcal{P}_{\text{LS}}$  is the class  
513 of pairs of distributions  $(\mathbf{P}_{\text{maj}}, \mathbf{P}_{\text{min}})$  that satisfy the assumptions in that section. The minimax excess  
514 risk over this class is lower bounded as follows:*

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

515 *Proof.* By Lemma B.1 we know that,

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

516 Next by the calculation in Lemma B.2 we have that

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

517 Setting  $K = \lceil n_{\text{min}}^{1/3} \rceil$  yields the result.  $\square$

## 518 B.2 Proof of Theorem 5.1

519 In this section, we derive an upper bound on the excess risk of the undersampled binning estimator  
520  $\mathcal{A}_{\text{USB}}$  (Eq. (5)) in the label shift setting. Recall that given a dataset  $\mathcal{S}$  this estimator first calculates  
521 the undersampled dataset  $\mathcal{S}_{\text{US}}$ , where the number of points from the minority group ( $n_{\text{min}}$ ) is equal to



522 the number of points from the majority group ( $n_{\min}$ ), and the size of the dataset is  $2n_{\min}$ . Throughout  
 523 this section,  $(P_{\text{maj}}, P_{\text{min}})$  shall be an arbitrary element of  $\mathcal{P}_{\text{LS}}$ .

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

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

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

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

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

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

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

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

534 Observe that, for  $x \in I_j$

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

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

537 Before proving the theorem we restate it here.

538 **Theorem 5.1.** *Consider the label shift setting described in Section 3.2.1. For any  $(P_{\text{maj}}, P_{\text{min}}) \in \mathcal{P}_{\text{LS}}$   
 539 the expected excess risk of the Undersampling Binning Estimator (Eq. (5)) with number of bins with  
 540  $K = c \lceil n_{\min}^{1/3} \rceil$  is upper bounded by*

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

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

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

542 By invoking [25, Theorem 1] we may upper bound the excess risk given a draw of  $\mathcal{S}$  by

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

543 Continuing using the definition of  $\widehat{\eta}^S$  above and because  $\eta = P_1/(P_1 + P_{-1})$  we have that,

$$\begin{aligned}
& R(\mathcal{A}_{\text{USB}}^S; P_{\text{test}}) - R(f^*; P_{\text{test}}) \\
&= 2 \int_0^1 \left| \frac{\widehat{P}_1^S(x)}{\widehat{P}_1^S(x) + \widehat{P}_{-1}^S(x)} - \frac{P_1(x)}{P_1(x) + P_{-1}(x)} \right| \left( \frac{P_1(x) + P_{-1}(x)}{2} \right) dx \\
&= \int_0^1 \left| \left( \frac{P_1(x) + P_{-1}(x)}{\widehat{P}_1^S(x) + \widehat{P}_{-1}^S(x)} \right) \widehat{P}_1^S(x) - P_1(x) \right| dx \\
&\stackrel{(i)}{\leq} \int_0^1 \left| \widehat{P}_1^S(x) - P_1(x) \right| dx + \int_0^1 \left| \frac{P_1(x) + P_{-1}(x)}{\widehat{P}_1^S(x) + \widehat{P}_{-1}^S(x)} - 1 \right| \widehat{P}_1^S(x) dx \\
&= \int_0^1 \left| \widehat{P}_1^S(x) - P_1(x) \right| dx + \int_0^1 \left| \widehat{P}_1^S(x) + \widehat{P}_{-1}^S(x) - P_1(x) - P_{-1}(x) \right| \frac{\widehat{P}_1^S(x)}{\widehat{P}_1^S(x) + \widehat{P}_{-1}^S(x)} dx \\
&\leq 2 \int_0^1 \left| \widehat{P}_1^S(x) - P_1(x) \right| dx + \int_0^1 \left| \widehat{P}_{-1}^S(x) - P_{-1}(x) \right| dx \\
&\stackrel{(ii)}{\leq} 2 \sqrt{\int_0^1 \left( \widehat{P}_1^S(x) - P_1(x) \right)^2 dx} + \sqrt{\int_0^1 \left( \widehat{P}_{-1}^S(x) - P_{-1}(x) \right)^2 dx},
\end{aligned}$$

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

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

$$\begin{aligned}
& \text{Excess Risk}(\mathcal{A}^S; (P_{\text{maj}}, P_{\text{min}})) \\
&= \mathbb{E}_{\mathcal{S}} [R(\mathcal{A}_{\text{USB}}^S; P_{\text{test}}) - R(f^*; P_{\text{test}})] \\
&\leq 2 \sqrt{\mathbb{E}_{\mathcal{S}} \left[ \int \left( \widehat{P}_1^S(x) - P_1(x) \right)^2 dx \right]} + \sqrt{\mathbb{E}_{\mathcal{S}} \left[ \int \left( \widehat{P}_{-1}^S(x) - P_{-1}(x) \right)^2 dx \right]}.
\end{aligned}$$

546 We note that  $\widehat{P}_j^S$  only depends on  $n_{\min}$  i.i.d. draws from class  $j$ . Thus by [9, Theorem 1.7], if

547  $K = c \lceil n_{\min} \rceil^{1/3}$  then

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

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

## 549 C Proof in the Group-Covariate Shift Setting

550 Throughout this section we operate in the group-covariate shift setting (Section 3.2.2).

551 First in Appendix C.1, we prove Theorem 4.2, the minimax lower bound through a sequence of  
552 lemmas. Second in Appendix C.2, we prove Theorem 5.2 that upper bound on the excess risk of the  
553 undersampled binning estimator with  $\lceil n_{\min} \rceil^{1/3}$  bins.

### 554 C.1 Proof of Theorem 4.2

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

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

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

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

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

559 and let

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

560 We can verify that

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

561 We continue to define

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

562 and

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

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

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

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

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

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

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

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

574 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  
575  $Q$  denotes the joint distribution over  $(V, S)$ .

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

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

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

579 where (i) follows by taking  $h$  to be the pointwise minimizer over  $x$ , (ii) follows since  $P_v(-1 | x) =$   
580  $1 - P_v(1 | x)$  and  $\min\{s, 1 - s\} = (1 - |1 - 2s|)/2$  for all  $s \in [0, 1]$ , and (iii) follows by Fubini's  
581 theorem which allows us to switch the order of the integrals.

582 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  
583 notation in place let us continue to upper bound integrand in the second term in the RHS above as  
584 follows:

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

585 where (i) follows since  $P_v(1 | x) = (1 + v_{j_x} \phi(x - (j_x + 1/2)/K))/2$  and by marginalizing  $Q(v | S)$   
586 over the indices  $j \neq j_x$ , (ii) follows by using Bayes' rule and (iii) follows since the total-variation  
587 distance is half the  $\ell_1$  distance. Now by the Bretagnolle–Huber inequality [see 4, Corollary 4] we get  
588 that,

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

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

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

590 **Upper bound on  $\mathfrak{B}_{\mathcal{V}}$ :** The Bayes error is

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

591 where (i) follows since  $P_v(1 | x) = 1 - P_v(-1 | x)$  and  $\min\{s, 1 - s\} = (1 - |1 - 2s|)/2$  for all  
592  $s \in [0, 1]$ , and (ii) follows by our construction of  $P_v$  above along with the fact that  $P_v(1 | x) =$   
593  $1 - P_v(-1 | x)$ .

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

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

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

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

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

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

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

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

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

605 Given samples  $S$ , let  $S = (S_1, \bar{S}_1)$  be a partition where  $S_1$  are the samples that fall in the interval  $I_1$ ,  
606 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  
607 first component and  $\bar{v}_1$  denotes the other components  $(2, \dots, K)$  of  $v$ .

608 First, we will show that

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

609 To see this, observe that

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

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

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

611 where (i) follows since by Bayes' rule

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

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

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

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

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

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

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

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

621 Now conditioned on  $v_1, n_{1,a}$  and  $n_{1,b}$ , samples in  $S_1$  are composed of 2 groups of samples ( $S_{1,a}, S_{1,b}$ ).  
622 The samples in each group ( $S_{1,a}, S_{1,b}$ ) are drawn independently from the distributions  $\mathbb{P}_a(x | x \in$   
623  $I_1)\mathbb{P}_v(y | x)$  and  $\mathbb{P}_b(x | x \in I_1)\mathbb{P}_v(y | x)$  respectively. Therefore,

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

624 where in (i) we let  $\mathbb{P}_{v_1}$  denote the conditional distribution of  $y$  for  $x \in I_1$  given  $v_1$ , (ii) follows since  
625 both  $\mathbb{P}_a$  and  $\mathbb{P}_b$  are constant in the interval, (iii) follows by our construction of  $\mathbb{P}_v$  above, and finally  
626 (iv) follows by invoking Lemma A.3 that ensures that the integral is bounded by  $1/3K^2$ .

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

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

628 Now there are  $n_{\text{maj}}$  samples from group  $a$  in  $S$  and  $n_{\text{min}}$  samples from group  $b$ . Therefore,

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

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

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

631

□

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

634 **Theorem 4.2.** Consider the group shift setting described in Section 3.2.2. Given any overlap  
635  $\tau \in [0, 1]$  recall that  $\mathcal{P}_{\text{GS}}(\tau)$  is the class of distributions such that  $\text{Overlap}(\mathbb{P}_{\text{maj}}, \mathbb{P}_{\text{min}}) \geq \tau$ . The  
636 minimax excess risk in this setting is lower bounded as follows:

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

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

638 *Proof.* First, by Lemma C.1 we know that

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

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

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

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

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

641 which completes the proof. □

642 **C.2 Proof of Theorem 5.2**

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

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

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

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

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

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

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

657 **Lemma C.3.** *The expected excess risk of undersampled binning estimator  $\mathcal{A}_{\text{USB}}$  can be decomposed*  
 658 *as follows*

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

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

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

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

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

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

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

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



663 Next, recall that the undersampled binning estimator is constant over the intervals  $I_j$  for  $j \in$   
664  $\{1, \dots, K\}$  where it takes the value  $\mathcal{A}_j^S$  (to ease notation let us simply denote it by  $\mathcal{A}_j$  below), and  
665 therefore

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

680 therefore

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

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

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

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

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

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

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

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

694 The probability of incurring risk is given by

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

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

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

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

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

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

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

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

703 Then,

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

704 where (i) follows by invoking Hoeffding's inequality[22, Proposition 2.5]. Combining this with  
705 Eqs. (28) and (29) we get that

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

706 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  
707 of  $n_{a,j}$  and  $n_{b,j}$  (the number of samples in this interval from each of the groups) are random variables  
708 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  
709  $P_b(I_j) = \int_{x \in I_j} P_b(x) dx$ . Therefore,  $n_j$  is distributed as a sum of two binomial distribution and is  
710 therefore Poisson binomially distributed [26]. Using the formula for the moment generating function  
711 (MGF) of a Poisson binomially distributed random variable we infer that,

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

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

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

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

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

714 and therefore

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

715 where (i) follows since  $|\max\{q_{j,1}^a, q_{j,1}^b\} - q_{j,1}| \leq 1/K$  by Eq. (26) and  $\gamma$  is such that  $|\gamma| \leq 1/K$ , and  
716 (ii) follows since  $(1+z)^b \leq \exp(bz)$ . Now the RHS above is maximized when  $(1 - 2q_{j,1} - 2\gamma)^2 =$   
717  $\frac{c}{n_{\text{min}}(\mathbb{P}_a(I_j) + \mathbb{P}_b(I_j))}$ , for some constant  $c$ . Plugging this into the equation above we get that

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

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

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

721 **Theorem 5.2.** Consider the group shift setting described in Section 3.2.2. For any overlap  $\tau \in [0, 1]$   
722 and for any  $(\mathbb{P}_{\text{maj}}, \mathbb{P}_{\text{min}}) \in \mathcal{P}_{\text{GS}}(\tau)$  the expected excess risk of the Undersampling Binning Estimator  
723 (Eq. (5)) with number of bins with  $K = \lceil n_{\text{min}}^{1/3} \rceil$  is

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

724 *Proof.* First by Lemma C.3 we know that

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

725 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 C.4 we get that,

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

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

728

## 729 D Experimental Details for Figure 2

730 We construct our label shift dataset from the original CIFAR10 dataset. We create a binary classi-  
731 fication task using the “cat” and “dog” classes. We use the official test examples as the balanced  
732 test set with 1000 cats and 1000 dogs. To form the initial train and validation sets, we use 2500 cat  
733 examples (half of the training set) and 500 dog examples, corresponding to a 5:1 label imbalance. We  
734 use 80% of those examples for training and the rest for validation. We are left with 2500 additional  
735 cat examples and 4500 dog examples from the original train set which we add into our training set to  
736 generate Figure 2.

737 We use the same convolutional neural network architecture as [3, 24] with random initializations for  
738 this dataset. We train this model using SGD for 400 epochs with batchsize 64, a constant learning  
739 rate 0.001 and momentum 0.9.

740 For the VS loss [13] we set  $\tau = 3$  and  $\gamma = 0.3$ , the best hyperparameters identified by Wang et al.  
741 [24] on this dataset for this neural network architecture. The importance weights used upweight the  
742 minority class samples in the training loss and validation loss is calculated to be  $\frac{\# \text{Cat Train Examples}}{\# \text{Dog Train Examples}}$ .

743 We note that all of the experiments were performed on an internal cluster on 8 GPUs.