

---

# Demystifying Local and Global Fairness Trade-offs in Federated Learning Using Information Theory

---

Faisal Hamman<sup>1</sup> Sanghamitra Dutta<sup>1</sup>

<sup>1</sup>University of Maryland, College Park

## Abstract

We present an information-theoretic perspective to group fairness trade-offs in federated learning (FL) with respect to sensitive attributes, such as, gender, race, etc. Existing works mostly focus on either *global fairness* (overall disparity of the model across all clients) or *local fairness* (disparity of the model at each individual client), without necessarily considering their trade-offs. There is a lack of understanding of the interplay between global and local fairness in FL, and if and when one implies the other. To address this gap, we leverage a body of work in information theory called partial information decomposition (PID) which first identifies three sources of unfairness in FL, namely, *Unique Disparity*, *Redundant Disparity*, and *Masked Disparity*. Using canonical examples, we demonstrate how these three disparities contribute to global and local fairness. This decomposition helps us derive fundamental limits and trade-offs between global or local fairness, particularly under data heterogeneity, as well as, derive conditions under which one implies the other. We also present experimental results on real-world datasets to support our theoretical findings. This work offers a more nuanced understanding of the sources of disparity in FL that can inform the use of local disparity mitigation techniques, and their convergence and effectiveness when deployed in practice.

## 1. Introduction

With the growing use of FL in various high-stakes applications, such as finance, healthcare, recommendation systems, etc., it is crucial to ensure that these models do not discriminate against any demographic group based on sensitive features (Smith et al., 2016). While there are several methods to achieve group fairness in the centralized settings (?), these methods do not directly apply to a FL setting since each client only has access to their own local dataset, and hence, is restricted to only performing local disparity mitigation.

Some recent works (Du et al., 2021; Abay et al., 2020; Ezzeldin et al., 2023) focus on developing models that are fair when evaluated on the entire dataset across all clients, a concept known as *global fairness*. For example, several banks have decided to engage in a FL process to train a model that will determine loan qualifications without exchanging data among them. A globally fair model is one that does not discriminate against any protected group, when evaluated on the entire dataset across all the banks. On the other hand, *local fairness* considers the disparity of the model at each individual client, i.e., when evaluated on a client’s local dataset. Local fairness is an important consideration, as the models will ultimately be deployed and used at the local client (Cui et al., 2021).

One might notice that global and local fairness evaluation can differ from each other when the local demographics at a client differ from the global demographics across the entire dataset (data heterogeneity, e.g., a bank with predominantly White customers). Previous research has mostly focused on achieving either global fairness (Du et al., 2021) or local fairness (Cui et al., 2021), without considering their trade-offs and interplay. We provide more related works in Appendix B. In this work, we aim to provide a fundamental understanding of group fairness trade-offs in the FL setting. Our main contributions can be summarized as follows:

- **Partial information decomposition (PID) of global and local disparity into three sources of unfairness:** We formalize the notion of global and local fairness in federated learning using information theory. We first define global disparity as the mutual information between a model’s prediction (denoted by  $\hat{Y}$ ) and the sensitive attribute (denoted by  $Z$ ), i.e.,  $I(Z; \hat{Y})$  (Definition 1). Then, we show that local disparity can be represented as the conditional mutual information  $I(Z; \hat{Y}|S)$  where  $S$  denotes the client (Definition 2). We also demonstrate relationships between these information-theoretic terms and well-known fairness metrics such as statistical parity (see Lemma 1).

Next, we propose a PID that breaks down the global and local disparity into three components: *Unique Disparity*, *Redundant Disparity*, and *Masked Disparity*. We provide canonical examples to help understand these three sources of disparities in the context of FL (see Section 3.1).

- Fundamental limits and trade-offs between local and global fairness:** With the use of the decomposed disparities, we have been able to uncover the fundamental information-theoretic limits and trade-offs between global and local disparities. We show the limitations of achieving global fairness using local fairness due to the redundant disparity (see Theorem 1) and the limitations of achieving local fairness using global fairness due to the masked disparity (see Theorem 2).
- Understanding scenarios where local fairness implies global fairness and vice versa:** We also identify the conditions under which one form of fairness (local or global) implies the other. Specifically, we have established conditions under which local fairness can result in global fairness (Theorem 4) and conditions under which global fairness can result in local fairness (Theorem 5).
- Experimental demonstrations:** We provide experimental evaluations on the Adult dataset to validate our theoretical findings. We demonstrate practical scenarios where unique, redundant, and masked disparities are prevalent.

## 2. Preliminaries

Let  $K$  be the total number of federating clients. A client is represented as  $S \in [K]$  where  $[K] = \{1, 2, \dots, K\}$ . A client  $S = k$  has a dataset  $\mathcal{D}_k = \{(x_i, y_i, z_i)\}_{i=1, \dots, n_k}$  where  $x_i$  denotes the input features,  $y_i \in \{0, 1\}$  is the true label,  $z_i \in \{0, 1\}$  is the sensitive attribute (1 for privileged group, 0 for unprivileged group), and  $n_k$  denotes the number of datapoints at client  $S = k$ . The entire dataset is given by  $\hat{\mathcal{D}} = \cup_{k=1}^K \mathcal{D}_k$ . When denoting a random variable drawn from this dataset, we let  $X$  denote the input features,  $Z$  denote the sensitive attribute, and  $Y$  denote the true label. We also let  $\hat{Y}$  represent the predictions of a model  $f_\theta(X)$ , parameterized by  $\theta$ .

Next, explore the concept of group fairness in the context of FL and formally discuss two prevalent perspectives of fairness: global and local.

**Definition 1** (Global Disparity). *The global disparity of a model  $f_\theta$  with respect to  $Z$  is defined as  $I(Z; \hat{Y})$ , the mutual information between  $Z$  and  $\hat{Y}$  (where  $\hat{Y} = f_\theta(X)$ ).*

This is related to a widely-used group fairness notion called statistical parity. Existing works (Hardt et al., 2016) define the global statistical parity of the model  $f_\theta$  as:  $\Pr(\hat{Y}=1|Z=1) = \Pr(\hat{Y}=1|Z=0)$ . Global statistical parity is satisfied when  $Z$  is independent of  $\hat{Y}$ , which is equivalent to zero mutual information  $I(Z; \hat{Y}) = 0$ . To further justify our choice of  $I(Z; \hat{Y})$  as a measure of global disparity, we provide a relationship between the absolute statistical parity gap and mutual information when they are non-zero in Lemma 1 (Proof in Appendix C).

**Lemma 1.** *Let  $Z$  and  $\hat{Y}$  be binary and  $\Pr(Z = 0) =$*

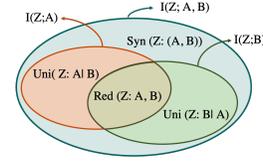


Figure 1. Venn diagram showing PID of  $I(Z; (A, B))$ .

$1 - \Pr(Z = 1) = \alpha$ . The global statistical parity gap  $SP_{global} = |\Pr(\hat{Y} = 1|Z = 1) - \Pr(\hat{Y} = 1|Z = 0)|$  is bounded by  $\frac{\sqrt{0.5 I(Z; \hat{Y})}}{2\alpha(1-\alpha)}$ .

Similarly, existing literature (Ezzeldin et al., 2023) defines local statistical parity at a client  $k$  as:  $\Pr(\hat{Y}=1|Z=1, S=k) = \Pr(\hat{Y}=1|Z=0, S=k)$ . A critical observation that we make in this work is that: *local client unfairness can be quantified as the conditional mutual information  $I(Z, \hat{Y}|S)$ .*

**Definition 2** (Local Disparity). *The local disparity is defined as  $I(Z; \hat{Y}|S)$ , the mutual information between  $Z$  and  $\hat{Y}$  conditioned on  $S$ .*

**Lemma 2.**  $I(Z, \hat{Y}|S) = 0$  if and only if  $\Pr(\hat{Y} = 1|Z = 1, S = k) = \Pr(\hat{Y} = 1|Z = 0, S = k)$  at all clients.

The proof (see Appendix C) uses the fact that  $I(Z, \hat{Y}|S) = \sum_{k=1}^K \Pr(S=k) I(Z, \hat{Y}|S=k)$  where  $I(Z, \hat{Y}|S=k)$  is the local mutual information at client  $k$ , and  $\Pr(S=k) = \frac{n_k}{n}$ , the proportion of data points at client  $k$ . For the rest of this paper, we use  $I(Z; \hat{Y})$  to denote the global disparity and  $I(Z, \hat{Y}|S)$  to denote the local disparity.

### 2.1. Background on Partial Information Decomposition

The Partial Information Decomposition (PID) decomposes the mutual information  $I(Z; (A, B))$  about a random variable  $Z$  contained in the tuple  $(A, B)$  into four non-negative terms as follows :

$$I(Z; (A, B)) = \text{Uni}(Z, A|B) + \text{Uni}(Z, B|A) + \text{Red}(Z; A, B) + \text{Syn}(Z; (A, B)) \quad (1)$$

Here,  $\text{Uni}(Z, A|B)$  denotes the unique information about  $Z$  that is present only in  $A$  and not in  $B$ ,  $\text{Red}(Z; (A, B))$  denotes the redundant information about  $Z$  that is present in both  $A$  and  $B$ , and  $\text{Syn}(Z; (A, B))$  denotes the synergistic information not present in either of  $A$  or  $B$  individually, but present jointly in  $(A, B)$ . Observe in Fig. 1 that  $\text{Uni}(Z; A|B)$  can be viewed as the information-theoretic sub-volume of the intersection between  $I(Z; A)$  and  $I(Z; A|B)$ . Similarly for  $\text{Red}(Z; (A, B))$ . Defining any one of the PID terms suffices to get the others. Hence, we include a popular definition of  $\text{Uni}(Z; A|B)$  from (Bertschinger et al., 2014) in Appendix C. We also include an example to better understand PID in Appendix C.

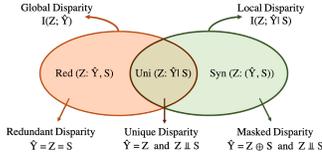


Figure 2. Venn diagram showing PID of Global and Local Disparity with canonical examples where each disparity is maximum.

### 3. Main Results

#### 3.1. Partial Information Decomposition of Disparity in Federated Learning

**Proposition 1.** *The global and local disparity in a federated setting can be decomposed into PID terms as follows:*

$$I(Z; \hat{Y}) = \text{Uni}(Z; \hat{Y}|S) + \text{Red}(Z; \hat{Y}, S). \quad (2)$$

$$I(Z; \hat{Y}|S) = \text{Uni}(Z; \hat{Y}|S) + \text{Syn}(Z; (\hat{Y}, S)). \quad (3)$$

We also refer to Fig. 2 for a pictorial illustration of this result. The proof follows directly from the relationship between different PID terms.

The term  $\text{Uni}(Z; \hat{Y}|S)$  captures the information about sensitive attribute  $Z$  that is present only in the model prediction  $\hat{Y}$  and not in the clients  $S$ . We refer to this as the **Unique Disparity**. It is important to note that this does not provide a complete picture of the model’s disparity, as  $S$  may also contain redundant information with  $\hat{Y}$  about  $Z$ .

The term  $\text{Red}(Z; \hat{Y}, S)$  denotes the redundant information about sensitive attribute  $Z$  that is present in both prediction  $\hat{Y}$  and client  $S$ . We call this as the **Redundant Disparity**. The unique and redundant disparities make up the global disparity  $I(Z; \hat{Y})$ .

The term  $\text{Syn}(Z; (\hat{Y}, S))$  represents the synergistic information about sensitive attribute  $Z$  that is not present in either  $\hat{Y}$  or  $S$  individually, but is present jointly in  $(\hat{Y}, S)$ . We refer to this as the **Masked Disparity**, as it is only observed when  $\hat{Y}$  and  $S$  are considered together. The unique and masked disparities make up the local disparity  $I(Z; \hat{Y}|S)$ .

**Canonical Examples.** Assume binary model predictions, sensitive attributes, and two clients, i.e.,  $\hat{Y}, Z, S \in \{0, 1\}$ . Note that  $I(Z; (\hat{Y}, S)) = H(Z) - H(Z|\hat{Y}, S) \leq H(Z) = 1$ , i.e., the maximum disparity is 1 bit for this case.

**Example 1 (Pure Uniqueness).** *Let  $\hat{Y} = Z$  and  $Z \perp S$ . The model accepts only males from each client dataset. This model is both locally and globally unfair. Here,  $\text{Red}(Z; \hat{Y}, S) = 0$ ,  $\text{Uni}(Z; \hat{Y}|S) = 1$ , and  $\text{Syn}(Z; (\hat{Y}, S)) = 0$*

Each client has the same proportion of privileged and un-

privileged groups,  $Z \perp S$ . If the model makes its predictions solely based on the sensitive attribute, i.e.,  $\hat{Y} = Z$ , the unique disparity,  $\text{Uni}(Z; \hat{Y}|S) = 1$ , since all information about  $Z$  is encoded in  $\hat{Y}$  and none is present in  $S$ .  $\text{Red}(Z; \hat{Y}, S) = 0$ , since  $\hat{Y}$  and  $S$  do not share any information about  $Z$ . Similarly,  $\text{Syn}(Z; (\hat{Y}, S)) = 0$ , since jointly,  $(\hat{Y}, S)$  do not contain any information about  $Z$ . As a result, both the global and local disparities,  $I(Z; \hat{Y}) = I(Z; \hat{Y}|S) = 1$ , indicating that the model is globally and locally unfair.

**Example 2 (Pure Redundancy).** *Let  $\hat{Y} = Z = S$ . Client  $S = 0$  has all females with negative predicted outcomes and client  $S = 1$  has all males with positive predicted outcomes. It is clear that this model achieves local fairness at each client however it is globally unfair. In terms of PID,  $\text{Red}(Z; \hat{Y}, S) = 1$ ,  $\text{Uni}(Z; \hat{Y}|S) = 0$ , and  $\text{Syn}(Z; (\hat{Y}, S)) = 0$ .*

The sensitive attributes are skewed across clients, with client  $S = 0$  containing only females ( $Z = 0$ ) and client  $S = 1$  containing only males ( $Z = 1$ ). The model makes its predictions based on these sensitive attributes. In this case, the redundant disparity,  $\text{Red}(Z; \hat{Y}, S) = 1$ , since all information about  $Z$  is contained in both  $\hat{Y}$  and  $S$ . The unique disparity,  $\text{Uni}(Z; \hat{Y}|S) = 0$ , since there is no information about  $Z$  in  $\hat{Y}$  that is not present in  $S$ . Similarly,  $\text{Syn}(Z; (\hat{Y}, S)) = 0$ , since jointly,  $(\hat{Y}, S)$  do not contain any information about  $Z$  that is not present in  $\hat{Y}$  and  $S$  individually. As a result, the global disparity,  $I(Z; \hat{Y}) = 1$ , and the local disparity,  $I(Z; \hat{Y}|S) = 0$ , indicate that the model is globally unfair but locally fair. It is not surprising that the model is globally unfair, as the predictions are based on the sensitive attributes. However, since each client has only one protected group, the model exhibits local fairness. In generally, redundant disparity is observed when  $Z - S - \hat{Y}$  forms a *Markov chain*, but  $Z$  and  $\hat{Y}$  are correlated.

**Example 3 (Pure Synergy).** *Let  $\hat{Y} = Z \oplus S$  and  $Z \perp S$ . The model accepts males from client  $S = 0$  and females from client  $S = 1$ , while others are rejected. This model is not locally fair but globally fair. Here,  $\text{Red}(Z; \hat{Y}, S) = 0$ ,  $\text{Uni}(Z; \hat{Y}|S) = 0$ , and  $\text{Syn}(Z; (\hat{Y}, S)) = 1$ .*

The sensitive attributes are identically distributed across clients, i.e.,  $Z \perp S$ . The model prediction is an XOR of the sensitive attribute  $Z$  and clients  $S$ , i.e.,  $\hat{Y} = Z \oplus S$ . Thus, the model accepts males ( $Z = 1$ ) from client  $S = 0$  and females ( $Z = 1$ ) from client  $S = 1$ , while rejecting all other individuals. The masked disparity  $\text{Syn}(Z; (\hat{Y}, S)) = 1$ , since  $(\hat{Y}, S)$  specifies information about  $Z$  that is not specified either  $\hat{Y}$  or  $S$ . With both  $\hat{Y}$  and  $S$  having no information about  $Z$ , i.e.,  $I(Z; S) = I(Z; \hat{Y}) = 0$ , it follows that there can not be any unique or redundant disparity. The model is locally unfair, with  $I(Z; \hat{Y}|S) = 1$ , but globally fair, with  $I(Z; \hat{Y}) = 0$ . The model achieves global fairness

by balancing the local unfairness at each client.

These examples demonstrate pure uniqueness, redundancy, and synergy. In practice, it is usually a combination of these cases as we show in the experimental section. These also extend to multiple clients and sensitive attributes.

### 3.2. Fundamental Limits and Tradeoffs

**Limitations in Achieving Global Fairness with Local Fairness.** As clients only have access to their own datasets, applying local disparity mitigation methods at each client can be convenient. Cui et al. (2021) argue for local fairness as models are deployed at the local client level. In Theorem 1, we formally demonstrate the limitations of this approach. Even if local clients are able to use some optimal local mitigation methods and model aggregation techniques to achieve local fairness, the global disparity may still be greater than zero.

**Theorem 1** (Impossibility of Using Local Fairness to Attain Global Fairness). *As long as redundant disparity  $\text{Red}(Z:\hat{Y}, S) > 0$ , the global disparity  $\text{I}(Z;\hat{Y}) > 0$  even if local disparity goes to 0.*

The proof leverages PID of global disparity as shown in (2). Recall example 2, where a locally fair model would fail to be globally fair due to a non-zero redundant disparity.

**Limitation in Achieving Local Fairness with Global Fairness.** We now consider the scenario where a model is trained to achieve global fairness and is subsequently deployed at the local client level.

**Theorem 2** (Global Fairness Does Not Imply Local Fairness). *As long as masked disparity  $\text{Syn}(Z:(\hat{Y}, S)) > 0$ , there exist scenarios where global fairness is attained but local fairness is not.*

The proof leverages the decomposition of local disparity as shown in (3). This demonstrates that while it is possible to train a model to achieve global fairness, it may still exhibit disparity when deployed at the local level due to the canceling of disparities between clients, recall example 3 (Pure Synergy).

**Definition 3** (Interaction Information). *The difference between global and local disparity is:  $\text{I}(Z;\hat{Y}) - \text{I}(Z;\hat{Y}|S) = \text{I}(Z;\hat{Y};S)$ . This term is the “interaction information.”*

Interaction information quantifies the redundancy and synergy present in a system. Positive interaction information indicates a system with high levels of redundancy and global disparity, while negative interaction information indicates a system with high levels of synergy and local disparity. Interaction information can inform the trade-off between local and global disparity.

**Towards Achieving Global Fairness with Local Fairness.**

**Theorem 3** (Necessary and Sufficient Condition to Achieve Global Fairness Using Local Fairness). *If local disparity  $\text{I}(Z,\hat{Y}|S)$  goes to zero, then the global disparity  $\text{I}(Z;\hat{Y})$  also goes to zero, if and only if the redundant disparity  $\text{Red}(Z:\hat{Y}, S) = 0$ .*

**Lemma 3.**  $Z \perp\!\!\!\perp S \implies \text{Red}(Z:\hat{Y}, S) = 0$ . *When  $Z$  and  $S$  are independent, the redundant disparity is zero.*

The results of Theorem 3 and Lemma 3 suggest that when the proportion of each protected group is equal across all clients, the redundant disparity will decrease to zero. Hence, when the local disparity goes to zero, the global disparity will also decrease to zero. However, in practice, this proportion is fixed since the dataset at each client cannot be changed, i.e.,  $\text{I}(Z;S)$  is fixed. Therefore, we explore another more controllable condition to eliminate redundant disparity even when  $\text{I}(Z;S) > 0$ .

**Lemma 4.** *If synergistic disparity  $\text{Syn}(Z:(\hat{Y}, S)) = 0$ , the redundant disparity  $\text{Red}(Z:\hat{Y}, S) = 0$  if  $\hat{Y}$  and  $S$  are independent  $\hat{Y} \perp\!\!\!\perp S$  or  $\text{I}(\hat{Y};S) = 0$ , even if  $\text{I}(Z;S) > 0$ .*

**Theorem 4.** *If local disparity goes to zero, then the global disparity also goes to zero, if the model prediction  $\hat{Y}$  is independent of  $S$ , i.e.,  $\text{I}(\hat{Y};S) = 0$ .*

Theorem 4 demonstrates that, in order to reduce redundant disparity and achieve global fairness when there is a strong correlation between  $Z$  and  $S$ , one potential solution is to enforce independence between  $\hat{Y}$  and  $S$ . This means that the model should make predictions at the same rate across all clients. The proofs are provided in Appendix D.

**Towards Achieving Local Fairness with Global Fairness.**

**Theorem 5.** *Local disparity will always be less than global disparity if masked disparity  $\text{Syn}(Z:(\hat{Y}, S)) = 0$ .*

**Corollary 1.** *The local disparity will always be less than global disparity if  $Z, \hat{Y}, S$  form a Markov chain  $Z - \hat{Y} - S$ .*

**Experimental Demonstration:** To validate our theoretical findings, we experiment on a real-world dataset (Adult) in Appendix A. We demonstrate practical scenarios where unique, redundant, and masked disparities are prevalent, for two or more clients. We analyze the PID under various data heterogeneity scenarios with varying sensitive attribute distributions and varying synergy levels across clients. We use PID functions from python `dit` package (James et al., 2018) to estimate the terms.

**Conclusions and Extended Work:** This work provides a more nuanced understanding of the sources of disparity in FL than no other unfairness measure provides. This can inform the use of bias mitigation techniques, and the effectiveness of models when deployed in practice. Extended work would investigate how PID decomposition could be directly estimated in a federated setting as well as extensions to other fairness metrics, such as equalized odds.

## References

- Abay, A., Zhou, Y., Baracaldo, N., Rajamoni, S., Chuba, E., and Ludwig, H. Mitigating bias in federated learning. *arXiv preprint arXiv:2012.02447*, 2020.
- Alghamdi, W., Hsu, H., Jeong, H., Wang, H., Michalak, P., Asoodeh, S., and Calmon, F. Beyond adult and compas: Fair multi-class prediction via information projection. In *Advances in Neural Information Processing Systems*, volume 35, pp. 38747–38760. Curran Associates, Inc., 2022.
- Baharlouei, S., Nouiehed, M., Beirami, A., and Raza-viyayn, M. Renyi fair inference. *arXiv preprint arXiv:1906.12005*, 2019.
- Bertschinger, N., Rauh, J., Olbrich, E., Jost, J., and Ay, N. Quantifying unique information. *Entropy*, 16(4):2161–2183, 2014.
- Calmon, F., Wei, D., Vinzamuri, B., Natesan Ramamurthy, K., and Varshney, K. R. Optimized pre-processing for discrimination prevention. *Advances in neural information processing systems*, 30, 2017.
- Cho, J., Hwang, G., and Suh, C. A fair classifier using mutual information. In *2020 IEEE international symposium on information theory (ISIT)*, pp. 2521–2526. IEEE, 2020.
- Chu, L., Wang, L., Dong, Y., Pei, J., Zhou, Z., and Zhang, Y. Fedfair: Training fair models in cross-silo federated learning. *arXiv preprint arXiv:2109.05662*, 2021.
- Cui, S., Pan, W., Liang, J., Zhang, C., and Wang, F. Addressing algorithmic disparity and performance inconsistency in federated learning. *Advances in Neural Information Processing Systems*, 34:26091–26102, 2021.
- Du, W., Xu, D., Wu, X., and Tong, H. Fairness-aware agnostic federated learning. In *Proceedings of the 2021 SIAM International Conference on Data Mining (SDM)*, pp. 181–189. SIAM, 2021.
- Dua, D. and Graff, C. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- Dutta, S. and Hamman, F. A review of partial information decomposition in algorithmic fairness and explainability. *Entropy*, 25(5):795, 2023.
- Dutta, S., Venkatesh, P., Mardziel, P., Datta, A., and Grover, P. An information-theoretic quantification of discrimination with exempt features. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020a.
- Dutta, S., Wei, D., Yueksel, H., Chen, P. Y., Liu, S., and Varshney, K. R. Is there a trade-off between fairness and accuracy? a perspective using mismatched hypothesis testing. In *International Conference on Machine Learning (ICML)*, pp. 2803–2813, 2020b.
- Dutta, S., Venkatesh, P., Mardziel, P., Datta, A., and Grover, P. Fairness under feature exemptions: Counterfactual and observational measures. *IEEE Transactions on Information Theory*, 67(10):6675–6710, 2021.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pp. 214–226, 2012.
- Ehrlich, D. A., Schneider, A. C., Wibral, M., Priesemann, V., and Makkeh, A. Partial information decomposition reveals the structure of neural representations. *arXiv preprint arXiv:2209.10438*, 2022.
- Ezzeldin, Y. H., Yan, S., He, C., Ferrara, E., and Avestimehr, A. S. Fairfed: Enabling group fairness in federated learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 7494–7502, 2023.
- Galhotra, S., Shanmugam, K., Sattigeri, P., and Varshney, K. R. Causal feature selection for algorithmic fairness. In *Proceedings of the 2022 International Conference on Management of Data*, pp. 276–285, 2022.
- Ghassami, A., Khodadadian, S., and Kiyavash, N. Fairness in supervised learning: An information theoretic approach. In *2018 IEEE international symposium on information theory (ISIT)*, pp. 176–180. IEEE, 2018.
- Grari, V., Ruf, B., Lamprier, S., and Detyniecki, M. Fairness-aware neural renyi minimization for continuous features. *arXiv preprint arXiv:1911.04929*, 2019.
- Hardt, M., Price, E., and Srebro, N. Equality of opportunity in supervised learning. *Advances in neural information processing systems*, 29, 2016.
- Hu, S., Wu, Z. S., and Smith, V. Provably fair federated learning via bounded group loss. *arXiv preprint arXiv:2203.10190*, 2022.
- James, R. G., Ellison, C. J., and Crutchfield, J. P. dit: a Python package for discrete information theory. *The Journal of Open Source Software*, 3(25):738, 2018. doi: <https://doi.org/10.21105/joss.00738>.
- Kairouz, P., Liao, J., Huang, C., Vyas, M., Welfert, M., and Sankar, L. Generating fair universal representations using adversarial models. *arXiv preprint arXiv:1910.00411*, 2019.

- Kamishima, T., Akaho, S., and Sakuma, J. Fairness-aware learning through regularization approach. In *2011 IEEE 11th International Conference on Data Mining Workshops*, pp. 643–650. IEEE, 2011.
- Kamishima, T., Akaho, S., Asoh, H., and Sakuma, J. Fairness-aware classifier with prejudice remover regularizer. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2012, Bristol, UK, September 24-28, 2012. Proceedings, Part II 23*, pp. 35–50. Springer, 2012.
- Li, T., Sanjabi, M., Beirami, A., and Smith, V. Fair resource allocation in federated learning. *arXiv preprint arXiv:1905.10497*, 2019.
- Liang, P. P., Cheng, Y., Fan, X., Ling, C. K., Nie, S., Chen, R., Deng, Z., Mahmood, F., Salakhutdinov, R., and Morency, L.-P. Quantifying & modeling feature interactions: An information decomposition framework. *arXiv preprint arXiv:2302.12247*, 2023.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and Arcas, B. A. y. Communication-Efficient Learning of Deep Networks from Decentralized Data. In Singh, A. and Zhu, J. (eds.), *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, volume 54 of *Proceedings of Machine Learning Research*, pp. 1273–1282. PMLR, 20–22 Apr 2017.
- Mohamadi, S., Doretto, G., and Adjeroh, D. A. More synergy, less redundancy: Exploiting joint mutual information for self-supervised learning. *arXiv preprint arXiv:2307.00651*, 2023.
- Pakman, A., Nejatbakhsh, A., Gilboa, D., Makkeh, A., Mazzucato, L., Wibral, M., and Schneidman, E. Estimating the unique information of continuous variables. *Advances in neural information processing systems*, 34: 20295–20307, 2021.
- Papadaki, A., Martinez, N., Bertran, M., Sapiro, G., and Rodrigues, M. Minimax demographic group fairness in federated learning. *arXiv preprint arXiv:2201.08304*, 2022.
- Pessach, D. and Shmueli, E. A review on fairness in machine learning. *ACM Comput. Surv.*, 55(3), feb 2022.
- Rodríguez-Gálvez, B., Granqvist, F., van Dalen, R., and Seigel, M. Enforcing fairness in private federated learning via the modified method of differential multipliers. *arXiv preprint arXiv:2109.08604*, 2021.
- Shi, Y., Yu, H., and Leung, C. A survey of fairness-aware federated learning. *arXiv preprint arXiv:2111.01872*, 2021.
- Smith, M., Munoz, C., and Patil, D. J. Big Risks, Big Opportunities: the Intersection of Big Data and Civil Rights, 2016.
- Tax, T. M., Mediano, P. A., and Shanahan, M. The partial information decomposition of generative neural network models. *Entropy*, 19(9):474, 2017.
- Wang, H., Hsu, H., Diaz, M., and Calmon, F. P. The impact of split classifiers on group fairness. In *2021 IEEE International Symposium on Information Theory (ISIT)*, pp. 3179–3184. IEEE Press, 2021.
- Wollstadt, P., Schmitt, S., and Wibral, M. A rigorous information-theoretic definition of redundancy and relevancy in feature selection based on (partial) information decomposition. *J. Mach. Learn. Res.*, 24:131–1, 2023.
- Zeng, Y., Chen, H., and Lee, K. Improving fairness via federated learning. *arXiv preprint arXiv:2110.15545*, 2021.
- Zhang, D. Y., Kou, Z., and Wang, D. Fairfl: A fair federated learning approach to reducing demographic bias in privacy-sensitive classification models. In *2020 IEEE International Conference on Big Data (Big Data)*, pp. 1051–1060, 2020.

## A. Experimental Demonstrations

In this section, we provide experimental evaluations on real-world datasets to validate our theoretical findings. We investigate the PID of global and local fairness under various conditions and scenarios.

**Dataset.** UCL *Adult* dataset (Dua & Graff, 2017), which comprises over 40,000 data points. The objective is to predict whether the annual earnings of an individual is more than 50K per year. We select *gender* as a sensitive attribute, with Male as  $Z = 1$  and Female as  $Z = 0$ .

**Evaluation.** We define global and local disparities as mutual information measures. To estimate these values, we use the python `dit` package (James et al., 2018), which includes PID functions that allows us to decompose the global and local disparities into unique, redundant, and masked disparities. We implement the definition of unique information from (Bertschinger et al., 2014).

**Demonstrating the Disparities.** First, we demonstrate scenarios with unique, redundant, and masked disparities for the *Adult* dataset using two clients,  $S = 0, 1$ . We do this by strategically splitting the dataset between the clients and training our federated model using *FedAvg* (McMahan et al., 2017). For context, a model in the centralized case, achieved an accuracy of 84.67% and disparity  $I(Z; \hat{Y}) = 0.03537$ .

**Scenario 1 Unique Disparity on Adult dataset.** Unique disparity is observed when the sensitive attribute is evenly distributed across clients ( $Z \perp\!\!\!\perp S$ ). To achieve this, we randomly distribute the dataset among the two clients to ensure  $I(Z; S) = 0$ . The global and local disparity is 0.0359 bits. Decomposing these we get the unique disparity as 0.0359 and zero redundant and masked disparity, indicating that the source of the disparity is solely from the dependence between the model predictions and sensitive attributes, and not from  $S$ . This matches the centralized case.

**Scenario 2 Redundant Disparity on Adult dataset.** Redundant disparity occurs when there is high heterogeneity of sensitive attributes across clients ( $Z \approx S$ ). To achieve this, we distribute the dataset so as client  $S = 0$  contains mainly females  $Z = 0$  and client  $S = 1$  contains mainly males,  $Z = 1$ , resulting in  $I(Z; S) = 0.8486$ . The model trained had a global disparity of 0.0431 and a local disparity of 0.0014. By decomposing these, we find that the redundant disparity is 0.0431, masked disparity is 0.0014, and zero unique disparity.

**Scenario 3 Masked Disparity on Adult dataset.** The masked disparity is observed when the model predictions  $\hat{Y} = Z \oplus S$ . To attain this, we distribute the dataset such that the first client dataset contains males ( $Z = 1$ ) with true labels  $Y = 1$  and females ( $Z = 0$ ) with true labels  $Y = 0$ . The second client dataset contains the remaining (males with  $Y = 0$  and females with  $Y = 1$ ). The trained model had a local disparity of 0.1761 and a global disparity of 0.0317. With a masked disparity of 0.1761, redundant disparity of 0.0317, and zero unique disparity. The non-zero redundant disparity is due to the way we split the data, which resulted in  $I(Z; S) = 0.2409$ .

We summarize the three scenarios in Fig. 3. Additionally, we evaluate the effects of using local fairness mitigation technique on the various disparities present. This is achieved by incorporating a statistical parity regularizer at each individual client. The results are presented in Table 1.

Table 1. Table illustrates the effects of using a naive local disparity mitigation technique on the various scenarios. It proved efficacious only when Unique disparity is present (scenario 1). However, with high redundancy or synergy (scenarios 2 & 3), the utilization of the disparity mitigation technique exacerbated disparities.

	Loc.	Glob.	Uniq.	Red.	Mas.
Scenario 1	0.0359	0.0359	0.0359	0.0000	0.0000
+ fairness	0.0062	0.0062	0.0062	0.0000	0.0000
Scenario 2	0.0014	0.0431	0.0000	0.0431	0.0014
+ fairness	0.0110	0.0626	0.0000	0.0626	0.0110
Scenario 3	0.1761	0.0317	0.0000	0.0317	0.1761
+ fairness	0.0935	0.0418	0.0053	0.0365	0.0882

**PID of Disparity under Heterogeneous Sensitive Attribute Distribution.** We analyze the PID of local and global disparities under different sensitive attribute distributions across clients. We train the model with two clients, each having equal-sized datasets. We use  $\alpha = \Pr(Z=0|S=0)$  to represent sensitive attribute heterogeneity. Note that for a fixed  $\alpha$ , the proportions of sensitive attributes at the other client are fixed. For example since  $\Pr(Z=0)=0.33$  for the *Adult* dataset,  $\alpha=0.33$  results in even distribution of sensitive attribute across the two clients. Our results are in Fig. 4 and Table 2.

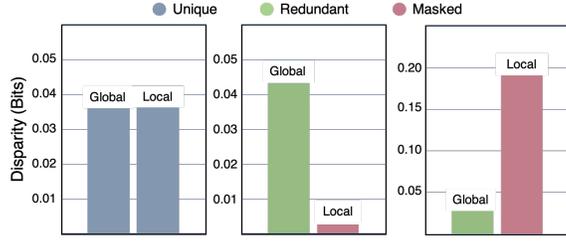


Figure 3. Plot demonstrating scenarios with unique, redundant, and masked disparities for the Adult dataset two client case. (left) Unique disparity when sensitive attributes are equally distributed across clients. (center) Redundant disparity when there is a dependency between clients and sensitive attributes. (right) Masked disparity when model predictions  $\hat{Y} \approx Z \oplus S$ .

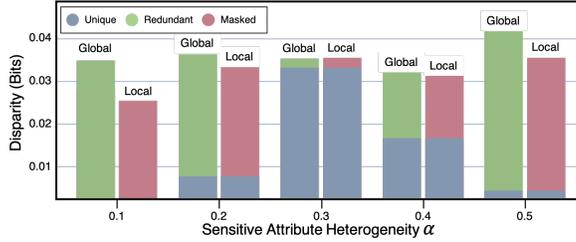


Figure 4. Plot illustrates PID of global and local fairness varying levels of sensitive attribute heterogeneity ( $\alpha$ ) and two clients case. When  $\alpha$  is close to 0.3, the data is split evenly across clients, resulting in a higher level of unique disparity. As  $\alpha$  deviates from 0.3, i.e., higher dependency between  $Z$  and  $S$ , the unique disparity decreases while redundant and masked disparity increase.

Table 2. The PID of global and local disparity for varying sensitive attribute heterogeneity  $\alpha$ .

$\alpha$	$I(Z; S)$	Loc.	Glob.	Uniq.	Red.	Mas.
0.1	0.1877	0.0262	0.0342	0.000	0.0342	0.0262
0.2	0.0575	0.0336	0.0364	0.0064	0.0301	0.0273
0.3	0.0032	0.0363	0.0365	0.0332	0.0032	0.0031
0.4	0.0154	0.0311	0.0319	0.0186	0.0133	0.0125
0.5	0.0957	0.0368	0.0413	0.0023	0.0390	0.0345

**Observing Levels of Masked Disparity.** Here, we aim to gain a deeper understanding of the circumstances leading to masked disparities. Through scenario 3, we showed how high masked disparities can occur. However, the level of synergy portrayed in the example may not always be present in reality. We attempt to quantify this using a metric *synergy level*. The synergy level ( $\lambda$ ) measures how closely the model prediction  $\hat{Y}$  aligns with the XOR of  $Z$  and  $S$ . A value of 1 represents perfect alignment, while a value of 0 indicates independence. To achieve a high synergy level, we apply the method outlined in Scenario 3. To decrease the level, we randomly switch data points between clients until the synergy level reaches 0. The number of data points switched controls the level of synergy, ranging from 0 to 1. We conduct experiments with varying levels of synergy to observe the impact on masked disparity. The results are summarized in Fig. 5 and Table 3.

Table 3. PID of global and local disparity under varying synergy levels  $\lambda$ .

$\lambda$	$I(Z; S)$	Loc.	Glob.	Uniq.	Red.	Mas.
0	0.0035	0.0402	0.0373	0.0338	0.0035	0.0063
0.25	0.0113	0.0486	0.0419	0.0308	0.0111	0.0178
0.5	0.0299	0.0536	0.0335	0.0127	0.0208	0.0410
0.75	0.0846	0.0932	0.0366	0.0023	0.0343	0.0909
1	0.2409	0.1644	0.0149	0.0000	0.0150	0.1644

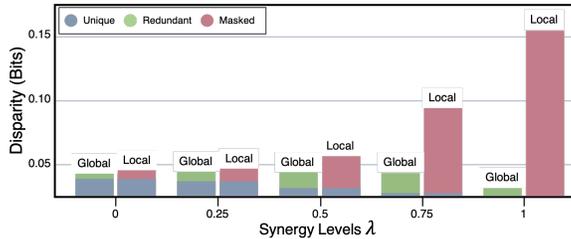


Figure 5. The plot demonstrates the relationship between the synergy level ( $\lambda$ ) and global and local fairness. As the synergy level increases, the masked disparity and local disparity also increases as expected.

**Experiments with Multiple Clients.** Here, we examine scenarios involving multiple clients. Observations are similar to the two-client case we previously studied. We experiment with  $K = 10$  clients. We study the disparities that occur when the data distribution is near i.i.d. distributed among the clients. We manipulate the proportion of sensitive attributes in the first half of the clients by using  $\alpha$ .  $\alpha = \Pr(Z = 0|S = 1, \dots, 5)$ .  $\alpha = 0.33$  would correspond to the case where the sensitive attribute is distributed independently among the clients. We choose values of  $\alpha$  that are close to 0.33. These distributions closely emulate the realistic scenarios that occur in the real world. In this scenario, we observe the presence of all types of disparities. Our results are summarized in Fig. 6 and Table 4.

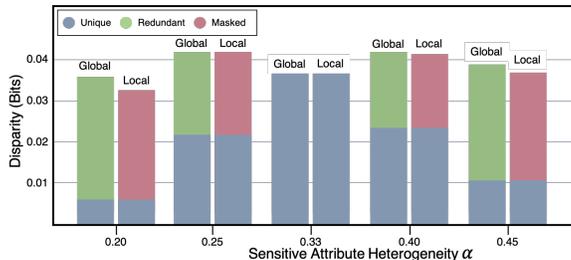


Figure 6. The plot shows the PID of disparities when the data is near i.i.d. among  $K = 10$  clients. All types of disparities can be observed. The value  $\alpha = 0.33$  represents the case where the data is i.i.d. and only unique disparity is observed.

Table 4. PID of global and local disparity for various sensitive attribute distributions across 10 clients.

$\alpha$	Unique	Redundant	Masked	Global	Local	Accuracy
0.25	0.0219	0.0190	0.0178	0.0409	0.0409	84.85%
0.33	0.0376	0.000	0.0000	0.0376	0.0376	85.58%
0.4	0.0268	0.0141	0.0137	0.0410	0.0405	84.85%
0.45	0.0107	0.0289	0.0270	0.039	0.0377	84.85%

**Setup.** In our federated learning model, both the server and clients had two hidden layers, each containing 32 hidden units. The activation function used was ReLU, with a binary cross-entropy loss function and Adam optimizer. The first round began with the server initializing the weights of the model and sharing them with all clients. Each client then trained their local model on their designated local dataset, which was carefully divided to observe various disparities. The training process for each client was done using a batch size of 64 and 2 epochs. After training, each client shared their weight parameters with the server. The server then used the FedAvg algorithm to aggregate the weights of all clients and update the model. The updated weights were then shared back with the clients. This process was repeated for several rounds until the loss converged. The final model is then used for our evaluations.

## B. Related works

There are various perspectives to fairness in FL (Shi et al., 2021). One such definition is *client-fairness* (Li et al., 2019), which aims to achieve equal performance across all client devices. In this work, we are instead interested in group fairness,

i.e., fairness with respect to demographic groups based on gender, race, etc. Methods for achieving group fairness in a centralized machine learning setting (Hardt et al., 2016; Dwork et al., 2012; Kamishima et al., 2011; Pessach & Shmueli, 2022) may not directly apply in a FL setting since each client only has access to their local dataset.

Existing works on group fairness in FL generally aim to develop models that achieve *global fairness*, without much consideration for the *local fairness* at each client (Ezzeldin et al., 2023). For instance, one approach to achieve global fairness in FL poses a constrained optimization problem to find the best model locally, while also ensuring that disparity at a client does not exceed a threshold and then aggregates those models (Chu et al., 2021; Rodríguez-Gálvez et al., 2021; Zhang et al., 2020). Other techniques involve bi-level optimization that aims to find the optimal global model (minimum loss) under the worst-case fairness violation (Papadaki et al., 2022; Hu et al., 2022; Zeng et al., 2021), or re-weighting mechanisms (Abay et al., 2020; Du et al., 2021), both of which often require sharing additional parameters with a server. More recently, (Cui et al., 2021) argue for local fairness, as the model will be deployed at the local client level, and propose constrained multi-objective optimization.

While previous works have made progress in attaining either global fairness (often with additional information sharing) or sometimes local fairness, their interplay has received less attention. In fact, the terms “global fairness” and “local fairness” have often been used somewhat loosely in the literature, without a clear understanding of their relationship, with the notable exceptions of (Ezzeldin et al., 2023; Cui et al., 2021) which claim that global and local fairness are equivalent when the dataset is i.i.d. across clients. Our work addresses this gap by formalizing both global and local fairness and deriving fundamental limits and trade-offs, particularly under data-heterogeneity, by leveraging a body of work in information theory called partial information decomposition (PID) (Bertschinger et al., 2014).

Information-theoretic measures have been widely used to express and handle group fairness in the fairness literature (Kamishima et al., 2012; Calmon et al., 2017; Ghassami et al., 2018; Dutta et al., 2021; 2020b;a; Cho et al., 2020; Baharlouei et al., 2019; Grari et al., 2019; Wang et al., 2021; Galhotra et al., 2022; Alghamdi et al., 2022; Kairouz et al., 2019). Kamishima et al. (2012) uses mutual information as a regularizer with the loss function to minimize the correlation between  $\hat{Y}$  and  $Z$ . Alternate Reyni-measures have also been explored in Baharlouei et al. (2019); Grari et al. (2019).

Dutta et al. (2020a; 2021) introduces PID into algorithmic fairness and explainability for the problem of selectively quantifying disparity that is not due to critical (core) features. We also refer to Dutta & Hamman (2023) for a survey of PID in fairness and explainability. PID is also generating interest in other machine learning applications (Ehrlich et al., 2022; Tax et al., 2017; Liang et al., 2023; Wollstadt et al., 2023; Mohamadi et al., 2023; Pakman et al., 2021). In this work, instead of trying to minimize mutual information as a regularizer, our goal is to quantify the fundamental trade-offs between local and global fairness in federated learning and develop insights on their interplay to better understand what is information-theoretically possible using any optimization technique.

## C. Additional Results and Proofs for Section 2

**Example 4** (Understanding PID). *Let  $Z = (Z_1, Z_2, Z_3)$  with  $Z_1, Z_2, Z_3 \sim \text{i.i.d. Bern}(1/2)$ . Let  $A = (Z_1, Z_2, Z_3 \oplus N)$ ,  $B = (Z_2, N)$ ,  $N \sim \text{Bern}(1/2)$  is independent of  $Z$ . Here,  $I(Z; (A, B)) = 3$  bits.*

The unique information about  $Z$  that is contained only in  $A$  and not in  $B$  is effectively contained in  $Z_1$  and is given by  $\text{Uni}(Z:A|B) = I(Z; Z_1) = 1$  bit. The redundant information about  $Z$  that is contained in both  $A$  and  $B$  is effectively contained in  $Z_2$  and is given by  $\text{Red}(Z:(A, B)) = I(Z; Z_2) = 1$  bit. Lastly, the synergistic information about  $Z$  that is not contained in either  $A$  or  $B$  alone, but is contained in both of them together is effectively contained in the tuple  $(Z_3 \oplus N, N)$ , and is given by  $\text{Syn}(Z:(A, B)) = I(Z; (Z_3 \oplus N, N)) = 1$  bit. This accounts for the 3 bits in  $I(Z; (A, B))$ . Here,  $B$  does not have any unique information about  $Z$  that is not contained in  $A$ , i.e.,  $\text{Uni}(Z:B|A) = 0$ .

**Definition 4** (Unique Information). *Let  $\Delta$  be the set of all joint distributions on  $(Z, A, B)$  and  $\Delta_p$  be the set of joint distributions with the same marginals on  $(Z, A)$  and  $(Z, B)$  as the true distribution, i.e.,  $\Delta_p = \{Q \in \Delta : q(z, a) = \Pr(Z = z, A = a) \text{ and } q(z, b) = \Pr(Z = z, B = b)\}$ . Then,  $\text{Uni}(Z:A|B) = \min_{Q \in \Delta_p} I_Q(Z; A | B)$ , where  $I_Q(Z; A | B)$  is the conditional mutual information when  $(Z, A, B)$  have joint distribution  $Q$ .*

**Lemma 1.** *Let  $Z$  and  $\hat{Y}$  be binary and  $\Pr(Z = 0) = 1 - \Pr(Z = 1) = \alpha$ . The global statistical parity gap  $SP_{\text{global}} = |\Pr(\hat{Y} = 1|Z = 1) - \Pr(\hat{Y} = 1|Z = 0)|$  is bounded by  $\frac{\sqrt{0.5 I(Z; \hat{Y})}}{2\alpha(1-\alpha)}$ .*

*Proof.* Mutual information can be expressed as KL divergence.

$$I(Z; \hat{Y}) = D\left(P(\hat{Y}, Z) \| P(\hat{Y}), P(Z)\right)$$

Using Pinskers Inequality,

$$d_{tv}(P, Q) \leq \sqrt{0.5D(P, Q)}$$

where,  $d_{tv}(P, Q)$  is the total variation between two probability distributions  $P, Q$ .

$$\begin{aligned} d_{tv}\left(\Pr(\hat{Y}, Z), \Pr(\hat{Y}) \Pr(Z)\right) &= \frac{1}{2} \sum_{\hat{y}, z} \left| \Pr_{\hat{Y}, Z}(\hat{y}, z) - \Pr_{\hat{Y}}(\hat{y}) \Pr_Z(z) \right| \\ &= \sum_z \Pr_Z(z) \sum_{\hat{y}} \frac{1}{2} \left| \Pr_{\hat{Y}|Z=z}(\hat{y}) - \Pr_{\hat{Y}}(\hat{y}) \right| \\ &= \frac{1}{2} \Pr(Z=1) \left[ |\Pr(\hat{Y}=1|Z=1) - \Pr(\hat{Y}=1)| + |\Pr(\hat{Y}=0|Z=1) - \Pr(\hat{Y}=0)| \right] \\ &\quad + \frac{1}{2} \Pr(Z=0) \left[ |\Pr(\hat{Y}=1|Z=0) - \Pr(\hat{Y}=1)| + |\Pr(\hat{Y}=0|Z=0) - \Pr(\hat{Y}=0)| \right] \\ &= \frac{1}{2} \alpha(1-\alpha)|SP1| + \frac{1}{2} \alpha(1-\alpha)|SP0| + \frac{1}{2} \alpha(1-\alpha)|SP1| + \frac{1}{2} \alpha(1-\alpha)|SP0| \\ &= \alpha(1-\alpha)|SP1| + \alpha(1-\alpha)|SP0| \end{aligned} \quad (4)$$

where,  $\Pr(Z=0) = 1 - \Pr(Z=1) = \alpha$  and

$$SPi = \Pr(\hat{Y}=i|Z=1) - \Pr(\hat{Y}=i|Z=0) = \Pr(\hat{Y}=i|Z=1) - \Pr(\hat{Y}=i).$$

To complete the proof, we show that  $|SP1| = |SP0|$

$$\begin{aligned} SP1 &= \Pr(\hat{Y}=1|Z=1) - \Pr(\hat{Y}=1) \\ &= \Pr(\hat{Y}=1|Z=1) - (1 - \Pr(\hat{Y}=0)) \\ &= -1 + \Pr(\hat{Y}=1|Z=1) + \Pr(\hat{Y}=0) \\ &= -\Pr(\hat{Y}=0|Z=1) + \Pr(\hat{Y}=0) = -SP0 \end{aligned}$$

Hence,  $|SP1| = |SP0|$  and from (4), we get

$$2\alpha(1-\alpha)|SP1| \leq \sqrt{0.5MI}$$

**Corollary 2.** *The statistical parity at each client  $k$  can be expressed as*

$$|SP_k| \leq \frac{\sqrt{0.5 I(Z, \hat{Y}|S=k)}}{2\alpha_k(1-\alpha_k)}$$

where,  $\alpha_k = \Pr(Z=0|S=k) = 1 - \Pr(Z=1|S=k)$ .

□

## D. Additional Results and Proofs for Section 3

**Theorem 3** (Necessary and Sufficient Condition to Achieve Global Fairness Using Local Fairness). *If local disparity  $I(Z, \hat{Y}|S)$  goes to zero, then the global disparity  $I(Z; \hat{Y})$  also goes to zero, if and only if the redundant disparity  $\text{Red}(Z; \hat{Y}, S) = 0$ .*

*Proof.* From the PID of local and global disparity,

$$\begin{aligned} I(Z; \hat{Y}) &= \text{Uni}(Z; \hat{Y}|S) + \text{Red}(Z; \hat{Y}, S). \\ I(Z; \hat{Y}|S) &= \text{Uni}(Z; \hat{Y}|S) + \text{Syn}(Z; (\hat{Y}, S)). \end{aligned}$$

Therefore if,  $I(Z; \hat{Y}|S) = 0$ , then  $\text{Uni}(Z:\hat{Y}|S) = 0$

Hence,

$$\begin{aligned} I(Z; \hat{Y}) &= \text{Red}(Z:\hat{Y}, S) \\ I(Z; \hat{Y}) = 0 &\iff \text{Red}(Z:\hat{Y}, S) = 0 \end{aligned}$$

□

**Lemma 3.**  $Z \perp\!\!\!\perp S \implies \text{Red}(Z:\hat{Y}, S) = 0$ . When  $Z$  and  $S$  are independent, the redundant disparity is zero.

*Proof.* By leveraging the PID of  $I(Z; S)$  and the non-negative property of the PID terms.

$$\begin{aligned} I(Z; S) &= \text{Uni}(Z:S|\hat{Y}) + \text{Red}(Z:\hat{Y}, S) \\ I(Z; S) &\geq \text{Red}(Z:\hat{Y}, S) \end{aligned}$$

Hence,  $Z \perp\!\!\!\perp S \implies \text{Red}(Z:\hat{Y}, S) = 0$ .

□

**Lemma 4.** If synergistic disparity  $\text{Syn}(Z:(\hat{Y}, S)) = 0$ , the redundant disparity  $\text{Red}(Z:\hat{Y}, S) = 0$  if  $\hat{Y}$  and  $S$  are independent  $\hat{Y} \perp\!\!\!\perp S$  or  $I(\hat{Y}; S) = 0$ , even if  $I(Z; S) > 0$ .

*Proof.* Interaction information expressed in PID terms:

$$\begin{aligned} I(Z; \hat{Y}; S) &= I(Z; \hat{Y}) - I(Z; \hat{Y}|S) \\ &= \text{Red}(Z:\hat{Y}, S) - \text{Syn}(Z; (\hat{Y}, S)) \end{aligned}$$

If synergistic information  $\text{Syn}(Z; (\hat{Y}, S)) = 0$ , we have:

$$\begin{aligned} I(Z; \hat{Y}; S) &= I(Z; \hat{Y}) - I(Z; \hat{Y}|S) \\ &= \text{Red}(Z:\hat{Y}, S) \geq 0 \end{aligned}$$

Since the interaction information is positive and symmetric,

$$I(\hat{Y}; S) \geq I(\hat{Y}; S) - I(\hat{Y}; S|Z) = \text{Red}(Z:\hat{Y}, S)$$

and therefore,  $\hat{Y} \perp\!\!\!\perp S \implies \text{Red}(Z:\hat{Y}, S) = 0$ .

□

**Remark 1.** It is worth noting that the independence between  $\hat{Y}$  and  $S$  can be approximately achieved if the true  $Y$  and  $S$  are independent, as  $\hat{Y}$  is an estimation of  $Y$ . However, it is often the case that  $Y \perp\!\!\!\perp S$  is fixed due to the fixed nature of datasets at each client. The mutual information  $I(Y; S)$  can provide insight into the expected value of  $I(\hat{Y}; S)$ , as the federated model typically aims to also achieve a reasonable level of accuracy. It may even be possible to enforce  $\hat{Y} \perp\!\!\!\perp S$  at the cost of accuracy.

**Corollary 1.** The local disparity will always be less than global disparity if  $Z, \hat{Y}, S$  form a Markov chain  $Z - \hat{Y} - S$ .

*Proof.* By leveraging the PID of  $I(Z; S|\hat{Y})$ ,

$$I(Z; S|\hat{Y}) = \text{Uni}(Z:S|\hat{Y}) + \text{Syn}(Z:(\hat{Y}, S))$$

Hence,  $I(Z; S|\hat{Y}) = 0 \implies \text{Syn}(Z:(\hat{Y}, S)) = 0$

□