# On the Relationship Between Binary Classification, Bipartite Ranking, and Binary Class Probability Estimation

## Appendix

## A Proof of Theorem 4

*Proof.* Assume w.l.o.g. that  $\text{Thresh}_{D,f,c}(u) = \text{sign}(u - t^*)$  for some  $t^* \in [-\infty, \infty]$ ; a similar analysis can be shown when  $\text{Thresh}_{D,f,c}(u) = \overline{\text{sign}}(u-t^*)$  for some  $t^*$ . We first recall the following result of Clémençon et al. [8] (adapted as in [26] to account for ties and conditioning on  $y \neq y'$ ).

$$\operatorname{regret}_{D}^{\operatorname{rank}}[f] = \frac{1}{2p(1-p)} \mathbf{E}_{x,x'} \Big[ |\eta(x) - \eta(x')| \Big( \mathbf{1} \big( (f(x) - f(x'))(\eta(x) - \eta(x')) < 0 \big) \\ + \frac{1}{2} \mathbf{1} \big( f(x) = f(x') \big) \Big) \Big].$$

Next, given a binary classifier  $h: X \to \{\pm 1\}$  and a cost parameter  $c \in (0, 1)$ , the cost-sensitive classification error can be rewritten as

$$\mathbf{e} r_D^{0-1,c}[h] = \mathbf{E}_x \big[ (1-c)\eta(x) \mathbf{1} \big( h(x) = -1 \big) + c \big( 1 - \eta(x) \big) \mathbf{1} \big( h(x) = 1 \big) \big]$$

and the corresponding regret can be expanded as

$$\operatorname{regret}_{D}^{0-1,c}[h] = \mathbf{E}_{x} \left[ (1-c)\eta(x) \mathbf{1} \left( h(x) = -1 \right) + c \left( 1 - \eta(x) \right) \mathbf{1} \left( h(x) = 1 \right) \right] \\ - \mathbf{E}_{x} \left[ (1-c)\eta(x) \mathbf{1} \left( \eta(x) \le c \right) + c \left( 1 - \eta(x) \right) \mathbf{1} \left( \eta(x) > c \right) \right] \\ = \mathbf{E}_{x} \left[ \left( c - \eta(x) \right) \mathbf{1} \left( h(x) = 1, \ \eta(x) \le c \right) \right] + \mathbf{E}_{x} \left[ \left( \eta(x) - c \right) \mathbf{1} \left( h(x) = -1, \ \eta(x) > c \right) \right].$$

For  $h = \operatorname{sign} \circ (f - t^*)$ ,

$$\operatorname{regret}_{D}^{0-1,c}[\operatorname{sign} \circ (f - t^{*})] = \mathbf{E}_{x}[(c - \eta(x))\mathbf{1}(f(x) > t^{*}, \eta(x) \le c)] + \mathbf{E}_{x}[(\eta(x) - c)\mathbf{1}(f(x) \le t^{*}, \eta(x) > c)] (1) = a + b \text{ (say).}$$

We then have

$$2p(1-p) \operatorname{regret}_{D}^{\operatorname{rank}}[f] \geq \frac{1}{2} \mathbf{E}_{x,x'} \Big[ |\eta(x) - \eta(x')| \Big( \mathbf{1} \big( (f(x) - f(x'))(\eta(x) - \eta(x')) \leq 0 \big) \Big) \Big]$$
(getting rid of the term accounting for ties)  

$$\geq \frac{1}{2} \mathbf{E}_{x,x'} \Big[ |\eta(x) - \eta(x')| \Big( \mathbf{1} \big( f(x) \geq f(x'), \ \eta(x) \leq c, \ \eta(x') > c \big) \\ + \mathbf{1} \big( f(x) \leq f(x'), \ \eta(x) > c, \ \eta(x') \leq c \big) \Big) \Big]$$

$$= \frac{2}{2} \mathbf{E}_{x,x'} \Big[ |\eta(x) - \eta(x')| \Big( \mathbf{1} \big( f(x) \geq f(x'), \ \eta(x) \leq c, \ \eta(x') > c \big) \Big) \Big]$$

$$= \operatorname{term}_{1} + \operatorname{term}_{2} + \operatorname{term}_{3},$$
(2)

where

$$\begin{split} \operatorname{term}_{1} &= \mathbf{E}_{x,x'} \Big[ \big| \eta(x) - \eta(x') \big| \Big( \mathbf{1} \big( f(x) \ge f(x') > t^{*}, \ \eta(x) \le c, \ \eta(x') > c \big) \Big) \Big], \\ \operatorname{term}_{2} &= \mathbf{E}_{x,x'} \Big[ \big| \eta(x) - \eta(x') \big| \Big( \mathbf{1} \big( t^{*} \ge f(x) \ge f(x'), \ \eta(x) \le c, \ \eta(x') > c \big) \Big) \Big] \quad \text{and} \\ \operatorname{term}_{3} &= \mathbf{E}_{x,x'} \Big[ \big| \eta(x) - \eta(x') \big| \Big( \mathbf{1} \big( f(x) > t^{*}, \ f(x') \le t^{*}, \ \eta(x) \le c, \ \eta(x') > c \big) \Big) \Big]. \end{split}$$

Each of the above terms corresponds to different sets of pairs of instances; term<sub>1</sub> corresponds to pairs where both instances are ranked by f above  $t^*$ ; term<sub>2</sub> corresponds to pairs where both instances are

ranked by f below (or at the same position as)  $t^*$ ; term<sub>3</sub> corresponds to pairs (x, x'), where x is ranked by f above  $t^*$ , while x' is ranked below (or at the same position as)  $t^*$ . We next bound each of these terms separately.

 $term_1$ 

$$= \mathbf{E}_{x,x'} \Big[ |\eta(x') - c + c - \eta(x)| \Big( \mathbf{1} \big( f(x) \ge f(x') > t^*, \ \eta(x) \le c, \ \eta(x') > c \big) \Big) \Big]$$

$$\ge \mathbf{E}_{x,x'} \Big[ 2 |\eta(x') - c| |c - \eta(x)| \Big( \mathbf{1} \big( f(x) \ge f(x') > t^*, \ \eta(x) \le c, \ \eta(x') > c \big) \Big) \Big]$$

$$(since \ u + v \ge 2\sqrt{uv} \ge 2uv, \ \forall u, v \in [0, 1])$$

$$= 2 \mathbf{E}_x \Big[ |c - \eta(x)| \mathbf{1} (f(x) > t^*, \ \eta(x) \le c) \mathbf{E}_{x'} \big[ |\eta(x') - c| \mathbf{1} \big( t^* < f(x') \le f(x), \ \eta(x') > c \big) \big] \Big].$$

$$(3)$$

By definition,  $t^*$  yields the minimum classification regret among all choices of thresholds  $t \in \mathbb{R}$ :

$$t^* = \underset{t \in [-\infty,\infty]}{\operatorname{argmin}} \left\{ \operatorname{regret}_D^{0-1,c} \left[ \operatorname{sign} \circ \left( f - t \right) \right] \right\}$$
$$= \underset{t \in [-\infty,\infty]}{\operatorname{argmin}} \mathbf{E}_{x'} \left[ \left( \eta(x') - c \right) \mathbf{1} \left( f(x') \le t, \ \eta(x') > c \right) + \left( c - \eta(x') \right) \mathbf{1} \left( f(x') > t, \ \eta(x') \le c \right) \right]$$

(from Eq. (1)).

It can hence be shown that for any  $t > t^*$ ,  $\mathbf{E}_{x'} \left[ |\eta(x') - c| \mathbf{1} \left( t^* < f(x') \le t, \ \eta(x') > c \right) \right] \ge \mathbf{E}_{x'} \left[ |c - \eta(x')| \mathbf{1} \left( t^* < f(x') \le t, \ \eta(x') \le c \right) \right].$ Applying the above inequality to Eq. (3) with t = f(x), we have  $\operatorname{term}_1$ 

Similarly, one can show

$$= 2\mathbf{E}_{x} \Big[ |c - \eta(x)| \mathbf{1} \big( f(x) > t^{*}, \ \eta(x) \le c \big) \Big] \mathbf{E}_{x'} \Big[ |\eta(x') - c| \mathbf{1} \big( f(x') \le t^{*}, \ \eta(x') > c \big) \Big]$$
  
= 2ab.

Applying the bounds on term<sub>1</sub>, term<sub>2</sub> and term<sub>3</sub> in Eq. (2), we have

$$\begin{aligned} 2p(1-p) \operatorname{regret}_{D}^{\operatorname{rank}}[f] &\geq a^{2} + b^{2} + 2ab \\ &= (a+b)^{2} \\ &= \left(\operatorname{regret}_{D}^{0,1,c}[\operatorname{sign} \circ (f-t^{*})]\right)^{2}. \end{aligned}$$

Hence the proof.

## **B Proof of Theorem 6**

Proof.

The second term in the above expression can be upper bounded in terms of the ranking regret of f using Theorem 4. We now derive a bound on the first term by using standard VC-dimension based uniform convergence result for binary classification. Note that the real-valued function f, when applied to each instance drawn from D, induces a distribution over  $\mathbb{R} \times \{\pm 1\}$ ; let us call this distribution  $D_f$ . Also, let  $S_f = \{(f(x_1), y_1), \ldots, (f(x_n), y_n)\}$  be the set constructed by applying f to each instance in S; given that S is drawn iid from D, it follows that  $S_f$  is also iid drawn from  $D_f$ . Recall that  $\mathcal{T}_{inc}$  is the set of all increasing functions from  $\mathbb{R}$  to  $\{\pm 1\}$  (see Section 3). One can now view the optimization problem in (OP1) as risk minimization over  $\mathcal{T}_{inc}$  w.r.t. the distribution  $D_f$  and the optimization problem in (OP2) as empirical risk minimization over  $\mathcal{T}_{inc}$  w.r.t. the training sample  $S_f$ . In other words,

$$\inf_{\theta \in \mathcal{T}_{\text{inc}}} \left\{ \text{er}_{D}^{0-1,c} \big[ \theta \circ f \big] \right\} = \inf_{\theta \in \mathcal{T}_{\text{inc}}} \left\{ \text{er}_{D_{f}}^{0-1,c} \big[ \theta \big] \right\} = \text{er}_{D_{f}}^{0-1,c} \big[ \theta^{*} \big]$$

and

$$\inf_{t\in\mathbb{R}}\left\{\mathrm{er}_{S}^{0-1,c}\left[\mathrm{sign}\circ\left(f-t\right)\right]\right\} = \inf_{\theta\in\mathcal{T}_{\mathrm{inc}}}\left\{\mathrm{er}_{S_{f}}^{0-1,c}\left[\theta\right]\right\} = \mathrm{er}_{S_{f}}^{0-1,c}\left[\widehat{\theta}\right].$$

Thus the first term in Eq. (4) evaluates to  $\operatorname{er}_{D_f}^{0\cdot 1,c}[\widehat{\theta}] - \operatorname{er}_{D_f}^{0\cdot 1,c}[\theta^*]$ . Using standard results, one can show that the following upper bound on this quantity holds with probability at least  $1 - \delta$  (over the draw of  $S \sim D^n$ ):

$$\operatorname{er}_{D_{f}}^{0.1,c}\left[\widehat{\theta}\right] - \operatorname{er}_{D_{f}}^{0.1,c}\left[\theta^{*}\right] \leq \sqrt{\frac{32\left(\operatorname{VC-dim}(\mathcal{T}_{\operatorname{inc}})\left(\ln(2n)+1\right) + \ln\left(\frac{4}{\delta}\right)\right)}{n}}$$

where VC-dim $(\mathcal{T}_{inc})$  is the VC dimension of  $\mathcal{T}_{inc}$ . Thus with probability at least  $1 - \delta$  (over the draw of  $S \sim D^n$ ), we have

$$\operatorname{regret}_{D}^{0,1,c}[\operatorname{sign} \circ (f - \widehat{t}_{S,f,c})] \\ \leq \sqrt{\frac{32\left(\operatorname{VC-dim}(\mathcal{T}_{\operatorname{inc}})\left(\ln(2n) + 1\right) + \ln\left(\frac{4}{\delta}\right)\right)}{n}} + \sqrt{2}\sqrt{p(1-p)\operatorname{regret}_{D}^{\operatorname{rank}}[f]}.$$

It is easy to see that VC-dim( $T_{inc}$ ) = 2; plugging this in the above expression completes the proof.

### C Proof of Theorem 10

Our proof for Theorem 10 is simpler than the one in [20] which holds for a more general result. We first state and prove two lemmas which will be useful in our proof.

**Lemma 20.** Let D be a distribution over  $X \times \{\pm 1\}$ . For any binary class probability estimator  $\widehat{\eta} : X \to [0, 1]$  calibrated w.r.t. D and threshold  $t \in [0, 1]$ ,

$$\operatorname{er}_{D}^{0,1,c}\left[\operatorname{sign}\circ(\widehat{\eta}-t)\right] = \mathbf{E}_{s_{\widehat{\eta}}}\left[(1-c)s_{\widehat{\eta}}\mathbf{1}(s_{\widehat{\eta}}\leq t) + c\left(1-s_{\widehat{\eta}}\right)\mathbf{1}(s_{\widehat{\eta}}>t)\right]$$

and

$$\mathbf{r}_{D}^{0-1,c}[\overline{\operatorname{sign}} \circ (\widehat{\eta} - t)] = \mathbf{E}_{s_{\widehat{\eta}}} \big[ (1 - c) s_{\widehat{\eta}} \mathbf{1} (s_{\widehat{\eta}} < t) + c \big( 1 - s_{\widehat{\eta}} \big) \mathbf{1} (s_{\widehat{\eta}} \ge t) \big],$$

where  $s_{\hat{\eta}}$  is the random variable associated with the score distribution of  $\hat{\eta}$  over [0, 1].

*Proof.* We give a proof for the first part of the result; the second part involving sign can be proved in a similar manner. For simplicity of notation, we omit the subscript on  $s_{\hat{\eta}}$ . For any  $c \in (0, 1)$ , we have

The next lemma states that for any binary class probability estimator  $\hat{\eta}$  calibrated w.r.t. D and a given cost parameter  $c \in (0, 1)$ , the optimal classification transform on  $\hat{\eta}$  that yields minimum cost-sensitive classification error is simply  $\theta(u) = \text{sign}(u - c)$ .

**Lemma 21.** Let D be a distribution over  $X \times \{\pm 1\}$ . For any binary class probability estimator  $\widehat{\eta}: X \to [0,1]$  calibrated w.r.t. D and cost parameter  $c \in (0,1)$ ,

Thresh<sub>D, $\hat{\eta}, c$ </sub> = sign  $\circ$  ( $\hat{\eta} - c$ ).

*Proof.* Let  $s_{\widehat{\eta}}$  denote the random variable associated with the score distribution of  $\widehat{\eta}$  over [0, 1]; for simplicity of notation, we omit the subscript on  $s_{\widehat{\eta}}$ . Let us start by considering functions  $\theta \in T_{\text{inc}}$  of the form  $\theta(u) = \text{sign}(u - t)$  for some  $t \in [0, 1]$ . For any  $c \in (0, 1)$ , we have

$$\operatorname{argmin}_{t \in [0,1]} \left\{ \operatorname{er}_{D}^{0 \cdot 1, c} \left[ \operatorname{sign} \circ (\widehat{\eta} - t) \right] \right\}$$
  
=  $\operatorname{argmin}_{t \in [0,1]} \left\{ \mathbf{E}_{s} \left[ \underbrace{(1 - c)s\mathbf{1}(s \le t) + c(1 - s)\mathbf{1}(s > t)}_{\operatorname{minimum at} t = c} \right] \right\}$  (from Lemma 20)

The last step follows from the fact that the point-wise minimum is attained at t = c; this implies that  $\theta(u) = \operatorname{sign}(u - c)$  yields the least possible value of  $\operatorname{er}_D^{0-1,c}[\theta \circ \widehat{\eta}]$  over all increasing functions in  $\mathcal{T}_{\operatorname{inc}}$ , and hence we have  $\operatorname{Thresh}_{D,\widehat{\eta},c} = \operatorname{sign} \circ (\widehat{\eta} - c)$ .

We are now ready to prove Theorem 10. As before, let  $s_{\hat{\eta}}$  denote the random variable associated with the score distribution of  $\hat{\eta}$  over [0, 1]; for simplicity of notation, let us omit the subscript on  $s_{\hat{\eta}}$ .

Proof of Theorem 10. Starting with the right hand side, we have

### **D** Proof of Lemma 11

*Proof.* Expanding the left hand side, we have

## E Proof of Lemma 13

We will find it useful to introduce a few notations. For a given ranking model  $f : X \to [a, b]$  and distribution D over  $X \times \{\pm 1\}$ , define  $\bar{\mu}_f(t) = \mathbf{P}(f(x) \le t)$  and  $\bar{\eta}_f(t) = \mathbf{P}(y = 1, f(x) \le t)$  for all  $t \in [a, b]$ ; as before,  $p = \mathbf{P}(y = 1)$ .

We first state a result of [27, 28] that characterizes the minimizer of (OP3).

**Theorem 22** ( [27,28]). Let  $f : X \to [a,b]$  (where  $a, b \in \mathbb{R}$ , a < b) be any bounded-range ranking model and D be any probability distribution over  $X \times \{\pm 1\}$  such that (D, f) satisfies Assumption A. Moreover assume that  $\mu_f$  (see Assumption A), if mixed, does not have a point mass at the end-points  $a, b, and that the function \eta_f : [a,b] \to [0,1]$  defined as  $\eta_f(t) = \mathbf{P}(y = 1 \mid f(x) = t)$  is squareintegrable w.r.t. the density of the continuous part of  $\mu_f$ . Then the minimizer  $\operatorname{Cal}_{D,f} : [a,b] \to [0,1]$ of (OP3) exists, and  $\operatorname{Cal}_{D,f}(\tau)$  for any  $\tau \in (a,b)$  is given by the right-continuous slope of the largest convex minorant<sup>5</sup> of following graph at  $t = \tau$ :

$$G[f] = \{ \left( \bar{\mu}_f(t), \ \bar{\eta}_f(t) \right) : t \in [a, b] \}.$$
(5)

Moreover,  $G[\operatorname{Cal}_{D,f} \circ f]$  is piece-wise linear on all portions where it disagrees with G[f]; in particular, there exists a collection of disjoint open intervals  $\{(a_{\alpha}, b_{\alpha}) \mid \alpha \in \Lambda\}$  in [a, b], where  $\Lambda$  is some index set, such that  $\operatorname{Cal}_{D,f}$  evaluates to a constant on each such interval (with the constant being distinct for each interval) and  $\operatorname{Cal}_{D,f}$  is equal to  $\eta_f$  everywhere else in [a, b]:

$$\operatorname{Cal}_{D,f}(t) = \begin{cases} \nu_{\alpha} & \text{if } t \in (a_{\alpha}, b_{\alpha}), \text{ for some } \alpha \in \Lambda \\ \eta_{f}(t) & \text{otherwise} \end{cases},$$

where

$$\nu_{\alpha} = \frac{\bar{\eta}_f(b_{\alpha}) - \bar{\eta}_f(a_{\alpha})}{\bar{\mu}_f(b_{\alpha}) - \bar{\mu}_f(a_{\alpha})},\tag{6}$$

with  $\nu_{\alpha} \neq \nu_{\alpha'}$  for any  $\alpha \neq \alpha'$ ,  $\alpha, \alpha' \in \Lambda$ .

While the proof for the above result in [27,28] assumes a continuous and strictly positive density  $\mu_f$  over [a, b], it can be extended to handle the slightly more general conditions considered here.

We are now ready to prove the two properties stated for  $\operatorname{Cal}_{D,f}$  in Lemma 13.

*Proof of Lemma 13.* We shall assume that the score distribution of f over [a, b] is continuous, and  $\mu_f$  denotes the corresponding probability density function; a similar proof can be shown when the score distribution is discrete or is mixed and satisfies conditions stated in the Lemma. For simplicity of notation, let us denote Cal<sub>D,f</sub> as Cal.

*Proof of (1):* We need to show that for any  $u \in \text{range}(\text{Cal} \circ f)$ ,  $\mathbf{P}(y = 1 | \text{Cal}(f(x)) = u) = u$ . There are three possible cases that we could consider: (i)  $u = \nu_{\alpha}$ , for some unique  $\alpha \in \Lambda$  (see

<sup>&</sup>lt;sup>5</sup>A real-valued function  $g_1$  is a minorant of another real-valued function  $g_2$  defined over the same domain, if  $g_1(z) \le g_2(z)$ ,  $\forall z$ ; similarly,  $g_1$  is a majorant of  $g_2$ , if  $g_1(z) \ge g_2(z)$ ,  $\forall z$ .

Eq. (6)), with  $\operatorname{Cal}(t) = u, \forall t \in (a_{\alpha}, b_{\alpha})$ , and  $\operatorname{Cal}(t) \neq u$ , for all  $t \notin (a_{\alpha}, b_{\alpha})$ ; (ii)  $u \neq \nu_{\alpha}$ , for any  $\alpha \in \Lambda$ ; (iii)  $u = \nu_{\alpha}$  for some unique  $\alpha \in \Lambda$ , and there exists  $t \notin \bigcup_{\alpha \in \Lambda} (a_{\alpha}, b_{\alpha})$  with  $\operatorname{Cal}(t) = u$ .

For any  $u \in \text{range}(\text{Cal} \circ f)$  satisfying case (i), there exists  $\alpha \in \Lambda$  s.t.  $\nu_{\alpha} = u$ . We have from Eq. (6),

$$u = \frac{\eta_f(b_\alpha) - \eta_f(a_\alpha)}{\bar{\mu}_f(b_\alpha) - \bar{\mu}_f(a_\alpha)}$$
$$= \frac{\int_{a_\alpha}^{b_\alpha} \eta_f(s)\mu_f(s)ds}{\int_{a_\alpha}^{b_\alpha} \mu_f(s)ds}$$
$$= \mathbf{P}(y = 1 \mid f(x) \in (a_\alpha, b_\alpha))$$
$$= \mathbf{P}(y = 1 \mid \mathbf{Cal}(f(x)) = u).$$

The last step follows from the fact that for all  $t \notin (a_{\alpha}, b_{\alpha})$ ,  $Cal(t) \neq u$ .

For any  $u \in \text{range}(\text{Cal} \circ f)$  satisfying case (ii), there exists no  $\alpha \in \Lambda$  with  $\nu_{\alpha} = u$ ; we thus have from Theorem 22 that  $\eta_f(t) = u$  for all t with Cal(t) = u. Then

$$\begin{aligned} \mathbf{P}(y = 1 \mid \operatorname{Cal}(f(x)) = u) &= \frac{\int_{\{s : \operatorname{Cal}(s) = u\}} \eta_f(s) \mu_f(s) ds}{\int_{\{s : \operatorname{Cal}(s) = u\}} \mu_f(s) ds} \\ &= \frac{\int_{\{s : \operatorname{Cal}(s) = u\}} u \mu_f(s) ds}{\int_{\{s : \operatorname{Cal}(s) = u\}} \mu_f(s) ds} \\ &= u. \end{aligned}$$

For any  $u \in \operatorname{range}(\operatorname{Cal} \circ f)$  satisfying case (iii), there exists a unique  $\alpha \in \Lambda$  for which  $\nu_{\alpha} = u$ , with  $\operatorname{Cal}(t) = u, \forall t \in (a_{\alpha}, b_{\alpha})$ , and there also exists  $t \notin \bigcup_{\alpha \in \Lambda} (a_{\alpha}, b_{\alpha})$ , for which  $\operatorname{Cal}(t) = \eta_f(t) = u$ .

$$\begin{aligned} \mathbf{P}(y=1 \mid \operatorname{Cal}(f(x)) &= u) &= \frac{\int_{\{s \,:\, \operatorname{Cal}(s)=u\}} \eta_f(s)\mu_f(s)ds}{\int_{\{s \,:\, \operatorname{Cal}(s)=u\}} \mu_f(s)ds} \\ &= \frac{\int_{a_\alpha}^{b_\alpha} \eta_f(s)\mu_f(s)ds + \int_{\{s \,:\, \operatorname{Cal}(s)=\eta_f(s)=u\}} \eta_f(s)\mu_f(s)ds}{\int_{\{s \,:\, \operatorname{Cal}(s)=u\}} \mu_f(s)ds} \\ &= \frac{u\int_{a_\alpha}^{b_\alpha} \mu_f(s)ds + u\int_{\{s \,:\, \operatorname{Cal}(s)=\eta_f(s)=u\}} \mu_f(s)ds}{\int_{\{s \,:\, \operatorname{Cal}(s)=u\}} \mu_f(s)ds} \\ &= u_t. \end{aligned}$$

*Proof of (2):* Recall that for a ranking model f,  $er_D^{rank}[f]$  is equivalent to one minus the area under the ROC curve<sup>6</sup> (AUC) of f. It is thus enough to show that the ROC curve of Cal  $\circ f$  is a majorant for the ROC curve of f. The ROC curve for f can be defined as

$$\operatorname{ROC}[f] = \left\{ \left( \mathbf{P}(f(x) \le t \mid y = -1), \ \mathbf{P}(f(x) > t \mid y = 1) \right) : t \in [a, b] \right\} \\ = \left\{ \left( \frac{1}{1-p} \int_{a}^{t} (1-\eta_{f}(s)) \mu_{f}(s) ds, \ \frac{1}{p} \int_{t}^{b} \eta_{f}(s) \mu_{f}(s) ds \right) : t \in [a, b] \right\}.$$
(7)

As illustrated in Figure 4, each point in the graph G[f] (defined in Eq. (5)) has a corresponding point in ROC[f]; similarly, each line segment in G[f] corresponds to a line segment in ROC[f]. Moreover, for any two given ranking models  $f_1$  and  $f_2$ , if a line segment in  $G[f_1]$  is a minorant for a certain portion of  $G[f_2]$ , the corresponding line segment in ROC[ $f_1$ ] is a majorant for the corresponding portion of ROC[ $f_2$ ] (see segments AB and A'B' in Figure 4). Since, from Theorem 22, we have that  $G[\text{Cal} \circ f]$  is a minorant for G[f], and  $G[\text{Cal} \circ f]$  is piece-wise linear on all portions where it disagrees with G[f], it follows that ROC[Cal  $\circ f$ ] is a majorant for ROC[f].

<sup>&</sup>lt;sup>6</sup>The ROC curve of a ranking model f is the plot of the true positive rate (probability of classifying a random positive example as positive) against the false positive rate (probability of classifying a random negative example as positive) of a classifier of the form sign  $\circ$  (f - t) for all thresholds  $t \in [a, b]$ .

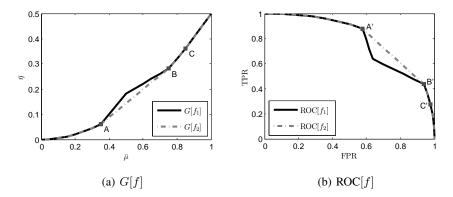


Figure 4: Sample plots illustrating the relationship between the graph G (plot of  $\bar{\eta}_f(t)$  against  $\bar{\mu}_f(t)$  for all  $t \in [a, b]$ ; see Eq. (5)) and the ROC curve (plot of true positive rate  $\text{TPR}_f(t) = \mathbf{P}(f(x) > t \mid y = 1)$  against false positive rate  $\text{FPR}_f(t) = \mathbf{P}(f(x) \le t \mid y = -1)$  for all  $t \in [a, b]$ ; see Eq. (7)). (a) Graph G for ranking models  $f_1$  and  $f_2$ : the graphs for  $f_1$  and  $f_2$  agree on all points except for the portion between points A and B, where the line segment AB in  $G[f_2]$  is a minorant for  $G[f_1]$ . (b) ROC curve for the ranking models  $f_1$  and  $f_2$ : the points A, B and C in the graph G for  $f_1$  and  $f_2$  correspond to points A', B' and C' respectively in the ROC curves for  $f_1$  and  $f_2$ ; the line segment AB in  $G[f_2]$  corresponds to the line segment A'B' in ROC[ $f_2$ ], which is a majorant for the corresponding portion in ROC[ $f_1$ ]. Moreover, while  $G[f_2]$  is a convex minorant for  $G[f_1]$ , the corresponding ROC curve ROC[ $f_2$ ] is a concave majorant for ROC[ $f_1$ ].

#### F Proof of Theorem 14

*Proof.* Using the fact that  $\operatorname{Cal}_{D,f} \circ f$  is calibrated (property 1 in Lemma 13), we have

$$\operatorname{regret}_{D}^{\operatorname{sq}}[\operatorname{Cal} \circ f] \leq \sqrt{8p(1-p)\operatorname{regret}_{D}^{\operatorname{rank}}[\operatorname{Cal}_{D,f} \circ f]} \quad (\text{from Lemma 11})$$
$$\leq \sqrt{8p(1-p)\operatorname{regret}_{D}^{\operatorname{rank}}[f]} \quad (\text{from property 2 in Lemma 13}).$$

#### G Proof of Theorem 16

Proof.

$$\operatorname{regret}_{D}^{\operatorname{sq}}[\widehat{\operatorname{Cal}}_{S,f} \circ f] = \operatorname{er}_{D}^{\operatorname{sq}}[\widehat{\operatorname{Cal}}_{S,f} \circ f] - \operatorname{er}_{D}^{\operatorname{sq}}[\eta]$$

$$= \operatorname{er}_{D}^{\operatorname{sq}}[\widehat{\operatorname{Cal}}_{S,f} \circ f] - \operatorname{er}_{D}^{\operatorname{sq}}[\operatorname{Cal}_{D,f} \circ f] + \operatorname{er}_{D}^{\operatorname{sq}}[\operatorname{Cal}_{D,f} \circ f] - \operatorname{er}_{D}^{\operatorname{sq}}[\eta]$$

$$= \left(\operatorname{er}_{D}^{\operatorname{sq}}[\widehat{\operatorname{Cal}}_{S,f} \circ f] - \operatorname{er}_{D}^{\operatorname{sq}}[\operatorname{Cal}_{D,f} \circ f]\right) + \operatorname{regret}_{D}^{\operatorname{sq}}[\operatorname{Cal}_{D,f} \circ f] \quad (8)$$

Using Theorem 14, the second term in the above expression can be upper bounded in terms of the ranking regret of f. We now focus on upper bounding the first term. As in the proof of Theorem 6, consider the distribution  $D_f$  induced by f over  $\mathbb{R} \times \{\pm 1\}$  and let  $S_f$  be the set obtained by applying f to each instance in S; clearly,  $S_f$  is iid drawn from  $D_f$ . One can then view the optimization problem in OP4 as empirical risk minimization over  $\mathcal{G}_{inc}$  w.r.t. the sample  $S_f$ . Using standard Rademacher averages based uniform convergence result for empirical risk minimization over a real-valued function class with the squared loss, we have that the following holds with probability at least  $1 - \delta$  (over the draw of  $S \sim D^n$ ):

$$\operatorname{er}_{D}^{\operatorname{sq}}[\widehat{\operatorname{Cal}}_{S,f} \circ f] - \inf_{g \in \mathcal{G}_{\operatorname{inc}}} \operatorname{er}_{D}^{\operatorname{sq}}[g \circ f] \leq 4R_{S_f}(\mathcal{G}_{\operatorname{inc}}) + 2\sqrt{\frac{2\ln\left(\frac{8}{\delta}\right)}{n}}$$

where  $R_{S_f}(\mathcal{G}_{inc})$  is the empirical Rademacher average of  $\mathcal{G}_{inc}$  w.r.t.  $S_f$ . Using Dudley's integral, and bounds on covering numbers of  $\mathcal{G}_{inc}$ , one can show  $R_{S_f}(\mathcal{G}_{inc}) \leq 24\sqrt{\frac{2\ln(n)}{n}}$  (see for example [21]); we thus have with probability at least  $1 - \delta$  (over the draw of  $S \sim D^n$ ),

$$\mathrm{er}_D^{\mathrm{sq}}[\widehat{\mathrm{Cal}}_{S,f} \circ f] - \inf_{g \in \mathcal{G}_{\mathrm{inc}}} \mathrm{er}_D^{\mathrm{sq}}[g \circ f] \leq 96 \sqrt{\frac{2\ln(n)}{n}} + 2 \sqrt{\frac{2\ln\left(\frac{8}{\delta}\right)}{n}}.$$

Plugging this into Eq. (8) (along with the upper bound on the second term) completes the proof.  $\Box$