

## A. Appendix

### A.1. Some notes on the code for reproducibility

In the code files that we provided, there are four important python files:

- *law\_data.py*: This file includes a function called *law\_data(seed)* which processes the law school admission dataset and splits the dataset randomly into training and test datasets (we keep 70% of the datapoints for training). Later, in our experiments, we set the *seed* equal to 0, 1, 2, 3, and 4 to get five different splits to repeat our experiments five times.
- *Adult\_data.py*: This file includes a function called *Adult\_dataset(seed)* which processes the adult income dataset and splits the dataset randomly into training and test datasets. Later, in our experiments, we set the *seed* equal to 0, 1, 2, 3, 4 to get five different splits to repeat our experiments five times.
- *Algorithms.py*: This file includes the following functions,
  - *ELminimizer(X0, Y0, X1, Y1, gamma, eta, model)*: This function implements ELminimizer algorithm.  $(X0, Y0)$  are the training datapoints belonging to group  $A = 0$  and  $(X1, Y1)$  are the datapoints belonging to group  $A = 1$ .  $\gamma$  is the fairness level for EL constraint.  $\eta$  is the regularizer parameter (in experiment 1 and experiment 2,  $\eta = 0.002$ ). *model* determines the model that we want to train. If *model* = "linear", then we are training a linear regression model. If *model* = "logistic", then we train a logistic regression model. This function returns five variables  $(w, b, l_0, l_1, l)$ .  $w, b$  are the weight vector and bias term of the trained model.  $l_0, l_1$  are the average training loss of group 0 and group 1, respectively.  $l$  is the overall training loss.
  - *Algorithm2(X0, Y0, X1, Y1, gamma, eta, model)*: This function implements Algorithm 2 which calls ELminimizer algorithm twice. This function also returns five variables  $(w, b, l_0, l_1, l)$ . These variables have been defined above.
  - *Algorithm3(X0, Y0, X1, Y1, gamma, eta, model)*: This function implements Algorithm 3 which find a sub-optimal solution under EL fairness. This function also returns five variables  $(w, b, l_0, l_1, l)$ . These variables have been defined above.
  - *solve\_constrained\_opt(X0, Y0, X1, Y1, eta, landa, model)*: This function uses the CVXPY package to solve the optimization problem (4). We set *landa* equal to  $\lambda_{mid}^{(i)}$  to solve the optimization problem (4) in iteration  $i$  of Algorithm 1.
  - *calculate\_loss(w, b, X0, Y0, X1, Y1, model)*: This function is used to find the test loss.  $w, b$  are model parameters (trained by Algorithm 2 or 3). It returns the average loss of group 0 and group 1 and the overall loss based on the given dataset.
  - *solve\_lin\_constrained\_opt(X0, Y0, X1, Y1, gamma, eta, model)*: This function is for solving optimization (8) after linear relaxation.
- *Baseline.py*: this file include the following functions,
  - *penalty\_method(method, X\_0, y\_0, X\_1, y\_1, num\_itr, lr, r, gamma, seed, epsilon)* where *method* can be either "linear" for linear regression or "logistic" for logistic regression. This function uses the penalty method and trains the model under EL using the Adam optimization. *num\_itr* is the maximum number of iterations. *r* is the regularization parameter (it is set to 0.002 in our experiment). *lr* is the learning rate and *gamma* is the fairness level.  $\epsilon$  is used for the stopping criterion. This function returns the trained model (which is a torch module), and training loss of group 0 and group 1, and the overall training loss. Please look at the next section for more details.
  - *fair\_batch(method, X\_0, y\_0, X\_1, y\_1, num\_itr, lr, r, alpha, gamma, seed, epsilon)*: This function is used to simulate the FairBAtch algorithm. The input parameters are similar to the input parameters of *penalty\_method* except for *alpha*. This parameter determines how to adjust the sub-sampling distribution for mini-batch formation. Please look at the next section for more details.

This function returns the trained model (which is a torch module), and training loss of group 0 and group 1, and the overall training loss.

We use the above functions to produce the results of experiment 1 and experiment 2. *table1.2.py* uses the above functions to reproduce the results in Table 1 and Table 2. *figure1.2.py* uses the above functions to reproduce figure 1 and figure 2. We

provide some comments in these files to make the code more readable. We have also provided code for training non-linear models. Please use *Table3.py* and *figure3.py* to generate the results in Table 3 and Figure 3, respectively.

Lastly, the following command to generate results in Table 1:

- `python3 table1_2.py --experiment=1 --gamma=0.0`
- `python3 table1_2.py --experiment=1 --gamma=0.1`

Use the following command to generate results in Table 2:

- `python3 table1_2.py --experiment=2 --gamma=0.0`
- `python3 table1_2.py --experiment=2 --gamma=0.1`

Use the following command to generate results in Table 3:

- `python3 table3.py --experiment=1 --gamma=0.0`
- `python3 table3.py --experiment=1 --gamma=0.1`

Use the following command to generate results in Figure 1:

- `python3 figure1_2.py --experiment=1`

Use the following command to generate results in Figure 2:

- `python3 figure1_2.py --experiment=2`

Use the following command to generate results in Figure 3:

- `python3 figure3.py --experiment=1`

Note that you need to install CVXPY, torch, and numpy.

## A.2. Notes on FairBatch (Roh et al., 2020)

This method has been proposed to find a predictor under equal opportunity, equalized odd or statistical parity. In each epoch, this method identifies the disadvantaged group and increases the subsampling rate corresponding to the disadvantaged group in mini-batch selection for the next epoch. We modify this approach for  $\gamma$ -EL as follows,

- We initialize the sub-sampling rate of group  $a$  (denoted by  $SR_a^{(0)}$ ) for mini-batch formation by  $SR_a^{(0)} = \frac{n_a}{n}, a = 0, 1$ . We Form the mini-batches using  $SR_0^{(0)}$  and  $SR_1^{(0)}$ .
- At epoch  $i$ , we run gradient descent using the mini-batches formed by  $SR_0^{(i-1)}$  and  $SR_1^{(i-1)}$ , and we obtain new model parameters  $\mathbf{w}_i$ .
- After epoch  $i$ , we calculate the empirical loss of each group. Then, we update  $SR_a^{(i)}$  as follows,

$$\begin{aligned} SR_a^{(i)} &\leftarrow SR_a^{(i-1)} + \alpha && \text{if } \hat{L}_a(\mathbf{w}_i) - \hat{L}_{1-a}(\mathbf{w}_i) > \gamma \\ SR_a^{(i)} &\leftarrow SR_a^{(i-1)} - \alpha && \text{if } \hat{L}_a(\mathbf{w}_i) - \hat{L}_{1-a}(\mathbf{w}_i) < -\gamma \\ SR_a^{(i)} &\leftarrow SR_a^{(i-1)} && \text{o.w.,} \end{aligned}$$

where  $\alpha$  is a hyperparameter and, in our experiment, is equal to 0.005 (used by (Roh et al., 2020)).

For the other hyperparameters in FairBatch, we also adopt those used in (Roh et al., 2020),

- Learning rate: 0.005
- Batch Size: 100
- Adam Optimization (Kingma and Ba, 2014) for gradient descent with default parameters in PyTorch

For stopping criteria, we stopped the learning process when the change in the objective function is less than  $10^{-6}$  between two consecutive epochs.

### A.3. Details of numerical experiments and additional numerical results

Due to the space limits of the main paper, we provide more details on our experiments here,

- Stopping criteria for penalty method and FairBatch: For stopping criteria, we stopped the learning process when the change in the objective function is less than  $10^{-6}$  between two consecutive epochs. The reason that we used  $10^{-6}$  was that we did not observe any significant change by choosing a smaller value.
- Learning rate for penalty method and FairBatch: We chose 0.005 for the learning rate. In the original implementation of FairBatch, the learning rate was 0.005 for linear models. Also, we realized that training a linear model is not sensitive to the learning rate.
- Stopping criteria for Algorithm 2 and Algorithm 3: As we stated in the main paper, we set  $\epsilon = 0.01$  in `ELminimizer` and Algorithm 3. Choosing smaller  $\epsilon$  did not change the performance significantly.
- Linear Relaxation: Note that equation (8) after linear relaxation is a convex optimization problem. We directly solve this optimization problem using CVXPY.

The experiment has been done on a system with the following configurations: 24 GB of RAM, 2 cores of P100-16GB GPU, and 2 cores of Intel Xeon CPU@2.3 GHz processor. We used GPUs for training FairBatch.

**Experiments with a non-linear model** To demonstrate how we can use our algorithms to fine-tune a non-linear model, we repeat the first experiment of Section 6 with nonlinear models. We work with Law School Admission dataset, and we train a Neural network with one hidden layer which consists of 125 neurons. We use sigmoid as the activation function for the hidden layer. We run the following algorithms for training,

- Penalty Method: We solve equation (10). Note that in this example,  $\hat{L}$  and  $\hat{L}_a$  are not convex anymore. The hyperparameters such as learning rate, penalty parameter, and stopping criterion, remain the same as in the first experiment.
- FairBatch: We train the whole network with mini-batch Adam optimization with a batch size of 100.
- Linear Relaxation: In order to take advantage of CVXPY, first, we train the network without any fairness constraint using Batch Adam optimization (batch size is equal to the size of the training dataset) and a learning rate of 0.005. Then, we fine-tune the parameters of the output layer. Note that the output layer has 126 parameters, and we fine tune those under relaxed EL fairness. In particular, we solve optimization problem (9) after linear relaxation.
- Algorithm 2 and Algorithm 3: We can run Algorithm 2 to fine-tune the network. After training the network without any constraint, we solve (9) using Algorithm 2 and Algorithm 3.

Table 3 illustrates the means and standard deviations of empirical loss and the loss difference between Black and White students. The first row specifies desired fairness level ( $\gamma = 0$  and  $\gamma = 0.1$ ) used as the input to each algorithm. Both Algorithm 2 and Algorithm 3 can achieve a fairness level (i.e.,  $|\hat{L}_0 - \hat{L}_1|$ ) close to desired fairness level  $\gamma$ . Also, we can see that the MSE of Algorithm 2 and Algorithm 3 under the nonlinear model is slightly lower than the MSE under the linear model.

Table 3: Neural Network training under EL fairness. The loss function in this example is the mean squared error loss.

		$\gamma = 0$	$\gamma = 0.1$
PM	test loss	$0.8358 \pm 0.0154$	$0.8376 \pm 0.0139$
	test $ \hat{L}_0 - \hat{L}_1 $	$0.6376 \pm 0.0948$	$0.7708 \pm 0.1207$
LinRe	test loss	$0.8430 \pm 0.0150$	$0.8190 \pm 0.0131$
	test $ \hat{L}_0 - \hat{L}_1 $	$0.6883 \pm 0.0254$	$0.5872 \pm 0.0344$
FairBatch	test loss	$0.8674 \pm 0.1032$	$0.8400 \pm 0.0522$
	test $ \hat{L}_0 - \hat{L}_1 $	$0.1808 \pm 0.0633$	$0.1995 \pm 0.1053$
ours Alg 2	test loss	$0.9116 \pm 0.0206$	$0.8552 \pm 0.0163$
	test $ \hat{L}_0 - \hat{L}_1 $	$0.0806 \pm 0.0547$	$0.1334 \pm 0.0841$
ours Alg 3	test loss	$0.9468 \pm 0.0198$	$0.8916 \pm 0.0183$
	test $ \hat{L}_0 - \hat{L}_1 $	$0.0864 \pm 0.0553$	$0.1422 \pm 0.0893$

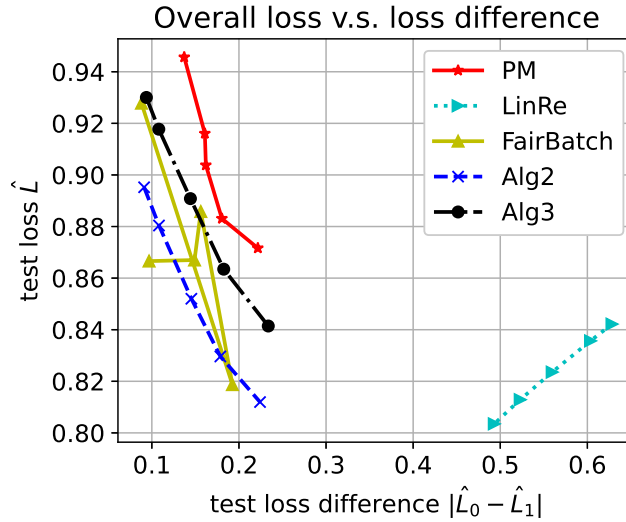


Figure 3: Trade-off between overall MSE and unfairness. A lower curve implies a better trade-off.

We also investigate how MSE  $\hat{L}$  changes as a function of fairness level  $|\hat{L}_1 - \hat{L}_0|$ . Figure 3 illustrates the MSE-fairness trade-off. To generate this plot, we change choose  $\gamma = [0.025, 0.05, 0.1, 0.15, 0.2]$ . For each  $\gamma$ , we ran the experiment 5 times and calculated the average of MSE  $\hat{L}$  and the average of MSE difference over 5 five runs using the test dataset. Based on Figure 3, we observe that FairBatch and LinRe are not very sensitive to the input  $\gamma$ . However, FairBatch may sometimes show a better trade-off than Algorithm 2. In this example, PM, Algorithm 2, and Algorithm 3 are very sensitive to  $\gamma$ , and as  $\gamma$  increases, MSE  $\hat{L}$  decreases and  $|\hat{L}_0 - \hat{L}_1|$  increases.

#### A.4. Notes on the Reduction Approach (Agarwal et al., 2018; 2019)

Let  $Q(f)$  be a distribution over  $\mathcal{F}$ . In order to find optimal  $Q(f)$  using the reduction approach, we have to solve the following optimization problem,

$$\begin{aligned} \min_Q \quad & \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \\ \text{s.t.}, \quad & \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X})) | A = 0\} = \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \\ & \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X})) | A = 1\} = \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \end{aligned}$$

Similar to (Agarwal et al., 2018; 2019), we can re-write the above optimization problem in the following form,

$$\begin{aligned} \min_Q \quad & \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \\ \text{s.t.}, \quad & \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X})) | A = 0\} - \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \leq 0 \\ & - \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X})) | A = 0\} + \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \leq 0 \\ & \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X})) | A = 1\} - \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \leq 0 \\ & - \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X})) | A = 1\} + \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \leq 0 \end{aligned}$$

Then, the reduction approach forms the Lagrangian function, in the following form,

$$\begin{aligned} L(Q, \mu) = & \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \\ & - \mu_1 \cdot \left( \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X})) | A = 0\} - \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \right) \\ & - \mu_2 \cdot \left( - \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X})) | A = 0\} + \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \right) \\ & - \mu_3 \cdot \left( \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X})) | A = 1\} - \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \right) \\ & - \mu_4 \cdot \left( - \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X})) | A = 1\} + \sum_{f \in \mathcal{F}} Q(f) \mathbb{E}\{l(Y, f(\mathbf{X}))\} \right), \\ & \mu_1 \geq 0, \mu_2 \geq 0, \mu_3 \geq 0, \mu_4 \geq 0. \end{aligned}$$

Since  $f$  is parametrized with  $\mathbf{w}$ , we can find distribution  $Q(\mathbf{w})$  over  $\mathbb{R}^{d_{\mathbf{w}}}$ . Therefore, we rewrite the problem in the following

form,

$$\begin{aligned}
 L(Q(\mathbf{w}), \mu_1, \mu_2, \mu_3, \mu_4) &= \sum_{\mathbf{w}} Q(\mathbf{w}) L(\mathbf{w}) \\
 &- \mu_1 \left( \sum_{\mathbf{w}} Q(\mathbf{w}) L_0(\mathbf{w}) - \sum_{\mathbf{w}} Q(\mathbf{w}) L(\mathbf{w}) \right) \\
 &- \mu_2 \left( - \sum_{\mathbf{w}} Q(\mathbf{w}) L_0(\mathbf{w}) + \sum_{\mathbf{w}} Q(\mathbf{w}) L(\mathbf{w}) \right) \\
 &- \mu_3 \left( \sum_{\mathbf{w}} Q(\mathbf{w}) L_1(\mathbf{w}) - \sum_{\mathbf{w}} Q(\mathbf{w}) L(\mathbf{w}) \right) \\
 &- \mu_4 \left( - \sum_{\mathbf{w}} Q(\mathbf{w}) L_1(\mathbf{w}) + \sum_{\mathbf{w}} Q(\mathbf{w}) L(\mathbf{w}) \right)
 \end{aligned}$$

The reduction approach updates  $Q(\mathbf{w})$  and  $(\mu_1, \mu_2, \mu_3, \mu_4)$  alternatively. Looking carefully at Algorithm 1 in (Agarwal et al., 2018), after updating  $(\mu_1, \mu_2, \mu_3, \mu_4)$ , we need to have access to an oracle that is able to solve the following optimization problem in each iteration,

$$\min_{\mathbf{w}} (1 + \mu_1 - \mu_2 + \mu_3 - \mu_4) L(\mathbf{w}) + (-\mu_1 + \mu_2) L_0(\mathbf{w}) + (-\mu_3 + \mu_4) L_1(\mathbf{w})$$

The above optimization problem is not convex for all  $\mu_1, \mu_2, \mu_3, \mu_4$ . Therefore, in order to use the reduction approach, we need to have access to an oracle that is able to solve the above non-convex optimization problem which is not available. Note that the original problem (1) is a non-convex optimization problem and using the reduction approach just leads to another non-convex optimization problem.

### A.5. Equalized Loss & Bounded Group Loss

In this section, we study the relation between EL and BGL fairness notions. It is straightforward to see that any predictor satisfying  $\gamma$ -BGL also satisfies the  $\gamma$ -EL. However, it is unclear to what extent an *optimal* fair predictor under  $\gamma$ -EL satisfies the BGL fairness notion. Next, we theoretically study the relation between BGL and EL fairness notions.

Let  $\mathbf{w}^*$  be denoted as the solution to (1) and  $f_{\mathbf{w}^*}$  the corresponding optimal  $\gamma$ -EL fair predictor. Theorem A.1 below shows that under certain conditions, it is impossible for both groups to experience a loss larger than  $2\gamma$  under the *optimal*  $\gamma$ -EL fair predictor.

**Theorem A.1.** *Suppose there exists a predictor that satisfies  $\gamma$ -BGL fairness notion. That is, the following optimization problem has at least one feasible point.*

$$\min_{\mathbf{w}} L(\mathbf{w}) \text{ s.t. } L_a(\mathbf{w}) \leq \gamma, \forall a \in \{0, 1\}. \quad (12)$$

Then, the followings hold,

$$\begin{aligned}
 \min\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} &\leq \gamma; \\
 \max\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} &\leq 2\gamma.
 \end{aligned}$$

Theorem A.1 shows that  $\gamma$ -EL implies  $2\gamma$ -BGL if  $\gamma$ -BGL is a feasible constraint. Therefore, if  $\gamma$  is not too small (e.g.,  $\gamma = 0$ ), then EL and BGL are not contradicting each other.

We emphasize that we are not claiming that whether EL fairness is better than BGL. Instead, these relations indicate the impacts the two fairness constraints could have on the model performance; the results may further provide the guidance for policy-makers.

## A.6. Proofs

In order to prove Theorem 3.3, we first introduce two lemmas.

**Lemma A.2.** Under assumption 3.2, there exists  $\bar{\mathbf{w}} \in \mathbb{R}^{d_{\mathbf{w}}}$  such that  $L_0(\bar{\mathbf{w}}) = L_1(\bar{\mathbf{w}}) = L(\bar{\mathbf{w}})$  and  $\lambda_{start}^{(0)} \leq L(\bar{\mathbf{w}}) \leq \lambda_{end}^{(0)}$ .

**Proof.** Let  $q_0(\beta) = L_0((1 - \beta)\mathbf{w}_{G_0} + \beta\mathbf{w}_{G_1})$  and  $q_1(\beta) = L_1((1 - \beta)\mathbf{w}_{G_0} + \beta\mathbf{w}_{G_1})$ , and  $q(\beta) = q_0(\beta) - q_1(\beta)$ ,  $\beta \in [0, 1]$ . Note that  $\nabla_{\mathbf{w}} L_a(\mathbf{w}_{G_a}) = 0$  because  $\mathbf{w}_{G_a}$  is the minimizer of  $L_a(\mathbf{w})$ .

First, we show that  $L_0((1 - \beta)\mathbf{w}_{G_0} + \beta\mathbf{w}_{G_1})$  is an increasing function in  $\beta$ , and  $L_1((1 - \beta)\mathbf{w}_{G_0} + \beta\mathbf{w}_{G_1})$  is a decreasing function in  $\beta$ . Note that  $q'_0(0) = (\mathbf{w}_{G_1} - \mathbf{w}_{G_0})^T \nabla_{\mathbf{w}} L_0(\mathbf{w}_{G_0}) = 0$ , and  $q_0(\beta)$  is convex because  $L_0(\mathbf{w})$  is convex. This implies that  $q'(\beta)$  is an increasing function, and  $q'_0(\beta) \geq 0, \forall \beta \in [0, 1]$ . Similarly, we can show that  $q'_1(\beta) \leq 0, \forall \beta \in [0, 1]$ .

Note that under Assumption (3.2),  $q(0) < 0$  and  $q(1) > 0$ . Therefore, by the intermediate value theorem, there exists  $\bar{\beta} \in (0, 1)$  such that  $q(\bar{\beta}) = 0$ . Define  $\bar{\mathbf{w}} = (1 - \bar{\beta})\mathbf{w}_{G_0} + \bar{\beta}\mathbf{w}_{G_1}$ . We have,

$$\begin{aligned} q(\bar{\beta}) &= 0 \implies L_0(\bar{\mathbf{w}}) = L_1(\bar{\mathbf{w}}) = L(\bar{\mathbf{w}}) \\ \mathbf{w}_{G_0} &\text{ is minimizer of } L_0 \implies \\ L(\bar{\mathbf{w}}) &= L_0(\bar{\mathbf{w}}) \geq \lambda_{start}^{(0)} \\ q'_0(\beta) &\geq 0, \forall \beta \in [0, 1] \implies q_0(1) \geq q_0(\bar{\beta}) \implies \\ \lambda_{end}^{(0)} &\geq L_0(\bar{\mathbf{w}}) = L(\bar{\mathbf{w}}) \end{aligned}$$

**Lemma A.3.**  $L_0(\mathbf{w}_i^*) = \lambda_{mid}^{(i)}$ , where  $\mathbf{w}_i^*$  is the solution to (4).

**Proof.** We proceed by contradiction. Assume that  $L_0(\mathbf{w}_i^*) < \lambda_{mid}^{(i)}$  (i.e.,  $\mathbf{w}_i^*$  is an interior point of the feasible set of (4)).

Notice that  $\mathbf{w}_{G_1}$  cannot be in the feasible set of (4) because  $L_0(\mathbf{w}_{G_1}) = \lambda_{end}^{(0)} > \lambda_{mid}^{(i)}$ . As a result,  $\nabla_{\mathbf{w}} L_1(\mathbf{w}_i^*) \neq 0$ . This is a contradiction because  $\mathbf{w}_i^*$  is an interior point of the feasible set of a convex optimization and cannot be optimal if  $\nabla_{\mathbf{w}} L_1(\mathbf{w}_i^*)$  is not equal to zero.

### Proof [Theorem 3.3]

Let  $I_i = [\lambda_{start}^{(i)}, \lambda_{end}^{(i)}]$  be a sequence of intervals. It is easy to see that  $I_1 \supseteq I_2 \supseteq \dots$  and  $\lambda_{end}^{(i)} - \lambda_{start}^{(i)} \rightarrow 0$  as  $i \rightarrow \infty$ . Therefore, by the Nested Interval Theorem,  $\cap_{i=1}^{\infty} I_i$  consists of exactly one real number  $\lambda^*$ , and both  $\lambda_{start}^{(i)}$  and  $\lambda_{end}^{(i)}$  converge to  $\lambda^*$ . Because  $\lambda_{mid}^{(i)} = \frac{\lambda_{start}^{(i)} + \lambda_{end}^{(i)}}{2}$ ,  $\lambda_{mid}^{(i)}$  also converges to  $\lambda^*$ .

Now, we show that  $L(\mathbf{w}^*) \in I_i$  for all  $i$  ( $\mathbf{w}^*$  is the solution to (1) when  $\gamma = 0$ ). As a result,  $L_0(\mathbf{w}^*) = L_1(\mathbf{w}^*) = L(\mathbf{w}^*)$ . Note that  $L(\mathbf{w}^*) = L_0(\mathbf{w}^*) \geq \lambda_{start}^{(0)}$  because  $\mathbf{w}_{G_0}$  is the minimizer of  $L_0$ . Moreover,  $\lambda_{end}^{(0)} \geq L(\mathbf{w}^*)$  otherwise  $L(\bar{\mathbf{w}}) < L(\mathbf{w}^*)$  ( $\bar{\mathbf{w}}$  is defined in Lemma A.2) and  $\mathbf{w}^*$  is not optimal solution under 0-EL. Therefore,  $L(\mathbf{w}^*) \in I_0$ .

Now we proceed by induction. Suppose  $L(\mathbf{w}^*) \in I_i$ . We show that  $L(\mathbf{w}^*) \in I_{i+1}$  as well. We consider two cases.

- $L(\mathbf{w}^*) \leq \lambda_{mid}^{(i)}$ . In this case  $\mathbf{w}^*$  is a feasible point for (4), and  $L_1(\mathbf{w}_i^*) = \lambda^{(i)} \leq L_1(\mathbf{w}^*) = L(\mathbf{w}^*) \leq \lambda_{mid}^{(i)}$ . Therefore,  $L(\mathbf{w}^*) \in I_{i+1}$ .
- $L(\mathbf{w}^*) > \lambda_{mid}^{(i)}$ . In this case, we proceed by contradiction to show that  $\lambda^{(i)} \geq \lambda_{mid}^{(i)}$ . Assume that  $\lambda^{(i)} < \lambda_{mid}^{(i)}$ . Define  $r(\beta) = r_0(\beta) - r_1(\beta)$ , where  $r_a(\beta) = L_a((1 - \beta)\mathbf{w}_{G_0} + \beta\mathbf{w}_i^*)$ . Note that  $\lambda^{(i)} = r_1(1)$  By Lemma A.3,  $r_0(1) = \lambda_{mid}^{(i)}$ . Therefore,  $r(1) = \lambda_{mid}^{(i)} - \lambda^{(i)} > 0$ . Moreover, under Assumption 3.2,  $r(0) < 0$ . Therefore, by the intermediate value theorem, there exists  $\bar{\beta}_0 \in (0, 1)$  such that  $r(\bar{\beta}_0) = 0$ . Similar to the proof of Lemma A.2, we can show that  $r_0(\beta)$  is an increasing function for all  $\beta \in [0, 1]$ . As a result  $r_0(\bar{\beta}_0) < r_0(1) = \lambda_{mid}^{(i)}$ . Define  $\bar{\mathbf{w}}_0 = (1 - \bar{\beta}_0)\mathbf{w}_{G_0} + \bar{\beta}_0\mathbf{w}_i^*$ . We have,

$$r_0(\bar{\beta}_0) = L_0(\bar{\mathbf{w}}_0) = L_1(\bar{\mathbf{w}}_0) = L(\bar{\mathbf{w}}_0) < \lambda_{mid}^{(i)} \quad (13)$$

$$L(\mathbf{w}^*) > \lambda_{mid}^{(i)} \quad (14)$$

The last two equations imply that  $\mathbf{w}^*$  is not a global optimal fair solution under 0-EL fairness constraint and it is not the global optimal solution to (1). This is a contradiction. Therefore, if  $L(\mathbf{w}^*) > \lambda_{mid}^{(i)}$ , then  $\lambda^{(i)} \geq \lambda_{mid}^{(i)}$ . As a result,  $L(\mathbf{w}^*) \in I_{i+1}$ .



By two above cases and the nested interval theorem, we conclude that,

$$L(\mathbf{w}^*) \in \cap_{i=1}^{\infty} I_i, \quad \lim_{i \rightarrow \infty} \lambda_{mid}^{(i)} = L(\mathbf{w}^*),$$

$$\text{define } \lambda_{mid}^{\infty} := \lim_{i \rightarrow \infty} \lambda_{mid}^{(i)}$$

Therefore,  $\lim_{i \rightarrow \infty} \mathbf{w}_i^*$  would be the solution to the following optimization problem,

$$\arg \min_{\mathbf{w}} L_1(\mathbf{w}) \text{ s.t., } L_0(\mathbf{w}) \leq \lambda_{mid}^{\infty} = L(\mathbf{w}^*)$$

By lemma A.3, the solution to above optimization problem (i.e.,  $\lim_{i \rightarrow \infty} \mathbf{w}_i^*$ ) satisfies the following,  $L_0(\lim_{i \rightarrow \infty} \mathbf{w}_i^*) = \lambda_{mid}^{\infty} = L(\mathbf{w}^*)$ . Therefore,  $\lim_{i \rightarrow \infty} \mathbf{w}_i^*$  is the global optimal solution to optimization problem (1).

**Proof [Theorem 3.4]** Let's assume that  $\mathbf{w}_O$  does not satisfy the  $\gamma$ -EL.<sup>11</sup> Let  $\mathbf{w}^*$  be the optimal weight vector under  $\gamma$ -EL. It is clear that  $\mathbf{w}^* \neq \mathbf{w}_O$ .

**Step 1.** we show that one of the following holds,

$$L_0(\mathbf{w}^*) - L_1(\mathbf{w}^*) = \gamma \quad (15)$$

$$L_0(\mathbf{w}^*) - L_1(\mathbf{w}^*) = -\gamma \quad (16)$$

Proof by contradiction. Assume  $-\gamma < L_0(\mathbf{w}^*) - L_1(\mathbf{w}^*) < \gamma$ . This implies that  $\mathbf{w}^*$  is an interior point of the feasible set of optimization problem (1). Since  $\mathbf{w}^* \neq \mathbf{w}_O$ , then  $\nabla L(\mathbf{w}^*) \neq 0$ . As a result, object function of (1) can be improved at  $\mathbf{w}^*$  by moving toward  $-\nabla L(\mathbf{w}^*)$ . This a contradiction. Therefore,  $|L_0(\mathbf{w}^*) - L_1(\mathbf{w}^*)| = \gamma$ .

**Step 2.** Function  $\mathbf{w}_{\gamma} = \text{ELminimizer}(\mathbf{w}_{G_0}, \mathbf{w}_{G_0}, \epsilon, \gamma)$  is the solution to the following optimization problem,

$$\begin{aligned} \min_{\mathbf{w}} \Pr\{A=0\}L_0(\mathbf{w}) + \Pr\{A=1\}L_1(\mathbf{w}), \\ \text{s.t., } L_0(\mathbf{w}) - L_1(\mathbf{w}) = \gamma \end{aligned} \quad (17)$$

To show the above claim, notice that the solution to optimization problem (17) is the same as the following,

$$\begin{aligned} \min_{\mathbf{w}} \Pr\{A=0\}L_0(\mathbf{w}) + \Pr\{A=1\}\tilde{L}_1(\mathbf{w}), \\ \text{s.t., } L_0(\mathbf{w}) - \tilde{L}_1(\mathbf{w}) = 0, \end{aligned} \quad (18)$$

where  $\tilde{L}_1(\mathbf{w}) = L_1(\mathbf{w}) + \gamma$ . Since  $L_0(\mathbf{w}_{G_0}) - \tilde{L}_1(\mathbf{w}_{G_0}) < 0$  and  $L_0(\mathbf{w}_{G_1}) - \tilde{L}_1(\mathbf{w}_{G_1}) > 0$ , by Theorem 3.3, we know that  $\mathbf{w}_{\gamma} = \text{ELminimizer}(\mathbf{w}_{G_0}, \mathbf{w}_{G_0}, \epsilon, \gamma)$  find the solution to (18) when  $\epsilon$  goes to zero.

Lastly, because  $|L_0(\mathbf{w}^*) - L_1(\mathbf{w}^*)| = \gamma$ , we have,

$$\mathbf{w}^* = \begin{cases} \mathbf{w}_{\gamma} & \text{if } L(\mathbf{w}_{\gamma}) \leq L(\mathbf{w}_{-\gamma}) \\ \mathbf{w}_{-\gamma} & \text{o.w.} \end{cases} \quad (19)$$

Thus, Algorithm 2 finds the solution to (1).

**Proof [Theorem 4.1]**

1. Under Assumption 3.2,  $g(1) < 0$ . Moreover,  $g(0) \geq 0$ . Therefore, by the intermediate value theorem, there exists  $\beta_0 \in [0, 1]$  such that  $g(\beta_0) = 0$ .
2. Since  $\mathbf{w}_O$  is the minimizer of  $L(\mathbf{w})$ ,  $h'(0) = 0$ . Moreover, since  $L(\mathbf{w})$  is strictly convex,  $h(\beta)$  is strictly convex and  $h'(\beta)$  is strictly increasing function. As a result,  $h'(\beta) > 0$  for  $\beta > 0$ , and  $h(\beta)$  is strictly increasing.
3. Similar to the above argument,  $s(\beta) = L_{\hat{a}}((1 - \beta)\mathbf{w}_O + \beta\mathbf{w}_{G_{\hat{a}}})$  is strictly decreasing function (notice that  $s'(1) = 0$  and  $s(\beta)$  is strictly convex).

<sup>11</sup>If  $\mathbf{w}_O$  satisfies  $\gamma$ -EL, it will be the optimal predictor under  $\gamma$ -EL fairness. Therefore, there is no need to solve any constrained optimization problem. Note that  $\mathbf{w}_O$  is the solution to problem (7).



Note that since  $h(\beta) = \Pr\{A = \hat{a}\}L_{\hat{a}}((1 - \beta)\mathbf{w}_O + \beta\mathbf{w}_{G_{\hat{a}}}) + \Pr\{A = 1 - \hat{a}\}L_{1-\hat{a}}((1 - \beta)\mathbf{w}_O + \beta\mathbf{w}_{G_{\hat{a}}})$  is strictly increasing and  $L_{\hat{a}}((1 - \beta)\mathbf{w}_O + \beta\mathbf{w}_{G_{\hat{a}}})$  is strictly decreasing. Therefore, we conclude that  $L_{1-\hat{a}}((1 - \beta)\mathbf{w}_O + \beta\mathbf{w}_{G_{\hat{a}}})$  is strictly increasing. As a result,  $g(\beta)$  should be strictly decreasing.

**Proof [Theorem 4.2]** First, we show that if  $g_\gamma(0) \leq 0$ , then  $\mathbf{w}_O$  satisfies  $\gamma$ -EL.

$$g_\gamma(0) \leq 0 \implies g(\beta) - \gamma \leq 0 \implies L_{\hat{a}}(\mathbf{w}_O) - L_{1-\hat{a}}(\mathbf{w}_O) \leq \gamma$$

Moreover,  $L_{\hat{a}}(\mathbf{w}_O) - L_{1-\hat{a}}(\mathbf{w}_O) \geq 0$  because  $\hat{a} = \arg \max_a L_a(\mathbf{w}_O)$ . Therefore,  $\gamma$ -EL is satisfied.

Now, let  $g_\gamma(0) > 0$ . Note that under Assumption 3.2,  $g_\gamma(1) = L_{\hat{a}}(\mathbf{w}_{G_{\hat{a}}}) - L_{1-\hat{a}}(\mathbf{w}_{G_{\hat{a}}}) - \gamma < 0$ . Therefore, by the intermediate value theorem, there exists  $\beta_0$  such that  $g_\gamma(\beta_0) = 0$ . Moreover, based on Theorem 4.2,  $g_\gamma$  is a strictly decreasing function. Therefore, the binary search proposed in Algorithm 3 converges to the root of  $g_\gamma(\beta)$ . As a result,  $(1 - \beta_{mid}^{(\infty)})\mathbf{w}_O + \beta_{mid}^{(\infty)}\mathbf{w}_{G_{\hat{a}}}$  satisfies  $\gamma$ -EL. Note that since  $g(\beta)$  is strictly decreasing, and  $g(\beta_{mid}^{(\infty)}) = \gamma$ , and  $\beta_{mid}^{(\infty)}$  is the smallest possible  $\beta$  under which  $(1 - \beta)\mathbf{w}_O + \beta\mathbf{w}_{G_{\hat{a}}}$  satisfies  $\gamma$ -EL. Since  $h$  is increasing, the smallest possible  $\beta$  gives us a better accuracy.

**Proof [Theorem 4.3]** If  $g_\gamma(0) \leq 0$ , then  $\mathbf{w}_O$  satisfies  $\gamma$ -EL, and  $\underline{\mathbf{w}} = \mathbf{w}_O$ . In this case, it is easy to see that  $L(\mathbf{w}_O) \leq \max_{a \in \{0,1\}} L_a(\mathbf{w}_O)$  (because  $L(\mathbf{w}_O)$  is a weighted average of  $L_0(\mathbf{w}_O)$  and  $L_1(\mathbf{w}_O)$ ).

Now assume that  $g_\gamma(0) > 0$ . Note that if we prove this theorem for  $\gamma = 0$ , then the theorem will hold for  $\gamma > 0$ . This is because the **optimal** predictor under 0-EL satisfies  $\gamma$ -EL condition as well. In other words, 0-EL is a stronger constraint than  $\gamma$ -EL.

Let  $\gamma = 0$ . In this case, Algorithm 3 finds  $\underline{\mathbf{w}} = (1 - \beta_0)\mathbf{w}_O + \beta_0\mathbf{w}_{G_{\hat{a}}}$ , where  $\beta_0$  is defined in Theorem 4.1. We have,

$$(*) \quad g(\beta_0) = 0 = L_{\hat{a}}(\underline{\mathbf{w}}) - L_{1-\hat{a}}(\underline{\mathbf{w}})$$

In the proof of theorem 4.1, we showed that  $L_{\hat{a}}((1 - \beta)\mathbf{w}_O + \beta\mathbf{w}_{G_{\hat{a}}})$  is decreasing in  $\beta$ . Therefore, we have,

$$(**) \quad L_{\hat{a}}(\underline{\mathbf{w}}) \leq L_{\hat{a}}(\mathbf{w}_O)$$

Therefore, we have,

$$L(\underline{\mathbf{w}}) = \Pr(A = 0) \cdot L_{\hat{a}}(\underline{\mathbf{w}}) + (1 - \Pr(A = 1)) \cdot L_{1-\hat{a}}(\underline{\mathbf{w}}) \quad (20)$$

$$\text{(By (*))} \quad = L_{\hat{a}}(\underline{\mathbf{w}}) \quad (21)$$

$$\text{(By (**))} \quad \leq L_{\hat{a}}(\mathbf{w}_O) \quad (22)$$

**Proof [Theorem 4.5]**

By the triangle inequality, the following holds,

$$\sup_{f_{\mathbf{w}} \in \mathcal{F}} ||L_0(\mathbf{w}) - L_1(\mathbf{w})| - |\hat{L}_0(\mathbf{w}) - \hat{L}_1(\mathbf{w})|| \leq \quad (23)$$

$$\sup_{f_{\mathbf{w}} \in \mathcal{F}} |L_0(\mathbf{w}) - \hat{L}_0(\mathbf{w})| + \sup_{f_{\mathbf{w}} \in \mathcal{F}} |L_1(\mathbf{w}) - \hat{L}_1(\mathbf{w})|. \quad (24)$$

Therefore, with probability at least  $1 - 2\delta$  we have,

$$\begin{aligned} \sup_{f_{\mathbf{w}} \in \mathcal{F}} ||L_0(\mathbf{w}) - L_1(\mathbf{w})| - |\hat{L}_0(\mathbf{w}) - \hat{L}_1(\mathbf{w})|| &\leq \\ &B(\delta, n_0, \mathcal{F}) + B(\delta, n_1, \mathcal{F}) \end{aligned} \quad (25)$$

As a result, with probability  $1 - 2\delta$  holds,

$$\begin{aligned} \{\mathbf{w} | f_{\mathbf{w}} \in \mathcal{F}, |L_0(\mathbf{w}) - L_1(\mathbf{w})| \leq \gamma\} &\subseteq \\ \{\mathbf{w} | f_{\mathbf{w}} \in \mathcal{F}, |\hat{L}_0(\mathbf{w}) - \hat{L}_1(\mathbf{w})| \leq \hat{\gamma}\} \end{aligned} \quad (26)$$

Now consider the following,

$$L(\hat{\mathbf{w}}) - L(\mathbf{w}^*) = L(\hat{\mathbf{w}}) - \hat{L}(\hat{\mathbf{w}}) + \hat{L}(\hat{\mathbf{w}}) - \hat{L}(\mathbf{w}^*) + \hat{L}(\mathbf{w}^*) - L(\mathbf{w}^*) \quad (27)$$

By (26),  $\hat{L}(\hat{\mathbf{w}}) - \hat{L}(\mathbf{w}^*) \leq 0$  with probability  $1 - 2\delta$ . Thus, with probability at least  $1 - 2\delta$ , we have,

$$L(\hat{\mathbf{w}}) - L(\mathbf{w}^*) \leq L(\hat{\mathbf{w}}) - \hat{L}(\hat{\mathbf{w}}) + \hat{L}(\mathbf{w}^*) - L(\mathbf{w}^*). \quad (28)$$

Therefore, under assumption 4.4, we can conclude with probability at least  $1 - 6\delta$ ,  $L(\hat{\mathbf{w}}) - L(\mathbf{w}^*) \leq 2B(\delta, n, \mathcal{F})$ . In addition, by (25), with probability at least  $1 - 2\delta$ , we have,

$$\begin{aligned} |L_0(\hat{\mathbf{w}}) - L_1(\hat{\mathbf{w}})| &\leq B(\delta, n_0, \mathcal{F}) + B(\delta, n_1, \mathcal{F}) + |\hat{L}_0(\mathbf{w}) - \hat{L}_1(\mathbf{w})| \\ &\leq \hat{\gamma} + B(\delta, n_0, \mathcal{F}) + B(\delta, n_1, \mathcal{F}) \\ &= \gamma + 2B(\delta, n_0, \mathcal{F}) + 2B(\delta, n_1, \mathcal{F}) \end{aligned}$$

**Proof [Theorem A.1]** Let  $\tilde{\mathbf{w}}$  be a feasible point of optimization problem (12), then  $\tilde{\mathbf{w}}$  is also a feasible point to (1).

We proceed by contradiction. We consider three cases,

- If  $\min\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} > \gamma$  and  $\max\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} > 2\gamma$ . In this case,  

$$L(\mathbf{w}^*) > \gamma \geq L(\tilde{\mathbf{w}}).$$

This is a contradiction because it implies that  $\mathbf{w}^*$  is not an optimal solution to (1), and  $\tilde{\mathbf{w}}$  is a better solution for (1).

- If  $\min\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} > \gamma$  and  $\max\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} \leq 2\gamma$ . This case is similar to above.  $\min\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} > \gamma$  implies that  $L(\mathbf{w}^*) > \gamma \geq L(\tilde{\mathbf{w}})$ . This is a contradiction because it implies that  $\mathbf{w}^*$  is not an optimal solution to (1).

- If  $\min\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} \leq \gamma$  and  $\max\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} > 2 \cdot \gamma$ . We have:  

$$\max\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} - \min\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} > \gamma,$$
 which shows that  $\mathbf{w}^*$  is not a feasible point for (1). This is a contradiction.

Therefore,  $\max\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} \leq 2\gamma$  and  $\min\{L_0(\mathbf{w}^*), L_1(\mathbf{w}^*)\} \leq \gamma$ .