
Supplementary File for: FedSR: A Simple and Effective Domain Generalization Method for Federated Learning

Anonymous Author(s)

Affiliation

Address

email

1 A Proofs

2 A.1 Lemma 1

3 *Proof.* Recall that the joint distribution of x, y, z for client i is:

$$p_i(x, y, z) = p_i(x, y)p(z|x) = p_i(y)p_i(x|y)p(z|x) \quad (1)$$

4 We have:

$$I_i(x, z|y) = \mathbb{E}_{p_i(x, y, z)} \left[\log \frac{p_i(x, z|y)}{p_i(x|y)p_i(z|y)} \right] \quad (2)$$

$$= \mathbb{E}_{p_i(x, y, z)} \left[\log \frac{p_i(x|y)p(z|x)}{p_i(x|y)p_i(z|y)} \right] \quad (3)$$

$$= \mathbb{E}_{p_i(x, y, z)} \left[\log \frac{p(z|x)}{p_i(z|y)} \right] \quad (4)$$

$$= \mathbb{E}_{p_i(x, y, z)} \left[\log \frac{p(z|x)}{p_i(z|y)} \right] \quad (5)$$

$$= \mathbb{E}_{p_i(x, y, z)} [\log p(z|x) - \log p_i(z|y)] \quad (6)$$

5 Notice that:

$$\mathbb{E}_{p_i(y)} [\text{KL}[p_i(z|y)|r(z|y)]] \geq 0 \quad (7)$$

$$\Rightarrow \mathbb{E}_{p_i(y)} [\mathbb{E}_{p_i(z|y)} [\log p_i(z|y) - \log r(z|y)]] \geq 0 \quad (8)$$

$$\Rightarrow \mathbb{E}_{p_i(y, z)} [\log p_i(z|y) - \log r(z|y)] \geq 0 \quad (9)$$

$$\Rightarrow \mathbb{E}_{p_i(x, y, z)} [\log p_i(z|y) - \log r(z|y)] \geq 0 \quad (10)$$

$$\Rightarrow \mathbb{E}_{p_i(x, y, z)} [\log p_i(z|y)] \geq \mathbb{E}_{p_i(x, y, z)} [\log r(z|y)] \quad (11)$$

$$(12)$$

6 Therefore:

$$I_i(x, z|y) = \mathbb{E}_{p_i(x, y, z)} [\log p(z|x) - \log p_i(z|y)] \quad (13)$$

$$\leq \mathbb{E}_{p_i(x, y, z)} [\log p(z|x) - \log r(z|y)] \quad (14)$$

$$= \mathbb{E}_{p_i(x, y)} [\mathbb{E}_{p(z|x)} [\log p(z|x) - \log r(z|y)]] \quad (15)$$

$$= \mathbb{E}_{p_i(x, y)} [\text{KL}[p(z|x)|r(z|y)]] \quad (16)$$

7

□

8 A.2 Lemma 2

9 *Proof.* We need to prove:

$$\mathbb{E}_{p_i(x,y)} [\text{KL}[p(z|x)|r(z|y)]] \geq \mathbb{E}_{p_i(y)} [\text{KL}[p_i(z|y)|r(z|y)]] \quad (17)$$

$$\Leftrightarrow \mathbb{E}_{p_i(x,y,z)} [\log p(z|x) - \log r(z|y)] \geq \mathbb{E}_{p_i(y,z)} [\log p_i(z|y) - \log r(z|y)] \quad (18)$$

$$\Leftrightarrow \mathbb{E}_{p_i(x,y,z)} [\log p(z|x)] \geq \mathbb{E}_{p_i(y,z)} [\log p_i(z|y)] \quad (19)$$

$$(\text{since } \mathbb{E}_{p_i(x,y,z)} [\log r(z|y)] = \mathbb{E}_{p_i(y,z)} [\log r(z|y)]) \quad (20)$$

$$\Leftrightarrow \mathbb{E}_{p_i(y)} [\mathbb{E}_{p_i(x,z|y)} [\log p(z|x)]] \geq \mathbb{E}_{p_i(y)} [\mathbb{E}_{p_i(z|y)} [\log p_i(z|y)]] \quad (21)$$

$$(22)$$

10 Therefore, we only need to prove that $\mathbb{E}_{p_i(x,z|y)} [\log p(z|x)] \geq \mathbb{E}_{p_i(z|y)} [\log p_i(z|y)] \forall y$.

11 This is equivalent to:

$$\int_z \int_x p_i(x, z|y) \log p(z|x) dx dz \geq \int_z p_i(z|y) \log p_i(z|y) dz \quad (23)$$

$$\Leftrightarrow \int_z \int_x p_i(x|y) [p(z|x) \log p(z|x)] dx dz \geq \int_z p_i(z|y) \log p_i(z|y) dz \quad (24)$$

$$\Leftrightarrow \int_z \mathbb{E}_{p_i(x|y)} [p(z|x) \log p(z|x)] dz \geq \int_z p_i(z|y) \log p_i(z|y) dz \quad (25)$$

12 Notice that the function $f(a) = a \log a$, $a > 0$ is a convex function since $f''(a) = \frac{1}{a} > 0 \forall a > 0$.

13 Hence, due to Jensen's Inequality we have:

$$\mathbb{E}_{p_i(x|y)} [p(z|x) \log p(z|x)] \geq b \log b \quad (26)$$

14 where $b = \mathbb{E}_{p_i(x|y)} [p(z|x)] = p_i(z|y)$.

15 Therefore, $\mathbb{E}_{p_i(x|y)} [p(z|x) \log p(z|x)] \geq p_i(z|y) \log p_i(z|y)$. Then, inequality 25 holds, and we
16 conclude our proof. \square

17 B Details on the ℓ_i^{CMI} objective

18 As mentioned in the main paper, we use a Gaussian distribution for both $p(z|x)$ and $r(z|y)$, so that
19 the KL term can be computed analytically.

20 Let w_1 be the parameter of the (probabilistic) representation network. Given an input x_0 (of a
21 datapoint (x_0, y_0)), the network output the mean and standard deviation of the Gaussian distribution
22 $p(z|x_0)$ (both from the final layer of the network), denoted $\mu_{w_1}(x_0)$ and $\sigma_{w_1}(x_0)$. This means that
23 $p(z|x_0) = \mathcal{N}(z; \mu_{w_1}(x_0), \sigma_{w_1}^2(x_0))$.

24 Also as pointed out in the paper, for a classification problem (where there is a finite number of the
25 labels), we set $r(z|y) = \mathcal{N}(z; \mu_y, \sigma_y^2)$, where μ_y, σ_y ($y \in \overline{1..C}$) are the variational parameters to be
26 optimized. Therefore, for the datapoint x_0, y_0 , we have $r(z|y_0) = \mathcal{N}(z; \mu_{y_0}, \sigma_{y_0}^2)$.

27 Note that when z is a K -dimensional vector, $\mu_{w_1}(x_0), \sigma_{w_1}(x_0), \mu_{y_0}$ and σ_{y_0} are all K -dimensional.

28 The KL term in ℓ^{CMI} can be computed analytically as:

$$\log \sigma_{y_0} - \log \sigma_{w_1}(y_0) + \frac{\sigma_{w_1}^2(x_0) + (\mu_{w_1}(x_0) - \mu_{y_0})^2}{2\sigma_{y_0}^2} - \frac{1}{2} \quad (27)$$

29 If $K > 1$ (i.e., z is high dimensional), the calculations in Eq. 27 are element-wise, and the result are
30 summed across the dimension K at the end.

31 As this computation is deterministic with respect to w_1 and μ_y, σ_y ($y \in \overline{1..C}$), the gradient w.r.t.
32 these parameters can be computed straight-forwardly.

33 C Experiments

34 C.1 Hyper-parameters:

35 As mentioned in the main text, we use random search to tune the hyper-parameters of our method,
 36 namely α^{CMI} and α^{L2R} . Specifically, the tuned value for those parameters are:

37 RotatedMNIST

- 38 • FedL2R: $\alpha^{L2R} = 0.1$
- 39 • FedCMI: $\alpha^{CMI} = 0.3$
- 40 • FedSR: $\alpha^{L2R} = 0.1, \alpha^{CMI} = 0.3$

41 PACS

- 42 • FedL2R: $\alpha^{L2R} = 0.01$
- 43 • FedCMI: $\alpha^{CMI} = 0.01$
- 44 • FedSR: $\alpha^{L2R} = 0.01, \alpha^{CMI} = 0.001$

45 OfficeHome

- 46 • FedL2R: $\alpha^{L2R} = 0.05$
- 47 • FedCMI: $\alpha^{CMI} = 0.0005$
- 48 • FedSR: $\alpha^{L2R} = 0.05, \alpha^{CMI} = 0.0005$

49 DomainNet

- 50 • FedL2R: $\alpha^{L2R} = 0.01$
- 51 • FedCMI: $\alpha^{CMI} = 0.005$
- 52 • FedSR: $\alpha^{L2R} = 0.01, \alpha^{CMI} = 0.0005$

53 C.2 Visualization of the Representation space

54 We also visualize the representation space of our method and observe that it aligns the representation
 55 much better than the conventional FL baseline FedAVG (which does not actively attempt to align
 56 the representation distributions). Specifically, Figure 1 show the t-SNE [1] visualization of the
 57 representation of FedSR and FedAVG in the RotatedMNIST experiment, with \mathcal{M}_0 (the target
 58 domain) and \mathcal{M}_{15} (one of the source domains). We can clearly see that both the marginal and the
 59 conditional distributions of the representation are aligned better between the two domains for FedSR.

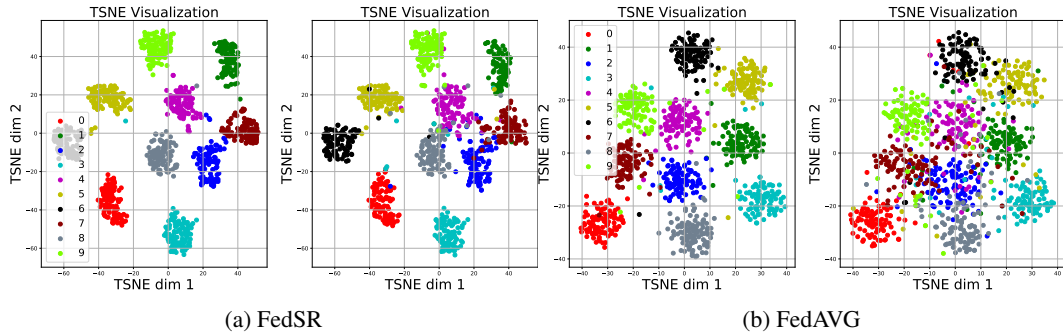


Figure 1: Visualization using t-SNE of the representation space of our method FedSR and the baselines FedAVG. For each method, the left subfigure corresponds to one source domain \mathcal{M}_{15} and the right one corresponds to the target domain \mathcal{M}_0 . Each color represents a digit class.

60 **References**

- 61 [1] L. Van der Maaten and G. Hinton. Visualizing data using t-sne. *Journal of machine learning*
62 *research*, 9(11), 2008.