#### Proof of the maximal inequality for IS weighted sequential empirical Α 538 processes 539

### A.1 Preliminary lemmas 540

For any sequence  $\tilde{g}_1, \ldots, \tilde{g}_T$  of conditional densities and any finite sequence  $\zeta_{1:T} := (\zeta_t)_{t=1}^T$  of 541  $\mathcal{O} \to \mathbb{R}$  functions, let 542

$$\rho_{T,\tilde{g}_{1:T}}(\zeta_{1:T}) := \left(\frac{1}{T}\sum_{t=1}^{T} \|\zeta_t\|_{2,\tilde{g}_t}^2\right)^{1/2}$$

For any conditional density  $\widetilde{g}: (a, x)\mathcal{A} \times \mathcal{X} \mapsto \widetilde{g}(a \mid x)$ , let 543

$$\rho_{T,\widetilde{g}}(f_{1:T}) := \rho_{T,\widetilde{g}_{1:T}}(f_{1:T}),$$

- where we set  $\widetilde{g}_t := \widetilde{g}$  for every  $t \in [T]$ . 544
- **Lemma 4.** Any  $\rho_{T,\tilde{g}_{1:T}}$  as defined above is a pseudonorm over the vector space  $(\mathcal{O} \to \mathbb{R})^T$ . 545
- **Lemma 5.** Consider  $g^*$  and  $g_1, \ldots, g_T$  as defined in the main text. Suppose that assumption 1 holds. Then, for any finite sequence of functions  $(\zeta_t)_{t=1}^T \in (\mathcal{O} \to \mathbb{R})^T$ , 546
- 547

$$\rho_{T,g_{1:T}}\left(\frac{g^*}{g_{1:T}}\zeta_{1:T}\right) \leq \gamma_T^{\max}\rho_{T,g^*}(\zeta_{1:T}).$$

If all elements of the sequence  $\zeta_t$  are the same, that is, if there exists  $\zeta : \mathcal{O} \to \mathbb{R}$  such that  $\zeta_t = \zeta$  for 548 every  $t \in [T]$ , then 549

$$\rho_{T,g_{1:T}}\left(\frac{g^*}{g_{1:T}}\zeta_{1:T}\right) \le \gamma_T^{avg} \|\zeta\|_{2,g^*}.$$

**Lemma 6.** Let  $\zeta_{1:T}^1, \ldots, \zeta_{1:T}^N$  be  $J \ \overline{O}_{1:T}$ -predictable sequences of  $\mathcal{O} \to \mathbb{R}$  functions, and let A be 550 an  $O_T$ -measurable event. Then, for any r > 0 and any b > 0 such that  $\max_{i \in [N], t \in [T]} \|\zeta_i^t\|_{\infty} \leq b$ , it 551 holds that 552

$$E\left[\max_{i\in[N]}\frac{1}{T}\sum_{t=1}^{N}(\delta_{O_t}-P_{g_t})\zeta_t^i \boldsymbol{I}(\rho_{T,g_{1:T}}(\zeta_{1:T}^i)\leq r) \mid A\right]$$
$$\lesssim r\sqrt{\frac{\log(1+N/P[A])}{t}}+\frac{B}{t}\log(1+N/P[A]).$$

### A.2 Proof of Theorem 1 553

*Proof of Theorem 1.* We treat together both the general case where, for each  $f, \xi_{1:T}(f)$  is an  $\overline{O}_{1:T}$ -554 predictable sequence, and the case where, for every f, there exists a deterministic  $\xi(f) : \mathcal{O} \to \mathbb{R}$ 555 such that  $\xi_t(f) = \xi(f)$  for every  $t \in [T]$ . We refer to the former as *case 1* and to the latter as *case 2* in the rest of the proof. In case 1, we let  $\tilde{\rho}_T := \rho_T^{\max}$ , and in case 2, we let  $\tilde{\rho}_T := \bar{\rho}_T$ . 556 557

From a conditional expectation bound to a high probability bound. Let x > 0. We introduce 558 the following event: 559

$$A := \left\{ \sup_{f \in \mathcal{F}} M_T(f) \ge \psi(x) \right\},\,$$

560 where

$$\begin{split} \psi(x) := & C \left\{ r^- + \sqrt{\frac{\widetilde{\gamma}_T}{T}} \int_{r^-}^r \sqrt{\log(1 + \mathcal{N}_{[]}(\epsilon, \Xi_T, \rho_{T,g^*}))} d\epsilon \right. \\ & \left. + \frac{B\gamma_T^{\max}}{T} \log(1 + \mathcal{N}_{[]}(r, \Xi_T, \rho_{T,g^*})) \right. \\ & \left. r \sqrt{\frac{x}{T}} + \frac{\gamma_T^{\max} x}{T} \right\}, \end{split}$$

where C is a universal constant to be discussed further down. Suppose we can show that

$$E\left[\sup_{f\in\mathcal{F}}M_T(f)\mid A\right] \leq \psi\left(\log\left(1+\frac{1}{P[A]}\right)\right).$$

- Then, we will have that  $\psi(x) \leq \psi(\log(2/P[A]))$ , that is  $P[A] \leq 2e^{-x}$ , which is the wished claim.
- Setting up the chaining decomposition. Let  $\epsilon_0 := r$ , and, for every  $j \ge 0$ , let  $\epsilon_j := \epsilon_0 2^{-j}$ . For any  $j \ge 0$ , let

$$\mathcal{B}^j := \left\{ (\lambda_s^{j,k}, v_t^{j,k})_{t=1}^T : k \in [N_j] \right\}$$

be a minimal  $(\epsilon_j, \rho_{T,g^*})$ -sequential bracketing of  $\Xi_T$ . For any  $f \in \mathcal{F}$ , let  $k(j, f) \in [N_j]$  be such that

$$\lambda_s^{j,k(j,f)} \le \xi_t(f) \le v^{j,k(j,f)}$$
 for every  $t \in [T]$ ,

and let  $\Delta_t^{j,f} := v_s^{j,k(j,f)} - \lambda_s^{j,k(j,f)}$  and  $u^{j,f} := v_s^{j,k(j,f)}$ . For any  $j \ge 0$ , let  $\bar{N}_j := \prod_{i=0}^j N_i$ . For any  $j \ge 0$ , and  $t \in [T]$  let

$$a_{j,t} := \epsilon_j \sqrt{\frac{T}{\log(1 + \bar{N}_j / P[A])}} \frac{\sqrt{\tilde{\gamma}_T}}{\gamma_t}$$

Let  $J \ge 0$  such that  $\epsilon_{J+1} < r^- \le \epsilon_J$ . The integer J will be the maximal depth of the chains in our chaining decomposition. For any  $t \in [T]$ ,  $f \in \mathcal{F}$ , let

$$\tau_t(f) := \inf\left\{j \ge 0 : \Delta_t^{j,f} > a_{j,t}\right\} \land J,$$

- be the depth at which we truncate the chains, adaptively depending on the value of  $\Delta_t^{j,f}$ , so that  $\Delta_t^{j,f} \mathbf{1}(\tau_t(f) > j)$  is no larger than  $a_{j,t}$  in supremum norm at any depth j.
- For any  $f \in \mathcal{F}$  and any  $t \in [T]$ , the following chaining decomposition holds:

573 Control of the tips.

• Case j = J. We have that

$$\frac{1}{T} \sum_{t=1}^{T} (\delta_{O_t} - P_{g_t}) \frac{g^*}{g_t} (\xi_t(f) - u_t^{J,f} \wedge u_t^{J-1,f}) \mathbf{1}(\tau_t(f) = J)$$

$$\leq \frac{1}{T} \sum_{t=1}^{T} P_{g_t} \frac{g^*}{g_t} \Delta_t^{J,f}$$

$$= \frac{1}{T} \sum_{t=1}^{T} \|\Delta_t^{J,f}\|_{1,g^*}$$

$$\leq \left(\frac{1}{T} \sum_{t=1}^{T} \|\Delta_t^{J,f}\|_{2,g^*}^2\right)^{1/2}$$

$$\leq \epsilon_J.$$

575

574

Therefore  

$$E\left[\sup_{f\in\mathcal{F}}\frac{1}{T}\sum_{t=1}^{T}(\delta_{O_t-P_{g_t}}\frac{g^*}{g_t})(\xi_t(f)-u_t^{J,f}\wedge u_t^{J-1,f})\mathbf{1}(\tau_t(f)=J)\mid A\right] \leq \epsilon_J.$$

576

• **Case** *j* < *J*.

$$\begin{split} \frac{1}{T} \sum_{t=1}^{T} (\delta_{O_t} - P_{g_t}) \frac{g^*}{g_t} (\xi_t(f) - u_t^{j,f} \wedge u_t^{j-1,f}) \mathbf{1}(\tau_t(f) = j) \\ \leq \frac{1}{T} \sum_{t=1}^{T} P_{g_t} \frac{g^*}{g_t} \Delta_t^{j,f} \mathbf{1}(\tau_t(f) = j) \\ \leq \frac{1}{T} \sum_{t=1}^{T} P_{g^*} \frac{(\Delta_t^{j,f})^2}{a_{j,t}} \\ \leq \epsilon_j^2 \frac{1}{T} \sum_{t=1}^{T} \frac{1}{a_{j,t}} \\ = \epsilon_j \sqrt{\frac{\log(1 + \bar{N}_j/P[A])}{T}} \frac{1}{\sqrt{\tilde{\gamma}_T}} \frac{1}{T} \sum_{t=1}^{T} \gamma_t \\ \leq \epsilon_j \sqrt{\frac{\tilde{\gamma}_T \log(1 + \bar{N}_j/P[A])}{T}}. \end{split}$$

577

(The last inequality is an equality in *case 2*).

578 **Control of the links.** We start with bounding the  $\rho_{T,g_{1:T}}$  pseudo-norm of the IS weighted links. 579 We have that

$$\rho_{T,g_{1:T}} \left( \left( \frac{g^*}{g_t} (u_t^{j,f} \wedge u_t^{j-1,f} - u_t^{j-1,f}) \right)_{t=1}^T \right) \\
\leq \rho_{T,g_{1:T}} \left( \left( \frac{g^*}{g_t} (u_t^{j,f} - u_t^{j-1,f}) \right)_{t=1}^T \right) \\
\leq \sqrt{\widetilde{\gamma}_T} \rho_{T,g^*} \left( u_{1:T}^{j,f} - u_{1:T}^{j-1,f} \right) \\
\leq \sqrt{\widetilde{\gamma}_T} \left\{ \rho_{T,g^*} \left( u_{1:T}^{j,f} - \xi_{1:T}(f) \right) + \rho_{T,g^*} \left( \xi_{1:T}(f) - u_{1:T}^{j-1,f} \right) \\
\lesssim \sqrt{\widetilde{\gamma}_T} \epsilon_j,$$

}

where we have used lemma 5 is the third line and where the fourth line above follows from the triangle inequality.

We now bound the supremum norm of the links. For every  $t \in [T]$ , 582

$$(u_t^{j,f} \wedge u_t^{j-1,f} - u_t^{j,f}) \mathbf{1}(\tau_t(f) = j)$$
  
= $(u_t^{j,f} \wedge u_t^{j-1,f} - \xi_t(f)) \mathbf{1}(\tau_t(f) = j)$   
- $(u_t^{j-1,f} - \xi(f)) \mathbf{1}(\tau_t(f) = j).$ 

Using the definition of  $\tau_t(f)$ , we obtain 583

$$0 \le (u_t^{j,f} \land u_t^{j-1,f} - \xi_t(f)) \mathbf{1}(\tau_t(f) = j) \le (u_t^{j-1,f} - \xi_t(f)) \mathbf{1}(\tau_t(f) = j) \le a_{j-1,t} \lesssim a_{j,t},$$

584 and

$$0 \le (u_t^{j-1,f} - \xi_t(f)) \mathbf{1}(\tau_t(f) = j) \le a_{j-1,t} \lesssim a_{j-1,t}.$$

Therefore, 585

$$\max_{t\in[T]} \left\| \frac{g^*}{g_t} \left( u_t^{j,f} \wedge u_t^{j-1,f} - u_t^{j-1,f} \right) \mathbf{1}(\tau_t(f) = j) \right\|_{\infty} \lesssim \gamma_t a_{j,t} = b_j$$

where 586

$$b_j := \epsilon_j \sqrt{\frac{T \widetilde{\gamma}_T}{\log(1 + \bar{N}_j / P[A])}}$$

Similarly, we have 587

> $0 \le (u_t^{j,f} - \xi_t(f)) \mathbf{1}(\tau_t(f) > j) \le a_{j,t} \qquad \text{and} \qquad 0 \le (u_t^{j-1,f} - \xi_t(f)) \mathbf{1}(\tau_t(f) > j) \le a_{j-1,t},$ and therefore, for every  $t \in [T]$

$$\left\|\frac{g^*}{g_t}\left(u_t^{j-1,f} - u_t^{j-1,f}\right)\mathbf{1}(\tau_t(f) > j)\right\|_{\infty} \lesssim \gamma_t a_{j,t} = b_j$$

Denote 589

588

$$w_t^{j,f} := \frac{g^*}{g_t} \left\{ (u_t^{j,f} \wedge u_t^{j-1,f} - u_t^{j,f}) \mathbf{1}(\tau_t(f) = j) + (u_t^{j,f} - u_t^{j-1,f}) \mathbf{1}(\tau_t(f) > j) \right\}.$$

Observe that as f varies over  $\mathcal{F}$ ,  $v_{1:T}^{j,f}$  varies over a collection of at most  $N_j \times N_{j-1} \leq \overline{N}_j$  elements. Therefore, lemma 6 yields 590

591

$$\begin{split} & E\left[\sup_{f\in\mathcal{F}}\frac{1}{T}\sum_{t=1}^{T}(\delta_{O_{t}}-P_{g_{t}})\frac{g^{*}}{g_{t}}v_{t}^{j,f}\right]\\ &\lesssim \epsilon_{j}\sqrt{\frac{\widetilde{\gamma}_{T}\log(1+\bar{N}_{j}/P[A])}{T}}+\frac{b_{j}}{T}\log(1+\bar{N}_{j}/P[A])\\ &\lesssim \epsilon_{j}\sqrt{\frac{\widetilde{\gamma}_{T}\log(1+\bar{N}_{j}/P[A])}{T}}. \end{split}$$

**Control of the root.** For any f such that  $\rho_{T,g^*}((\xi_t(f))_{t=1}^T) \leq r$ , we have that 592

$$\rho_{T,g_{1:T}}(((g^*/g_t)u_t^{0,f})_{t=1}^T) \\ \leq \sqrt{\widetilde{\gamma}_T}\rho_{T,g^*}(u_{1:T}^{0,f}) \\ \leq \sqrt{\widetilde{\gamma}_T}(\rho_{T,g^*}(u_{1:T}^{0,f} - \xi_{1:T}(f)) + \rho_{T,g^*}(\xi_{1:T}(f)))$$

Without loss of generality, we can assume that  $\max_{t \in [T]} \|u_t^{0,f}\|_{\infty} \leq B$ , since thresholding to B preserves the bracketing property. Therefore,  $\max_{t \in [T]} \|(g^*/g_t)u_t^{0,f}\|_{\infty} \leq \gamma_T^{\max}B\epsilon$ . 593 594

Then, from lemma 6, 595

$$E\left[\sup\left\{\frac{1}{T}\sum_{t=1}^{T}(\delta_{O_t} - P_{g_t})\xi_t(f) : f \in \mathcal{F}, \rho_{T,g^*}((\xi_t(f))_{t=1}^T) \le r\right\}\right]$$
$$\leq \sqrt{\frac{\widetilde{\gamma}_T}{T}}\sqrt{\log\left(\left(1 + \frac{\overline{N}_0}{P[A]}\right)} + \frac{B\gamma_T^{\max}}{T}\log\left(1 + \frac{\overline{N}_0}{P[A]}\right)$$

## 596 Adding up the bounds. We obtain

E

$$\begin{bmatrix} \sup_{f \in \mathcal{F}} M_T(f) \mid A \end{bmatrix} \lesssim \underbrace{\sqrt{\frac{\tilde{\gamma}_T}{T}} \sqrt{\log\left(\left(1 + \frac{\bar{N}_0}{P[A]}\right)} + \frac{B}{\delta T} \log\left(1 + \frac{\bar{N}_0}{P[A]}\right)}_{\text{root contribution}} + \underbrace{\sqrt{\frac{\tilde{\gamma}_T}{T}} \sum_{j=1}^J \epsilon_j \log\left(1 + \frac{\bar{N}_j}{P[A]}\right)}_{\text{links contribution}} + \underbrace{\sqrt{\frac{\tilde{\gamma}_T}{T}} \sum_{j=0}^{J-1} \epsilon_j \log\left(1 + \frac{\bar{N}_j}{P[A]}\right) + \epsilon_J}_{\text{tip contribution}} \\ \lesssim \epsilon_J + \sqrt{\frac{\tilde{\gamma}_T}{T}} \sum_{j=0}^J \epsilon_j \log\left(1 + \frac{\bar{N}_j}{P[A]}\right) + \frac{B\gamma_T^{\max}}{T} \log\left(1 + \frac{\bar{N}_0}{P[A]}\right)$$

We use the classical technique from finite adaptive chaining proofs to bound the sum in the second term with an integral [see e.g. 8, 57]. We obtain

$$\sum_{j=0}^{J} \epsilon_j \log\left(1 + \frac{\bar{N}_j}{P[A]}\right) \lesssim \int_{r^-}^r \sqrt{\log(1 + N_{[]}(\epsilon, \Xi_T, \rho_{T,g^*}))} d\epsilon + \log\left(1 + \frac{1}{P[A]}\right).$$

599 Therefore,

$$E\left[\sup_{f\in\mathcal{F}}M_{T}(f)\mid A\right] \lesssim r^{-} + \sqrt{\frac{\tilde{\gamma}_{T}}{T}} \int_{r^{-}}^{r} \sqrt{\log(1+N_{[]}(\epsilon,\Xi_{T},\rho_{T,g^{*}}))} d\epsilon$$
$$+ \frac{B\gamma_{T}^{\max}}{T} \log(1+N_{[]}(r,\Xi_{T},\rho_{T,g^{*}}))$$
$$+ \sqrt{\frac{\tilde{\gamma}_{T}}{T}} \sqrt{\log\left(1+\frac{1}{P[A]}\right)} + \frac{B\gamma_{T}^{\max}}{T} \log\left(1+\frac{1}{P[A]}\right)$$

Therefore, for an appropriate choice of the universal constant C in the definition of  $\psi$ , we have that

$$E\left[\sup_{f\in\mathcal{F}}M_T(f)\mid A\right] \leq \psi\left(\log\left(1+\frac{1}{P[A]}\right)\right),$$

which, from the first paragraph of the proof, implies the wished claim.

# 602 **B** Proof of theorem 4

*Proof.* We begin with stating a few basic facts. From the range condition on  $\mathcal{Y}$  and  $\mathcal{F}$ , the loss diameters assumption (Assumption 3) holds with  $r_0 = B = M$ , and  $\sqrt{M}$  and 4M are envelopes for  $\mathcal{F}$  and  $\ell(\mathcal{F})$ , respectively. From the assumption 6 on the entropy of  $\mathcal{F}$ , we then have

$$\log N_{[]}(\sqrt{M}\epsilon, \mathcal{F}, \|\cdot\|_{2,g^*}) \lesssim \epsilon^{-p}$$

Since  $f \mapsto \ell(f, \cdot)$  is  $\sqrt{M}$ -Lispchitz (Lemma 2), Lemma 1 yields that

$$\log N_{[]}(4M\epsilon, \ell(\mathcal{F}), \|\cdot\|_{2,g^*}) \lesssim N_{[]}(\sqrt{M}\epsilon, \ell(F), \|\cdot\|_{2,g^*}) \lesssim \epsilon^{-p}.$$

- 603 Therefore, Assumption 2 holds with envelope 4M.
- 604 Let, for any  $f \in \mathcal{F}$ ,

$$M_T(f) := \frac{1}{T} \sum_{t=1}^T (P_{Q_0,g_t} - \delta_{O_t})(\ell(f, \cdot) - \ell(f_1, \cdot)).$$

605 We distinguish the case  $p \in (0, 2)$  and the case p > 2.

**Case**  $p \in (0,2)$ . From convexity of  $f \mapsto \ell(f,\cdot)$ , and from convexity of the set  $\mathcal{F}$ , the following implication holds: for any r > 0,

$$\exists f \in \mathcal{F}, \ R^*(f) - R^*(f_1) \ge r^2 \quad \text{and} \quad \widehat{R}_T(f) - \widehat{R}_T(f_1) \le 0, \\ \Longrightarrow \exists f \in \mathcal{F}, \ R^*(f) - R^*(f_1) = r^2 \quad \text{and} \quad \widehat{R}_T(f) - \widehat{R}_T(f_1) \le 0.$$

For any r > 0, let  $\rho := r/\sqrt{M}$ , so that if  $r^2$  is an excess risk,  $\rho^2$  is the corresponding envelopestandardized excess risk. For any  $\rho > 0$ , we have

$$P\left[R^{*}(f_{T}) - R^{*}(f_{1}) \ge M\rho^{2}\right]$$
  

$$\leq P\left[\sup\left\{M_{T}(f) : f \in \mathcal{F}, R^{*}(f) - R^{*}(f_{1}) = M\rho^{2}\right\} \ge M\rho^{2}\right]$$
  

$$\leq P\left[\sup\left\{M_{T}(f) : f \in \mathcal{F}, \|\ell(f) - \ell(f_{1})\|_{2,g^{*}} \lesssim M\rho\right\} \ge M\rho^{2}\right],$$

where we have used the variance bound from Lemma 2 to obtain the last inequality. From Theorem 1, there exists C > 0 such that, for any r > 0, x > 0, it holds with probability at most  $2e^{-x}$  that

there exists C > 0 such that, for any r > 0, x > 0, it holds with probability at most  $2e^{-x}$  that  $\sup \{M_T(f) : f \in \mathcal{F}, \|\ell(f) - \ell(f_1)\|_{2,q^*} \leq M\rho\}$ 

$$\begin{split} &\leq \psi_t(M\rho) \\ &:= C\left(\sqrt{\frac{\gamma_T^{\text{avg}}}{T}}M\int_0^\rho \sqrt{\log(1+N_{[]}(\epsilon,\ell(\mathcal{F}),\|\cdot\|_{2,g^*}))}d\epsilon \\ &+ \frac{M\gamma^{\max}}{T}\log(1+N_{[]}(M\rho,\ell(\mathcal{F}),\|\cdot\|_{2,g^*})) + M\rho\sqrt{\frac{\gamma_T^{\text{avg}}x}{T}} + \frac{M\gamma_T^{\max}x}{T}\right). \end{split}$$

<sup>612</sup> Therefore, if  $\rho$  is such that  $M\rho^2$  is larger than  $\Psi_t(M\rho)$ , then, it holds that  $R^*(f) - R^*(f_1) \le M\rho^2$ <sup>613</sup> with probability at least  $1 - 2e^{-x}$ . Using the entropy bound, we obtain

$$\frac{\Psi_t(M\rho)}{M} \lesssim \sqrt{\frac{\gamma_T^{\mathrm{avg}}}{T}} \rho^{1-\frac{p}{2}} + \frac{\gamma_T^{\mathrm{max}}\rho^{-p}}{T} + \rho \sqrt{\frac{\gamma_T^{\mathrm{avg}}x}{T}} + \frac{\gamma_T^{\mathrm{max}}x}{T}.$$

614 Therefore, to have  $ho^2 \gtrsim \Psi_t(M
ho)/M$ , it suffices to have

$$\rho^2 \gtrsim \max\left\{\sqrt{\frac{\gamma_T^{\text{avg}}}{T}}\rho^{1-\frac{p}{2}}, \frac{\gamma_T^{\max}\rho^{-p}}{T}, \rho\sqrt{\frac{\gamma_T^{\text{avg}}x}{T}}, \frac{\gamma_T^{\max}x}{T}\right\},$$

615 that is

$$\rho \geq \max\left\{ \left(\frac{\gamma_T^{\max}}{T}\right)^{\frac{1}{2+p}}, \sqrt{\frac{\gamma_T^{\max}x}{T}} \right\}.$$

616 **Case** p > 2. We have that

$$\begin{split} & R^*(\widehat{f}_T) - R^*(f_1) \\ \leq \sup_{f \in \mathcal{F}} M_T(f) \\ \lesssim & M\rho^- + \sqrt{\frac{\gamma_T^{\text{avg}}}{T}} \int_{M\rho^-}^M \sqrt{\log(1 + N_{[]}(\epsilon, \ell(\mathcal{F}), \|\cdot\|_{2,g^*}))} d\epsilon \\ & + \frac{\gamma_T^{\max} M}{T} \log(1 + N_{[]}(M, \ell(\mathcal{F}), \|\cdot\|_{2,g^*})) + M\sqrt{\frac{\gamma_T^{\text{avg}} x}{T}} + \frac{M\gamma_T^{\max} x}{T} \\ \lesssim & r^- + \sqrt{\frac{\gamma_T^{\text{avg}}}{T}} M \int_{\rho^-}^1 \sqrt{\log(1 + N_{[]}(Mu, \ell(\mathcal{F}), \|\cdot\|_{2,g^*}))} du \\ & + \frac{\gamma_T^{\max} M}{T} + M\sqrt{\frac{\gamma_T^{\text{avg}} x}{T}} + \frac{M\gamma_T^{\max} x}{T} \\ \lesssim & M\rho^- + \sqrt{\frac{\gamma_T^{\text{avg}}}{T}} M \times (\rho^-)^{-(p/2-1)} + \frac{\gamma_T^{\max} M}{T} + M\sqrt{\frac{\gamma_T^{\text{avg}} x}{T}} + M\frac{\gamma_T^{\max} x}{T} \\ \lesssim & M \left(\frac{\gamma_T^{\text{avg}}}{T}\right)^{\frac{1}{p}} + \frac{\gamma_T^{\max} M}{T} + M\sqrt{\frac{\gamma_T^{\text{avg}} x}{T}} + \frac{M\gamma_T^{\max} x}{T}, \end{split}$$

<sup>617</sup> where we obtained the last line by optimizing  $\rho^-$ .

## 618 C Proof of the variance bound under margin condition

619 Proof of Lemma 3. By assumption there exists  $f_1 \in \mathcal{F}$  such that  $R^*(f_1) = \mathbb{E}_{p_X} \mu^*(X)$ . Applying

Assumption 7 with u = 0 shows that we necessarily have  $|\operatorname{argmin}_{a \in \mathcal{A}} \mu(X, a)| = 1$  almost surely.

Therefore, almost surely,  $f_1(X, a^*(X)) = 1$  and  $f_1(X, a) = 0$  for  $a \neq a^*(X)$ .

Now fix any  $f \in \mathcal{F}$ . Given X, let  $A \in \mathcal{A}$  be random variable draw from  $f(X, \cdot)$ . We will henceforth denote expectations and probabilities as wrt  $(X, A) \sim p_X \times f$ . For brevity we will also denote  $A^* = a^*(X)$ . Note that

$$\|\ell(f,\cdot) - \ell(f_1,\cdot)\|_{2,g^*}^2 \le M^2 \mathbb{P}(A^* \ne A)$$

and that

$$\|\Lambda\|_{2,g^*}^2 \left(\frac{R^*(f) - R^*(f_1)}{\|\Lambda\|_{2,g^*}}\right)^{\alpha} = M^2 (\mathbb{E}\left[\mu(X, A) - \mu(X, A^*)\right] / M)^{\nu/(\nu+1)}.$$

Denoting  $\Delta = \min_{a \in \mathcal{A} \setminus \{a^*(X)\}} \mu(X, a) - \mu^*(X)$ , Assumption 7 says that for some  $\kappa > 0$  we have  $\mathbb{P}(\Delta \le u) \le (\kappa u/M)^{\nu}$ , where  $1^{\infty} = 1$  and  $x^{\infty} = 0$  for  $x \in [0, 1)$ .

624 Fix u > 0. Then

$$\begin{split} \mathbb{E}\left[\mu(X,A) - \mu(X,A^*)\right] &= \mathbb{E}\left[(\mu(X,A) - \mu(X,A^*))\mathbf{1}(A \neq A^*)\right] \\ &\geq \mathbb{E}\left[(\mu(X,A) - \mu(X,A^*))\mathbf{1}(A \neq A^*, \Delta > u)\right] \\ &\geq u\mathbb{P}\left(A \neq A^*, \Delta > u\right) \\ &= u\left(\mathbb{P}\left(A \neq A^*\right) - \mathbb{P}\left(A \neq A^*, \Delta \le u\right)\right) \\ &\geq u\left(\mathbb{P}\left(A \neq A^*\right) - \mathbb{P}\left(\Delta \le u\right)\right) \\ &\geq u\left(\mathbb{P}\left(A \neq A^*\right) - (\kappa u/M)^{\nu}\right). \end{split}$$

Set  $u = ((\nu + 1)\kappa/M)^{-1/\nu} \mathbb{P} (A \neq A^*)^{1/\nu}$  and obtain

$$\mathbb{E}\left[\mu(X,A) - \mu(X,A^*)\right] \ge \nu(\nu+1)^{-(\nu+1)/\nu} (\kappa/M)^{-1} \mathbb{P}\left(A \neq A^*\right)^{(\nu+1)/\nu}$$

whence

$$\mathbb{P}(A \neq A^*) \le \nu^{-\nu/(\nu+1)}(\nu+1) \left( (\kappa/M) \mathbb{E} \left[ \mu(X,A) - \mu(X,A^*) \right] \right)^{\nu/(\nu+1)}.$$

We conclude that

$$\|\ell(f,\cdot) - \ell(f_1,\cdot)\|_{2,g^*}^2 \lesssim M^2 \left(\mathbb{E}\left[\mu(X,A) - \mu(X,A^*)\right]/M\right)^{\nu/(\nu+1)}$$

625 as desired.

## 626 D Additional Details and Results for the Empirical Investigation

<sup>627</sup> Here we provide additional details and results for Section 6.

## 628 D.1 Contextual Bandit Data from Multi-Class Classification Datasets

To construct our data, we turn K-class classification tasks into a K-armed contextual bandit problems 629 [15, 17, 51], which has the benefits of reproducibility using public datasets and being able to make 630 uncontroversial comparisons using actual ground truth data with counterfactuals. We use the public 631 OpenML Curated Classification benchmarking suite 2018 (OpenML-CC18; BSD 3-Clause license) 632 [11], which has datasets that vary in domain, number of observations, number of classes and number 633 of features. Among these, we select the classification datasets which have less than 60 features. This 634 results in 51 classification datasets from OpenML-CC18 used for evaluation. Table 1 summarizes the 635 characteristics of the 51 OpenML datasets used. 636

Each dataset is a collection of pairs of covariates X and labels  $L \in \{1, ..., K\}$ . We transform each dataset to the contextual bandit problem as follows. At each round, we draw  $X_t, L_t$  uniformly at random with replacement from the dataset. We reveal the context  $X_t$  to the agent, and given an arm pull  $A_t$ , we draw and return the reward  $Y_t \sim \mathcal{N}(\mathbf{1}\{A_t = L_t\}, 1)$ . To generate our data, we set T = 100000 and use the following  $\epsilon$ -greedy procedure. We pull arms uniformly at random until each



0.75 0.50 WERM 0.75 0.75 0.75 0.75 0.75 0.50 0.50 0.50 0.50 0.50 0.25 0.25 0.25 0.2 0.2 0.25 0.00 1.0 0.5 1.0

(c) CART outcome model with unrestricted tree depth.

Figure 2: Comparison of weighted regression run on contextual-bandit-collected data. Each dot is one of 51 OpenML-CC18 datasets. Lines denote  $\pm 1$  standard error. Dots are blue when ISWERM is clearly better, red when clearly worse, and black when indistinguishable within one standard error.

arm has been pulled at least once. Then at each subsequent round t, we fit  $\hat{\mu}_{t-1}$  using the data up to that time. Specifically, for each a, we take the data  $\{(X_s, Y_s) : 1 \le s \le t-1, A_s = a\}$  and pass it to a regression algorithm in order to construct  $\hat{\mu}_{t-1}(\cdot, a)$ . In Section 6, we presented results where we use sklearn's LinearRegression to fit  $\hat{\mu}_{t-1}(\cdot, a)$  (using sklearn defaults). In Appendix D.2, we repeat the experiments where we instead use sklearn's DecisionTreeRegressor (using sklearn defaults). We set  $\tilde{A}_t(x) = \operatorname{argmax}_{a=1,\ldots,K} \hat{\mu}_{t-1}(a, x)$  and  $\epsilon_t = t^{-1/3}$ . We then let  $g_t(a \mid x) =$  $\epsilon_t/K$  for  $a \neq \tilde{A}_t(x)$  and  $g_t(\tilde{A}_t(x) \mid x) = 1 - \epsilon_t + \epsilon_t/K$ . That is, with probability  $\epsilon_t$  we pull a random arm, and otherwise we pull  $\tilde{A}_t(X_t)$ .

## 650 D.2 Additional Results

In Section 6, we presented results where we use a linear-contextual  $\epsilon$ -greedy bandit algorithm to collect the data. Here, we repeat our experiments when the data are instead collected by a treecontextual  $\epsilon$ -greedy bandit algorithm, as described in Appendix D.1 above. The results are shown in Fig. 2. The conclusions are generally the same: ISWERM compares favorably for fitting linear models, while all methods perform similarly for fitting tree models.

## 656 D.3 Code and Execution Details

The IPython notebook to reproduce the experimental results of the main paper and the appendix is included as an attachment in the Supplemental Material. One needs to obtain an OpenML API key to run this code (instructions can be found at https://docs.openml.org/Python-guide/) and replace the string 'YOURKEY' in summarize\_openmlcc18() and in download\_openmlcc18() functions with it. After that, if the notebook is executed as is, it reproduces Figure 1 (38h 26min on a single Intel Xeon machine with 32 physical cores/64 CPUs). Changing variable bandit\_model from 'linear' to 'tree' reproduces Figure 2 (56h 45min on a single Intel Xeon machine with 32 physical cores/64

664 CPUs).